

# Overcoming numerical shockwave anomalies using energy balanced numerical schemes. Application to the Shallow Water Equations with discontinuous topography.

A. Navas-Montilla , J. Murillo <sup>1</sup>

*anavas@unizar.es, Javier.Murillo@unizar.es, Fluid Mechanics-LIFTEC,  
CSIC-Universidad de Zaragoza. Zaragoza, Spain*

---

## Abstract

When designing a numerical scheme for the resolution of conservation laws, the selection of a particular source term discretization (STD) may seem irrelevant whenever it ensures convergence with mesh refinement, but it has a decisive impact on the solution. In the framework of the Shallow Water Equations (SWE), well-balanced STD based on quiescent equilibrium are unable to converge to physically based solutions, which can be constructed considering energy arguments. Energy based discretizations can be designed assuming dissipation or conservation, but in any case, the STD procedure required should not be merely based on ad hoc approximations. The STD proposed in this work is derived from the Generalized Hugoniot Locus obtained from the Generalized Rankine Hugoniot conditions and the Integral Curve across the contact wave associated to the bed step. In any case, the STD must allow energy-dissipative solutions: steady and unsteady hydraulic jumps, for which some numerical anomalies have been documented in the literature. These anomalies are the incorrect positioning of steady jumps and the presence of a spurious spike of discharge inside the cell containing the jump. The former issue can be addressed by proposing a modification of the energy-conservative STD that ensures a correct dissipation rate across the hydraulic jump, whereas the latter is of greater complexity and cannot be fixed by simply choosing a suitable STD, as there are more variables involved. The problem concerning the spike of discharge is a well-known problem in the scientific community, also known as slowly-moving shock anomaly, it is produced by a non-linearity of the Hugoniot locus connecting the states at

both sides of the jump. However, it seems that this issue is more a feature than a problem when considering steady solutions of the SWE containing hydraulic jumps. The presence of the spurious spike in the discharge has been taken for granted and has become a feature of the solution. Even though it does not disturb the rest of the solution in steady cases, when considering transient cases it produces a very undesirable shedding of spurious oscillations downstream that should be circumvented. Based on spike-reducing techniques (originally designed for homogeneous Euler equations) that propose the construction of interpolated fluxes in the untrustworthy regions, we design a novel Roe-type scheme for the SWE with discontinuous topography that reduces the presence of the aforementioned spurious spike. The resulting spike-reducing method in combination with the proposed STD ensures an accurate positioning of steady jumps, provides convergence with mesh refinement, which was not possible for previous methods that cannot avoid the spike.

*Keywords:* Roe solver, Energy balanced, Shallow water, Source terms, Hydraulic jump, Postshock oscillations

---

## 1. Introduction

There is a wide variety of physical problems modelled by non-homogeneous hyperbolic systems of conservation laws that are dominated by source terms. For such problems, the treatment of the source terms when designing a numerical scheme is of utmost importance in order to provide realistic and physically feasible solutions. Depending on the nature of the source term, different numerical techniques may be required. In this work, we focus on a certain type of source term, called geometric source term, present in many physical one-dimensional (1D) problems. This kind of source makes the conserved quantities account for the variation in space of a geometric variable, which is provided in the problem. Examples of mathematical models including geometric source terms are, for instance, the SWE with discontinuous topography, which is the object of study in the present work, the 1D Euler equations in a duct of variable cross section [1] and 1D flow in collapsible vessels [2].

Most popular methods for the resolution of homogeneous hyperbolic problems are within the framework of finite volume Godunov's numerical schemes [3], which aim to provide a numerical solution to the problem by means of

19 a prior discretization of the domain into volume cells and integration of the  
20 information and governing equations inside these cells. After integration,  
21 simple algebraic evolution equations for the conserved variables, that de-  
22 pend upon the same variables at a previous time step and the fluxes at cell  
23 interfaces, arise. The keystone in Godunov's schemes is the computation of  
24 the numerical fluxes at cell interfaces, which is carried out by means of the  
25 resolution of the so-called Riemann Problems (RPs). RPs are initial value  
26 problems defined at cell interfaces, whose initial data is piecewise constant  
27 data given by the cell-averaged variables at each side of the discontinuity.  
28 They may be regarded as first order approach to the more general Cauchy  
29 problem [4].

30 When dealing with geometric source terms, it is necessary to account  
31 for the jump of the geometric quantity across cell interfaces when defining  
32 numerical fluxes at cell interfaces. To this end, augmented solvers were intro-  
33 duced [5, 6, 7]. When using augmented solvers, the source term is accounted  
34 for in the solution of the RP as an extra stationary wave at the interface. Due  
35 to the presence of the new wave, two solutions appear now at each side of the  
36 initial discontinuity instead of having a single homogeneous solution. Aug-  
37 mented versions of the traditional Roe [8] (ARoe) and HLLC [9, 10] solvers  
38 were presented by Murillo in [11] and [12] respectively. An extense review of  
39 the ARoe method can be found in [13].

40 If examining the system of equations in the so-called non-conservative  
41 form, the contribution of the source term is modelled as an additional sta-  
42 tionary wave at the interface, which allows to include the effect of the source  
43 term in the eigenstructure of the system. This way, it can be noticed that the  
44 presence of a jump in the geometric variable gives rise to a contact wave and  
45 furthermore, that Riemann invariants are not necessarily conserved across  
46 such a wave, as pointed out by Rosatti et al. [14]. This issue will be recalled  
47 when designing the numerical scheme.

48 In the early stages of the design of numerical schemes for hyperbolic prob-  
49 lems with source terms, the main effort was put on how to modify the original  
50 schemes, initially designed for homogeneous equations, so that they maintain  
51 the discrete equilibrium between fluxes and source term under steady state.  
52 When considering realistic applications, such goal was translated into the  
53 preservation of physical steady situations of quiescent equilibrium. For in-  
54 stance, in the framework of the SWE, the preservation of the steadiness of  
55 the solution for still water at rest. Numerical schemes satisfying this property  
56 were called well-balanced schemes [15, 16, 17, 18, 19].

57 When considering steady states with moving water over a irregular bed  
58 profile, the preservation of the C-property (exact conservation property) [16]  
59 is also of utmost importance in order to provide an exact equilibrium between  
60 fluxes and source terms. Numerical methods preserving the C-property are  
61 able to ensure a uniform discharge under steady conditions and can be con-  
62 structed using flux-type definitions of the source terms [20, 6].

63 We can still take the well-balanced and C-property a step further by con-  
64 sidering the conservation of the discrete specific mechanical energy in the  
65 scheme, enhancing in this way the performance of the numerical method.  
66 When friction is not considered in the SWE, mechanical energy is con-  
67 served under steady conditions in absence of hydraulic jumps. Such idea  
68 of energy conservation can be integrated in the numerical scheme, allowing  
69 the extension of well-balanced methods to exactly well-balanced methods  
70 [21, 22, 23, 24, 25, 26], hereafter referred to as E-schemes. Numerical meth-  
71 ods defined as E-schemes will always satisfy the energy conservation property  
72 in the discrete level, hereafter referred to as E-property. Arbitrary order aug-  
73 mented Roe and HLL schemes preserving the E-property, called AR-ADER  
74 and HLLS-ADER E-schemes respectively, were presented by the authors of  
75 this work in [27, 28] and applied to the SWE. As a result of preserving the  
76 E-property, the aforementioned schemes were able to provide the exact solu-  
77 tion in transient cases with independence of the grid and also to converge to  
78 the exact solution in transient problems at a high rate as the grid is refined.

79 For transient problems in the framework of the SWE, different approaches  
80 can be found in the literature regarding the treatment of the source term  
81 contact discontinuity. Two main tendencies are observed in the literature:  
82 one is based on energy and mass conservation and the other one based on  
83 mass and momentum conservation. For instance, some authors [29, 30] claim  
84 that energy must always be conserved since the bed step discontinuity is a  
85 contact wave and Riemann invariants, namely mass and energy for the bed  
86 step discontinuity, are conserved across contact waves. Alcrudo et al. [31] also  
87 state that the use of the mass-energy approach is necessary, specially when  
88 the slope of the bed profile becomes infinite (e.g. in the bed step), however,  
89 they allow for the possibility of some dissipation across the bed step, due to  
90 recirculation. On the other hand, Bernetti et al. [32] hold that the relation  
91 among variables across a bed discontinuity must be calculated by means of  
92 the Generalized Rankine-Hugoniot (GRH) conditions for the full system of  
93 equations. As an effort to unify all the previous approaches, Rosatti et al. [14]  
94 proposed a novel technique, based on the GRH conditions and using energy

95 as a constraint to rule out solutions that are not physically admissible. They  
96 show that in nonconservative systems, such as the SWE, unlike in standard  
97 conservative systems, Riemann invariants are generally not constant across  
98 a contact discontinuity whose relevant eigenvalue is independent from the  
99 problem variables, and use this statement to design a numerical scheme that  
100 allows for the presence of dissipation due to recirculation at the bed step.

101 In the present work, the authors are faithful to the original SW system  
102 and do not include any dissipation mechanism (e.g. recirculation at bed  
103 step), as the original equations do not consider friction terms. Dissipation  
104 will only take place in certain conditions, such as a sudden change of flow  
105 regime (hydraulic jump), according to the physical behavior described by the  
106 equations. A theoretical study on the relations among states across the bed  
107 step contact wave is included in the text, leading to the particular conditions  
108 that ensure conservation of energy across the step: the Generalized Hugoniot  
109 Locus (GHL) derived from the GRH must coincide with the Integral Curve  
110 (IC). In other words, not only the GRH conditions must be fulfilled but also  
111 Riemann invariants should be conserved, as the specific mechanical energy is  
112 one of the relevant invariants for the characteristic field of the contact wave.

113 The AR-ADER and HLLS-ADER methods in [28], proposed by the au-  
114 thors of this work, are based on a particular energy conservative STD which  
115 is computed as a linear combination of a differential and integral approxi-  
116 mation of the integral of the source term at cell interfaces. The method was  
117 presented in [25] for the first time and allowed to enhance the capabilities of  
118 augmented solvers in the framework of the SWE. Very high order methods  
119 are truly desirable as they have the ability of reducing dramatically numeri-  
120 cal diffusion, allowing to provide predictions that would not be affordable by  
121 first order numerical schemes [33]. This can be done at the cost of replacing  
122 time derivatives by spatial derivatives. As a result, the strengths and also  
123 the weaknesses of the approximate solver used are enhanced.

124 E-schemes in [28] have desirable properties: they provide the exact solu-  
125 tion for steady cases and are convergent to the exact solution with arbitrary  
126 order for transient cases including non-resonant and resonant cases. But  
127 there is still room for improvement. A recent study on the convergence of  
128 several schemes, including first order ARoe E-scheme, to steady shocks (hy-  
129 draulic jumps) [34] proved that this scheme leads to a displacement of the  
130 hydraulic jump. When moving to very high order, integration of the source  
131 term must be done using a quadrature rule that matches with the order of  
132 convergence of the numerical scheme. This could be seen as an opportu-

133 nity to improve numerical results regarding the positioning of the hydraulic  
134 jump, but contrary to intuition, the same issue observed in the first order  
135 scheme is repeated when using the high order methods in [28]. This issue is  
136 deeply studied and addressed here, proposing a STD that makes the scheme  
137 unequivocally identify the position of the hydraulic jump and dissipate the  
138 exact amount of energy across it. This technique will be referred to as selec-  
139 tive energy balanced formulation (SEBF) of the integral of the source term  
140 and is applied to the ARoe and HLLS solvers, and their high order versions.

141 High order also preserves the effect of undesirable numerical shockwave  
142 anomalies. The utilization of high order numerical schemes in presence  
143 of spurious oscillations prevents numerical diffusion from dissipating those  
144 oscillations as fast as they would be dissipated if a first order scheme was  
145 used. It has been widely reported in the literature that significant numerical  
146 anomalies arise in presence of shock waves. An example of such problems are  
147 the Carbuncle [35, 36], the slowly-moving shock [37, 40] and the wall-heating  
148 phenomenon [41], all of them leading to spurious numerical solutions. An-  
149 other major point addressed in this work is the study of such anomalies in  
150 the framework of SWE with and without bed variations and the extension of  
151 a spike-reducing scheme for non-homogeneous systems that avoids the pres-  
152 ence of spurious oscillations due to numerical shocks. Shockwaves are typical  
153 solutions for nonlinear hyperbolic systems of conservation laws and their nu-  
154 merical treatment is of utmost importance to provide accurate solutions. As  
155 mentioned by Zaide and Roe [42], physical shockwaves have a finite width  
156 which is determined by the physical dissipation processes, however, when  
157 considering numerical shockwaves, a numerical width, usually much greater  
158 than the physical width, is enforced. This leads to the appearance of inter-  
159 mediate states which cannot be given a direct physical interpretation. Such  
160 states cannot be removed even when refining the grid, therefore we find in  
161 the literature that a special emphasis is put on this issue when designing a  
162 numerical scheme. Up to the present time, most studies have been carried  
163 out in the framework of Euler equations. In this work we will focus on the  
164 SWE.

165 Some of the problems related to numerical shockwave anomalies were first  
166 identified by Cameron and Emery [43, 44], who proposed some improvements  
167 based on the addition of artificial viscosity and modification of the grid.  
168 Here, we focus on the slowly-moving shock problem, which is associated to  
169 hydraulic jumps in the SWE. The slowly-moving shock problem was first  
170 investigated by Roberts in [37], who defined it as numerical noise generated

171 in the discrete shock transition layer which is transported downstream. Such  
172 noise will be hereafter referred to as post-shock oscillations. In [37], the  
173 schemes of Godunov, Roe, and Osher were examined and the source of this  
174 error as also provided by using the Hugoniot locus. It was also observed  
175 that the slowly-moving shock problem only appears for systems of equations  
176 and not for scalar equations, where such schemes perform correctly. It is  
177 worth pointing out that even for non-linear systems, the slowly-moving shock  
178 problem does not appear if the Hugoniot curves are linear [38], as happens in  
179 the system in [39]. Later on, Arora and Roe [40] carried out a thorough study  
180 on this problem and evidenced that it can be ruinous when, for instance,  
181 making calculations of shock-sound interaction.

182 The spike-reducing techniques presented in this work are of first order of  
183 accuracy and one could think that by increasing the order of the scheme the  
184 slowly-moving shock problem could be circumvented. However, as mentioned  
185 by other authors [38, 45, 46], the slowly-moving shock problem will only be  
186 accentuated when increasing the accuracy of the scheme. Such an increase  
187 of accuracy will be translated into a longer preservation of post-shock os-  
188 cillations as they provide a better resolution of the spurious physics. When  
189 using a high order scheme, the order is reduced to first order in the vicinity of  
190 the shock and the numerical solution within this region will behave accord-  
191 ing to what is expected from a first order scheme [47, 48]. Away from the  
192 shock, the order of accuracy is higher and therefore the spurious oscillations  
193 will be better resolved, preventing them from vanishing as one would desire.  
194 It must be borne in mind that even when using high order interpolations  
195 with limiting techniques, such as Total Variation Diminishing (TVD) inter-  
196 polations and Essentially Non-oscillatory (ENO) schemes, the slowly-moving  
197 shock problem is accentuated [46].

198 The slowly-moving shock problem has been deeply studied for homoge-  
199 neous systems of equations (e.g. the Euler equations) but scarcely studied  
200 for systems dominated by source terms. In [46], numerical results for the  
201 computation of a 1D compressible flow through a divergent nozzle by means  
202 of different first and high order schemes were presented, showing the inabil-  
203 ity of all schemes to converge to the exact solution in presence of shocks.  
204 The authors outline that this is due to the appearance of a spike in the  
205 momentum and the shedding of spurious oscillations downstream. This is  
206 the slowly-moving shock problem in the limit when shock speed is nil. The  
207 SWE are analogous to the 1D compressible flow with varying area, hence the  
208 slowly-moving shock problem is also likely to appear.

209 Here we focus on the slowly-moving shock problem in the SWE. To this  
210 end, we identify the conditions for the aforementioned problem to appear by  
211 studying the Hugoniot locus of the SWE and by seeking slowly-moving shock-  
212 type waves. We notice that they are only produced when dealing with a kind  
213 of transcritical shocks called hydraulic jumps, characterized by a change of  
214 sign of the relevant eigenvalue across them. A complete description of such  
215 kind of waves is provided and a thorough study on the shock structure,  
216 comparing exact and Godunov type solutions, is carried out in phase space.  
217 The slowly-moving shock problem in the SWE is a well-known problem in  
218 the scientific community, characterized by a spike in the discharge at the cell  
219 where the hydraulic jump is contained. In fact, it seems that this problem  
220 is more a feature than a problem when considering steady solutions of the  
221 SWE containing hydraulic jumps. The presence of the spurious spike in  
222 the discharge has been taken for granted as it does not perturb the rest  
223 of the solution. However, when considering transient cases, it produces a  
224 very undesirable shedding of spurious oscillations downstream that should  
225 be avoided.

226 When designing numerical schemes for the computation of slowly-moving  
227 shocks, the addition of extra artificial viscosity seems to be the most preferred  
228 technique in the scientific community [43, 44, 37, 40, 45, 49, 50]. If we  
229 want to avoid extra diffusion, another suitable possibility is the use of inter-  
230 polation of fluxes, which avoids using the evaluation of the physical fluxes in  
231 the untrustworthy intermediate cells corresponding to the shock discontinuity.  
232 This idea of flux interpolation was first presented by Zaide and Roe [42],  
233 who proposed to find the fluxes in the intermediate cells by extrapolation  
234 from trustworthy neighbors. The authors claim that, by enforcing a linear  
235 shock structure and unambiguous sub-cell shock position, numerical shock-  
236 wave anomalies are dramatically reduced. It could be said that their method  
237 is also based on the addition of artificial viscosity, as their flux functions can  
238 be regarded as the traditional Roe flux plus a viscosity term. However, the  
239 flux interpolation functions use dissipation to control shock structure rather  
240 than to approach the true viscous solution and therefore they do not expand  
241 the shock profile [38].

242 In this work, we use the approach in [42] to propose a novel spike-reducing  
243 flux function for the SWE with varying bed. Prior to the presentation of  
244 the proposed technique, the flux functions in [42] are applied to the SWE  
245 with flat bed, showing their spike-reducing nature. The proposed technique  
246 is assessed in a variety of situations, including steady and transient cases,



247 with continuous and discontinuous bed profiles, proving the expected spike-  
 248 reducing behavior. The analogous SWE problem of the 1D nozzle problem  
 249 in [46], which is the steady flow over a hump, is reproduced in this work,  
 250 showing that the proposed scheme leads to a convergent solution, even when  
 251 measured with  $L_\infty$  error norm.

252 The outline of the paper is next presented. In section 2, an introduction to  
 253 nonlinear systems of conservation laws with source terms is provided and the  
 254 definition of geometric source terms and derivation of the GRH conditions for  
 255 such systems are recalled. In this section, the description of non-conservative  
 256 systems and the treatment of contact waves in this kind of systems is also  
 257 recalled following [14]. In Section 3, we briefly describe Godunov type finite  
 258 volume schemes. Section 4 is devoted to the description of the SWE, both in  
 259 conservative and non-conservative form, including a thorough study on the  
 260 bed step contact wave. In this section, the numerical treatment of the source  
 261 term in the SWE is also described and the novel SEBF discretization method  
 262 is presented. At the end of this section, numerical results for the computation  
 263 of steady flows are displayed. Section 5 is entirely devoted to the study  
 264 of numerical shockwave anomalies in the SWE. A thorough description of  
 265 the slowly-moving shock problem arising from the hydraulic jump, using  
 266 the phase-space representation, is presented. In Section 6, numerical fixes  
 267 addressing the aforementioned problem are studied. First, numerical results  
 268 for the computation of several homogeneous test cases using the flux functions  
 269 A and B in [42] are shown. Then, the novel spike-reducing technique for the  
 270 SWE with source term is presented and a set of tests are carried out to  
 271 evidence the capabilities of the proposed method. Finally, in Section 7 we  
 272 present a summary of the work and the concluding remarks.

## 273 **2. Nonlinear systems of equations with source term**

274 The basic ideas underlying this work can be illustrated by examining  
 275 hyperbolic nonlinear systems of equations with source terms in 1D, that can  
 276 be expressed in integral form as

$$\frac{\partial}{\partial t} \int_{x_1}^{x_2} \mathbf{U} dx + \mathbf{F}|_{x_2} - \mathbf{F}|_{x_1} - \int_{x_1}^{x_2} \mathbf{S} dx = 0, \quad (1)$$

277 where  $x_1, x_2$  are the limits of a generic control volume and with  $N_\lambda$  equations.  
 278 Such systems arise naturally from the conservation laws for certain physical  
 279 quantities in nature. The differential formulation is obtained when assuming

280 a smooth variation of the variables and an infinitesimal width of the control  
 281 volume, yielding

$$\frac{\partial \mathbf{U}}{\partial t} + \frac{\partial \mathbf{F}}{\partial x} = \mathbf{S}, \quad (2)$$

282 where  $\mathbf{U} = \mathbf{U}(x, t) \in \mathcal{C} \subset \mathbb{R}^{N_\lambda}$  is the vector of conserved quantities that  
 283 takes values on  $\mathcal{C}$ , the set of admissible states of  $\mathbf{U}$ ,  $\mathbf{F} = \mathbf{F}(\mathbf{U})$  is the flux  
 284 function that represents a nonlinear mapping of the conserved quantities from  
 285  $\mathcal{C}$  to  $\mathbb{R}^{N_\lambda}$  and  $\mathbf{S}$  is the source term, that will be considered a function of the  
 286 conserved quantities and spatial coordinate as  $\mathbf{S} = \mathbf{S}(\mathbf{U}, x)$ . In this work,  
 287 we put a special emphasis on the so-called *geometric source terms*, that are  
 288 expressed as

$$\mathbf{S}(\mathbf{U}, x) = \mathbf{S}_s(\mathbf{U}) \frac{d}{dx} \mathbf{S}_g(x), \quad (3)$$

289 with  $\mathbf{S}_s(\mathbf{U})$  a function of the conserved quantities and  $\mathbf{S}_g(x)$  the geometric  
 290 function that depends upon the position  $x$  and can be discontinuous [28].

291 From (2), the Jacobian matrix of the convective part is defined as

$$\mathbf{J} = \frac{d\mathbf{F}(\mathbf{U})}{d\mathbf{U}}. \quad (4)$$

292 Assuming that the convective part in (2) is strictly hyperbolic, with  $N_\lambda$   
 293 real eigenvalues  $\lambda^1, \dots, \lambda^{N_\lambda}$  and eigenvectors  $\mathbf{e}^1, \dots, \mathbf{e}^{N_\lambda}$ , it is possible to de-  
 294 fine the matrices  $\mathbf{P} = (\mathbf{e}^1, \dots, \mathbf{e}^{N_\lambda})$  and  $\mathbf{P}^{-1}$  with the property that they  
 295 diagonalize the Jacobian as

$$\mathbf{J} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^{-1}. \quad (5)$$

### 296 2.1. Conservative vs non-conservative form

297 For the sake of simplicity, dependency of variables upon the conserved  
 298 quantities is hereafter omitted. A generic homogeneous conservative system  
 299 is written as

$$\frac{\partial \mathbf{U}}{\partial t} + \frac{\partial \mathbf{F}}{\partial x} = 0, \quad (6)$$

300 where  $\mathbf{U}$  is the vector of conserved quantities and  $\mathbf{F}$  the vector of conservative  
 301 fluxes. It can be expressed in its quasilinear form as

$$\frac{\partial \mathbf{U}}{\partial t} + \mathbf{J} \frac{\partial \mathbf{U}}{\partial x} = 0, \quad (7)$$

302 where the Jacobian matrix  $\mathbf{J} = d\mathbf{F}/d\mathbf{U}$  can be diagonalized with  $N_\lambda$  eigenval-  
 303 ues by means of  $N_\lambda$  linearly independent eigenvectors. The following relation  
 304 is worth being shown

$$\mathbf{J} \cdot \mathbf{e}^m - \lambda^m \mathbf{e}^m = 0, \quad (8)$$

305 where  $\lambda^m$  and  $\mathbf{e}^m$  are the eigenvalues and right eigenvectors of matrix  $\mathbf{J}$ .

306 Non-homogeneous hyperbolic conservation laws (2) cannot be expressed  
 307 in the strict conservative form of (6) due to the presence of the source term.  
 308 When having geometric source terms of the type of (3), they can be expressed  
 309 in non-conservative form as

$$\frac{\partial \hat{\mathbf{U}}}{\partial t} + \frac{\partial \hat{\mathbf{F}}(\hat{\mathbf{U}})}{\partial x} + \mathbf{H} \frac{\partial \hat{\mathbf{U}}}{\partial x} = 0, \quad (9)$$

310 where  $\hat{\mathbf{U}} \in \mathcal{C} \subset \mathbb{R}^{N_\lambda + N_S}$  is the new vector of variables composed of the  $N_\lambda$   
 311 conserved variables in (2) plus additional  $N_S$  variables related to the source  
 312 term,  $\hat{\mathbf{F}}(\hat{\mathbf{U}}) : \mathcal{C} \rightarrow \mathbb{R}^{N_\lambda + N_S}$  is the vector of conservative fluxes and  $\mathbf{H}$  the  
 313 matrix of non-conservative fluxes.

314 In this work, we will focus on physical problems (e.g. the shallow water  
 315 model with bed topography) with a geometric source term like (3) that only  
 316 involves a single geometric quantity,  $s_g(x)$ , as follows

$$\mathbf{S}_g(x) = (0, \dots, s_g(x), \dots, 0)^T. \quad (10)$$

317 In this case, the new vector of variables will be constructed as  $\hat{\mathbf{U}} =$   
 318  $(\mathbf{U}, s_g)^T$ , hence  $N_S = 1$ , with  $\lambda^s = 0$ , the speed of the wave associated to the  
 319 source equal to zero as the geometric quantity does not evolve in time. This  
 320 is depicted in Figure 1, for an arbitrary system with  $N_\lambda = 3$  and a single  
 321 geometric variable, that is  $N_S = 1$ .

322 Also notice that the evolution equation corresponding to the geometric  
 323 quantity,  $s_g$ , reads

$$\frac{\partial s_g}{\partial t} = 0, \quad (11)$$

324 which stands for the conservation of this quantity in time, as it only depends  
 325 upon the spatial position  $x$ .

326 The non-conservative system in (9) can be more compactly expressed as

$$\frac{\partial \hat{\mathbf{U}}}{\partial t} + \mathbf{A} \frac{\partial \hat{\mathbf{U}}}{\partial x} = 0, \quad (12)$$

327 where  $\mathbf{A} = \mathbf{J} + \mathbf{H}$  and with  $\mathbf{J} = d\hat{\mathbf{F}}/d\hat{\mathbf{U}}$ . Relation in (8) is now written as

$$\mathbf{J} \cdot \hat{\mathbf{e}}^m - \hat{\lambda}^m \hat{\mathbf{e}}^m = -\mathbf{H} \cdot \hat{\mathbf{e}}^m, \quad (13)$$

328 where  $\hat{\lambda}^m$  and  $\hat{\mathbf{e}}^m$  are the eigenvalues and right eigenvectors of matrix  $\mathbf{A}$ .

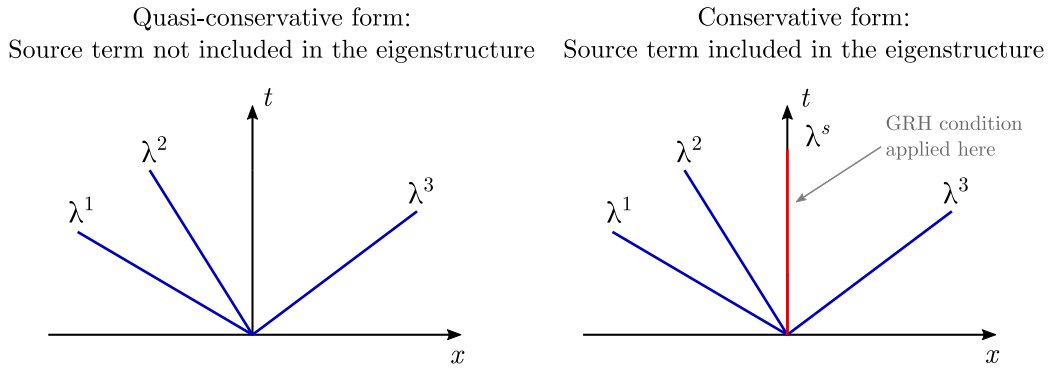


Figure 1: Difference in eigenstructure between the quasi-conservative system (2) and the non-conservative system (9).

329 For the sake of clarity, it is worth recalling that the system in (6) will  
 330 be hereafter referred to as conservative system, the system in (2) as quasi-  
 331 conservative system and the system in (12) as non-conservative system. This  
 332 work focuses on the study of hyperbolic equations with source term, therefore  
 333 (6) will be useless in what follows.

### 334 2.2. Integral relations in discontinuous solutions

335 It is of utmost importance to mention that there exists a certain re-  
 336 lation between the wave speed and the jump of conserved quantities and  
 337 fluxes across the discontinuities carried by the waves. This relation is called  
 338 *Rankine-Hugoniot (RH) condition* or *jump condition*. When dealing with  
 339 non-homogeneous systems of equations, such condition must be extended to  
 340 account for the contribution of the source term, leading to the *Generalized*  
 341 *Rankine-Hugoniot (GRH) condition*.

342 Initial system in (2) is composed of  $N_\lambda$  waves, nevertheless, none of these  
 343 waves are related to the source term and only conventional RH conditions

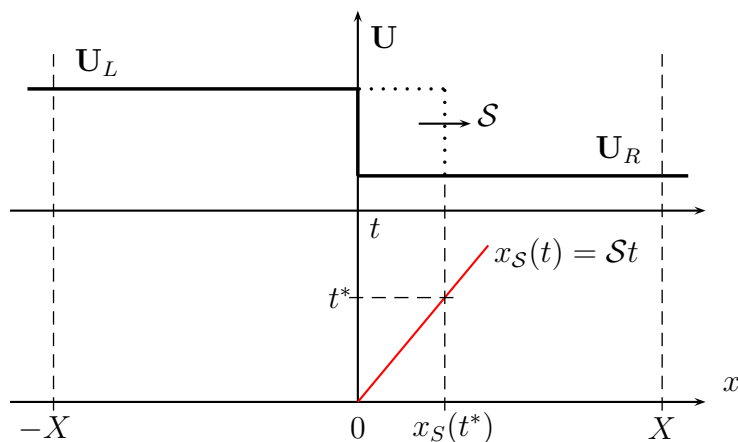


Figure 2: Discontinuity propagation in a non-linear system. The integration domain for the derivation of the Rankine-Hugoniot condition is depicted.

344 could be defined across them. In order to study the more general case,  
 345 where GRH can be defined, it is necessary to express the system in (2) in its  
 346 non-conservative form according to Equation (9). In this way, the system is  
 347 not only characterized by the  $N_\lambda$  eigenvalues associated to the conservative  
 348 fluxes but also by other  $N_S$  eigenvalues, related to extra variables modelling  
 349 the source term, as the dynamics of the source term is included, in some way,  
 350 in the set of characteristic fields. For the sake of simplicity,  $N_S$  is hereafter  
 351 set to 1.

352 The derivation of the GRH condition for the system in (2) with a geo-  
 353 metric source term, or (9) equivalently, can be derived in two different ways.  
 354 The first one would be using equation (2) and considering the source term as  
 355 a Dirac delta that moves with the wave [51]. The second option, the one we  
 356 use here, is to derive the GRH condition from the non-conservative system  
 357 of equations in (9). It is done by integrating (9) over an arbitrary domain  
 358  $[-X, X]$  with  $X$  sufficiently large, as depicted in Figure 2. Notice that the  
 359 displacement of the discontinuity represented in Figure 2 is done from  $t = t_0$   
 360 to  $t = t^* = t_0 + \delta t$ , with  $\delta t$  of differential size. For each  $\lambda^m$  wave defining a  
 361 characteristic field, the left and right states of the solution at each side of the  
 362 discontinuity carried by wave  $\lambda^m$  are denoted by  $\mathbf{U}_L$  and  $\mathbf{U}_R$ , and the speed  
 363 of the discontinuity is denoted by  $\mathcal{S}^m$ . The integral of (9) over  $[-X, X]$  reads

$$\int_{-X}^X \frac{\partial \hat{\mathbf{U}}}{\partial t} dx + \int_{-X}^X \frac{\partial \hat{\mathbf{F}}}{\partial x} dx + \int_{-X}^X \mathbf{H} \frac{\partial \hat{\mathbf{U}}}{\partial x} dx = 0. \quad (14)$$

364 Considering that the integration domain does not change in time, Equation  
365 (14) is rewritten as

$$\frac{d}{dt} \int_{-X}^X \hat{\mathbf{U}} dx + [\hat{\mathbf{F}}]_{-X}^X + \int_{-X}^X \mathbf{H} \frac{\partial \hat{\mathbf{U}}}{\partial x} dx = 0. \quad (15)$$

366 If separating the first term on the left hand side of Equation (15) as

$$\frac{d}{dt} \left( \int_{-X}^{x_S(t)} \hat{\mathbf{U}} dx + \int_{x_S(t)}^X \hat{\mathbf{U}} dx \right) = \frac{d}{dt} \left( \hat{\mathbf{U}}_L(X + \mathcal{S}^m t) + \hat{\mathbf{U}}_R(X - \mathcal{S}^m t) \right) \quad (16)$$

367 and taking the time derivative of the previous result, Equation (16) is rewrit-  
368 ten as

$$\frac{d}{dt} \int_{-X}^X \hat{\mathbf{U}} dx = \mathcal{S}^m \left( \hat{\mathbf{U}}_L - \hat{\mathbf{U}}_R \right). \quad (17)$$

369 When combining the results obtained in (15) and (17), the following condition  
370 for the jump is obtained

$$\hat{\mathbf{F}}_R - \hat{\mathbf{F}}_L - \hat{\mathbf{D}} = \mathcal{S}^m \left( \hat{\mathbf{U}}_R - \hat{\mathbf{U}}_L \right), \quad (18)$$

371 where

$$\hat{\mathbf{D}} = - \int_{-X}^X \mathbf{H} \frac{\partial \hat{\mathbf{U}}}{\partial x} dx \quad (19)$$

372 is a suitable approximation of the integral of the source term. Notice that  
373 the case  $\mathbf{D} = 0$  corresponds to the traditional RH condition.

374 When using this formulation, it must be borne in mind that the geometric  
375 variable is known and is considered to only change at fixed positions, that is to  
376 say, discontinuities on the geometric variable remain at a fixed location. This  
377 helps to understand the conditions for the application of the GRH condition.

378 Let us consider a discontinuity traveling at speed  $\mathcal{S}^m \neq 0$ . Application of  
379 the GRH condition in (18) for the geometric variable yields

$$\mathcal{S}^m ([s_g]_R - [s_g]_L) = 0, \quad (20)$$

380 according to (11). It is observed that  $[s_g]_R = [s_g]_L$  for any  $\mathcal{S}^m \neq 0$ , which  
 381 agrees with the aforementioned consideration saying that variations on the  
 382 geometric variable only take place at fixed positions. This implies that

$$\hat{\mathbf{D}} = 0, \quad (21)$$

383 recovering the traditional RH condition

$$\mathbf{F}_R - \mathbf{F}_L = \mathcal{S}^m (\mathbf{U}_R - \mathbf{U}_L) \quad (22)$$

384 for all  $\mathcal{S}^m \neq 0$ . Notice that the vectors of fluxes and variables in (22) do not  
 385 include the source term as its contribution is nil at this point.

386 On the other hand, if  $\mathcal{S}^m = 0$ , application of the GRH condition in (18)  
 387 for the geometric variable yields

$$0 \cdot ([s_g]_R - [s_g]_L) = 0, \quad (23)$$

388 which holds for any combination of  $[s_g]_R$  and  $[s_g]_L$ . Therefore, for  $\mathcal{S}^m = 0$ ,  
 389 the GRH condition always applies and is written as

$$\hat{\mathbf{F}}_R - \hat{\mathbf{F}}_L = \hat{\mathbf{D}}. \quad (24)$$

390 Here, the last component of the equation, corresponding to the source  
 391 variable, is useless again, therefore we can rewrite (24) as

$$\mathbf{F}_R - \mathbf{F}_L = \mathbf{D}, \quad (25)$$

392 with  $\hat{\mathbf{D}} = (\mathbf{D}, 0)^T$  and due to the nature of the source in (3), the integral of  
 393 this source can be expressed as

$$\mathbf{D} = \int_{[s_g]_L}^{[s_g]_R} \mathbf{S}_s d\mathbf{S}_g, \quad (26)$$

394 with  $\delta [\mathbf{S}_g]_L^R$  the jump in the geometric variable across the wave.

395 It is worth recalling that the set of right (left) states that can be connected  
 396 to a given left (right) state by means of a discontinuous solution describe a  
 397 curve in the phase space called Hugoniot Locus (HL), or Generalized Hugo-  
 398 niot Locus (GHL).

399 *2.3. Integral curves and Riemann invariants*

400 Let us consider a hyperbolic system expressed in non-conservative form as  
 401 (9)

$$\frac{\partial \hat{\mathbf{U}}}{\partial t} + \mathbf{A} \frac{\partial \hat{\mathbf{U}}}{\partial x} = 0, \quad (27)$$

402 where matrix  $\mathbf{A}$  can be diagonalized with  $N_\lambda + N_S$  eigenvalues by means  
 403 of  $N_\lambda + N_S$  linearly independent eigenvectors. For the sake of clarity, hat  
 404 symbol in vectors standing for the extended vectors that include the equa-  
 405 tion of the source term is hereafter omitted. Each eigenvalue  $\lambda^m(\mathbf{U})$ , or  
 406 eigenvector  $\mathbf{e}^m(\mathbf{U})$  equivalently, defines a *characteristic field* associated to  
 407 it, for  $m = 1, \dots, N_\lambda + N_S$ . The properties of the characteristic fields will  
 408 provide useful information about the solution. Prior to the analysis of the  
 409 characteristic fields, it is worth introducing the concepts of Integral Curves  
 410 and state space. The state space, or phase plane, is the representation of a  
 411 component of the state vector with respect to the other components. For in-  
 412 stance, if considering a system of  $N_\lambda + N_S = 2$  equations, with  $\mathbf{U} = (u_1, u_2)$ ,  
 413 the state space representation will be given by the representation of  $u_1$ - $u_2$  in  
 414 a Cartesian coordinate system.

415 **Definition 1.** (*Integral Curve*). Let  $\mathbf{U}(\xi)$  be a smooth curve through state  
 416 space parametrized by the scalar  $\xi$ . This curve is said to be an *Integral*  
 417 *Curve (IC)* of the vector field  $\mathbf{e}^m$  if at each point, the tangent vector to the  
 418 curve,  $d\mathbf{U}(\xi)/d\xi$  is an eigenvector of  $\mathbf{J}(\mathbf{U}(\xi))$  corresponding to the eigenvalue  
 419  $\lambda^m(\mathbf{U}(\xi))$ . When considering a particular set of eigenvectors, the *integral*  
 420 *curve* for  $\mathbf{e}^m$  field is given by

$$\frac{d\mathbf{U}(\xi)}{d\xi} = \nu(\xi) \cdot \mathbf{e}^m(\mathbf{U}(\xi)), \quad (28)$$

421 with  $\nu(\xi)$  a constant parameter that depends on the normalization of the  
 422 eigenvectors [51].

423 When analyzing the solution of hyperbolic systems of conservation laws,  
 424 it is observed that the wave pattern present in the solution is related to the  
 425 variation of the characteristic speed,  $\lambda^m(\mathbf{U})$ , along the integral curve of the  
 426 vector field  $\mathbf{e}^m$ . This variation can be expressed as the directional derivative  
 427 of  $\lambda^m(\mathbf{U})$  in the direction of the eigenvector [51]



$$\frac{d}{d\xi}\lambda^m(\mathbf{U}(\xi)) = \nabla_u \lambda^m(\mathbf{U}(\xi)) \cdot \mathbf{e}^m(\mathbf{U}(\xi)). \quad (29)$$

428 When  $\lambda^m(\mathbf{U})$  is constant along the integral curve, that is (29) is equal  
 429 to zero, the characteristic field is said to be *linearly degenerate*. On the  
 430 other hand, if  $\lambda^m(\mathbf{U})$  varies along the integral curve, which means that the  
 431 characteristic curves are compressing or expanding, the characteristic field is  
 432 said to be *genuinely nonlinear*.

433 Along each integral curve, there are certain quantities that remain con-  
 434 stant. Such quantities are called Riemann invariants.

435 **Definition 2.** (*Riemann invariant*). *The scalar  $w^m$  is said to be a  $m$ -*  
 436 *Riemann invariant when*

$$\nabla_u w^m(\mathbf{U}) \cdot \mathbf{e}^m(\mathbf{U}) \neq 0, \quad \forall \mathbf{U} \in \mathcal{C}, \quad (30)$$

437 with  $\mathcal{C} \subseteq \mathbb{R}^{N_\lambda}$  and where  $\nabla_u$  stands for the gradient with respect to the  
 438 components of vector  $\mathbf{U}$ .

#### 439 2.4. The solution of non-linear hyperbolic systems

440 Non-linear hyperbolic systems of the type of (2) admit complex solutions  
 441 including shocks, rarefaction waves or contact waves. For the sake of brevity,  
 442 the latter are only described here, as they have important implications in  
 443 the design of numerical schemes in presence of geometric source terms. A  
 444 more detailed study on shocks and rarefactions can be found in [52]. Contact  
 445 waves in conservative and non-conservative systems are described below:

446 • **Contact wave in conservative (homogeneous) systems:** If  $\lambda^m$   
 447 defines a *linearly degenerate field* and the following conditions apply:

448 – RH condition:

$$\mathbf{F}(\mathbf{U}_L) - \mathbf{F}(\mathbf{U}_R) = \mathcal{S}^m (\mathbf{U}_L - \mathbf{U}_R) \quad (31)$$

449 – Parallel characteristic condition:

$$\lambda^m(\mathbf{U}_L) = \mathcal{S}^m = \lambda^m(\mathbf{U}_R) \quad (32)$$

450 – Conservation of the Riemann Invariants across the discontinuity.

451 then left and right states  $\mathbf{U}_L$  and  $\mathbf{U}_R$  will be connected by a single  
 452 jump discontinuity wave of speed  $\mathcal{S}^m$  called contact wave.

- 453 • **Contact wave in non-conservative systems** (with geometric source  
 454 term) where the relevant eigenvalue does not depend upon  $\mathbf{U}$  [14]:

455 The presence of contact discontinuities in RPs given by non-homogeneous  
 456 systems of conservation laws has to be taken into account when con-  
 457 structing augmented solvers. In this work, we consider contact waves  
 458 whose relevant eigenvalue does not depend upon  $\mathbf{U}$ . This would be the  
 459 case of a system like (9) where  $\mathbf{H}\mathbf{U}_x$  includes the contribution of the  
 460 geometric source term (3). For such case, given a initial left state,  $\mathbf{U}_L$ ,  
 461 the right state, hereafter denoted by  $\mathbf{U}(\xi)$ , does not necessarily lie on  
 462 the integral curve, while it will always be related to the left state by  
 463 means of the GRH condition [4, 14], as all discontinuous solutions do  
 464 satisfy this relation. Recall that  $\mathbf{U}_L = \mathbf{U}(\xi = 0)$ .

465 Let us consider the non-conservative system in (9) and assume that the  
 466  $m$ -th characteristic field, associated to eigenvalue  $\lambda^m$  and eigenvector  
 467  $\mathbf{e}^m$ , is linearly degenerate. Then, the associated contact wave is given  
 468 by

$$\mathbf{U}(x, t) = \begin{cases} \mathbf{U}_L & x < \mathcal{S}^m t \\ \mathbf{U}(\xi) & x > \mathcal{S}^m t \end{cases} \quad (33)$$

469 with constant speed  $\mathcal{S}^m = \lambda^m(\mathbf{U}(\xi)) = \lambda^m(\mathbf{U}_L)$ . All possible  $\mathbf{U}(\xi)$   
 470 states can be found by means of the GHL. From (18) we have

$$\mathbf{F}(\mathbf{U}(\xi)) - \mathbf{F}(\mathbf{U}_L) - \mathcal{S}^m(\mathbf{U}(\xi) - \mathbf{U}_L) = \mathbf{D}. \quad (34)$$

471 In this way,  $\mathbf{U}(\xi)$  will satisfy the GRH condition, however, we have  
 472 not imposed yet any condition for the conservation of the relevant  $m$ -  
 473 Riemann invariants across the contact discontinuity, hence IC and GHL  
 474 may not coincide. To find the condition so that such sets of states  
 475 coincide, following [14], let us consider the differential form of (34)

$$\frac{d}{d\xi} [\mathbf{F}(\mathbf{U}(\xi)) - \mathcal{S}^m \mathbf{U}(\xi)] = \frac{d}{d\xi} \mathbf{D} \quad (35)$$

476 that can be rewritten as

$$\frac{d\mathbf{F}}{d\mathbf{U}} \frac{d\mathbf{U}(\xi)}{d\xi} - \mathcal{S}^m \frac{d\mathbf{U}(\xi)}{d\xi} = \frac{d}{d\xi} \mathbf{D}. \quad (36)$$

477 To enforce the solution to lie on both the IC and the GHL, we set  
 478  $\mathbf{U} = \mathbf{U}^m(\xi)$  to be the set of states lying on the IC according to (28),  
 479 yielding

$$\mathbf{J} \frac{d\mathbf{U}^m(\xi)}{d\xi} - \mathcal{S}^m \frac{d\mathbf{U}^m(\xi)}{d\xi} = \frac{d}{d\xi} \mathbf{D}, \quad (37)$$

480 where  $d\mathbf{U}^m(\xi)/d\xi$  can be substituted by  $\mathbf{e}^m$  as the solution follows the  
 481 IC, and  $\mathcal{S}^m$  by  $\lambda^m$ , leading to

$$\mathbf{J} \cdot \mathbf{e}^m - \lambda^m \cdot \mathbf{e}^m = \frac{d}{d\xi} \mathbf{D}, \quad (38)$$

482 that can be rewritten by means of (13) as

$$-\mathbf{H} \cdot \mathbf{e}^m = \frac{d}{d\xi} \mathbf{D}. \quad (39)$$

483 Only when relation in (39) is satisfied, the IC and GHL coincide and  
 484 the Riemann invariants are conserved across the contact wave. This  
 485 property will be used later to design an E-scheme for the SWE.

### 486 3. Finite volume discretization

487 In the present framework, problems of interest are defined as initial value  
 488 boundary problems (IVBP) that can be expressed as

$$\left\{ \begin{array}{l} \text{PDEs: } \frac{\partial \mathbf{U}}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{U})}{\partial x} = \mathbf{S} \\ \text{IC: } \mathbf{U}(x, 0) = \mathring{\mathbf{U}}(x) \\ \text{BC: } \mathbf{U}(a, t) = \mathbf{U}_a(t) \quad \mathbf{U}(b, t) = \mathbf{U}_b(t) \end{array} \right. \quad (40)$$

489 defined inside the domain  $[a, b] \times [0, T]$ , with  $\mathring{\mathbf{U}}(x)$  the initial condition and  
 490  $\mathbf{U}_a(t)$  and  $\mathbf{U}_b(t)$  the left and right boundary conditions. When using a first

491 order finite volume approach, the domain is discretized in computational  
 492 cells and the conserved variables and governing equations are integrated in-  
 493 side those cells, leading to algebraic equations that depend upon piecewise  
 494 constant data. In this work, the following computational grid composed of  
 495  $N$  cells is used

$$a = x_{\frac{1}{2}} < x_{\frac{3}{2}} < \dots < x_{N-\frac{1}{2}} < x_{N+\frac{1}{2}} = b, \quad (41)$$

496 as shown in Figure 3, with cells and cell sizes defined as

$$\Omega_i = \left[ x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}} \right], \quad \Delta x_i = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}, \quad i = 1, \dots, N \quad (42)$$

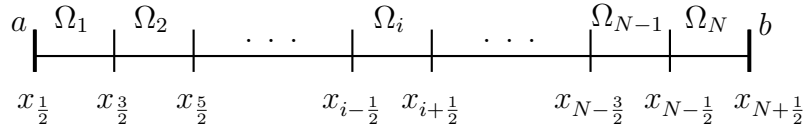


Figure 3: Mesh discretization

497 Inside each cell, conserved quantities at time  $t^n$  are defined as cell averages  
 498 as

$$\mathbf{U}_i^n = \frac{1}{\Delta x_i} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \mathbf{U}(x, t^n) dx, \quad i = 1, \dots, N. \quad (43)$$

499 Following the approach proposed by Godunov, the finite volume dis-  
 500 cretization of the system in (2) inside  $[x_{i-1/2}, x_{i+1/2}] \times [t^n, t^{n+1}]$  is straight-  
 501 forward derived from integration of (2) in this volume, leading to

$$\mathbf{U}_i^{n+1} = \mathbf{U}_i^n - \frac{\Delta t}{\Delta x} [\mathbf{F}_{i+1/2}^- - \mathbf{F}_{i-1/2}^+], \quad (44)$$

502 where  $\mathbf{F}_{i+1/2}^-$  and  $\mathbf{F}_{i-1/2}^+$  are the numerical fluxes, which are computed solv-  
 503 ing the Riemann Problems (RPs) at the interfaces by means of a suitable  
 504 Riemann solver.

505 Analogously, equation (44) can be rewritten in terms of fluctuations, gen-  
 506 erally denoted by  $\delta \mathbf{M}$ , leading to

$$\mathbf{U}_i^{n+1} = \mathbf{U}_i^n - \frac{\Delta t}{\Delta x} [\delta \mathbf{M}_{i+1/2}^- + \delta \mathbf{M}_{i-1/2}^+], \quad (45)$$

507 where

$$\begin{aligned}\delta\mathbf{M}_{i+1/2}^- &= \mathbf{F}_{i+1/2}^- - \mathbf{F}_i, \\ \delta\mathbf{M}_{i-1/2}^+ &= \mathbf{F}_i - \mathbf{F}_{i-1/2}^+, \end{aligned}\tag{46}$$

508 represent the contribution of the incoming waves to the right and left edges,  
509 respectively. The Riemann solver selected here is called the augmented Roe  
510 Riemann solver (ARoe) and is detailed in Appendix A.

#### 511 4. Application to the Shallow Water Equations (SWE)

512 The SWE can be expressed in matrix form as

$$\frac{\partial\mathbf{U}}{\partial t} + \frac{\partial\mathbf{F}(\mathbf{U})}{\partial x} = \mathbf{S}.\tag{47}$$

513 where

$$\mathbf{U} = \begin{pmatrix} h \\ hu \end{pmatrix}, \quad \mathbf{F} = \begin{pmatrix} hu \\ hu^2 + \frac{1}{2}gh^2 \end{pmatrix}, \quad \mathbf{S} = \begin{pmatrix} 0 \\ S_z \end{pmatrix},\tag{48}$$

514 where  $h$  is the water depth,  $u$  is the depth averaged velocity,  $hu$  the discharge  
515 and  $g$  is the acceleration of gravity. The source term  $S_z$  involves the variations  
516 in bed geometry  $S_z$

$$S_z = -gh\frac{dz}{dx},\tag{49}$$

517 where  $z$  stands for the bed elevation.

518 In order to design a suitable numerical scheme that mimics the physical  
519 behavior of (47), these equations must be thoroughly analyzed. In physics,  
520 invariance of certain quantities is usually present in systems. In the SWE,  
521 the mechanical energy is an example. From the analysis of (47) under steady  
522 regime and considering a smooth solution, we obtain that

$$\frac{\partial}{\partial x} \left( \frac{u^2}{2g} + h + z \right) = 0,\tag{50}$$

523 where  $E = \frac{u^2}{2g} + h + z$  is the specific mechanical energy. By looking at this  
524 quantity when designing the numerical scheme, the well-balanced property  
525 can be extended to the so-called energy-balanced property, which allows the

526 numerical scheme to provide the exact solution in steady cases with moving  
 527 water.

528 It is worth pointing out that, unlike in previous publications [14], the  
 529 authors in this work are faithful to the original system in (47) and do not  
 530 include any dissipation mechanism (for instance, across shocks), as the orig-  
 531 inal equations do not consider extra friction terms. When neglecting shear  
 532 stress, dissipation will only take place in certain conditions, such as a sudden  
 533 change of flow regime, according to the physical behavior described by the  
 534 original equations.

535 For system in (47), the discretization of the source term is not a triv-  
 536 ial task and additional information must be taken into account in order to  
 537 construct a trustworthy numerical solution and eventually obtain an energy-  
 538 balanced scheme. The analysis of the system of equations in non-conservative  
 539 form is useful to this end as it provides information on the physical nature  
 540 of the additional wave associated to the source term.

#### 541 4.1. Characteristic analysis of the SWE system in its non-conservative form

542 According to Equation (9), system in (47) can be expressed in non-  
 543 conservative form

$$\frac{\partial \mathbf{U}}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{U})}{\partial x} + \mathbf{H}(\mathbf{U}) \frac{\partial \mathbf{U}}{\partial x} = 0, \quad (51)$$

544 where

$$\mathbf{U} = \begin{pmatrix} h \\ hu \\ z \end{pmatrix}, \quad \mathbf{F} = \begin{pmatrix} hu \\ hu^2 + \frac{1}{2}gh^2 \\ 0 \end{pmatrix}, \quad \mathbf{H} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & gh \\ 0 & 0 & 0 \end{pmatrix}. \quad (52)$$

545 The Jacobian matrix of the flux reads

$$\mathbf{J} = \begin{pmatrix} 0 & 1 & 0 \\ c^2 - u^2 & 2u & 0 \\ 0 & 0 & 0 \end{pmatrix}, \quad (53)$$

546 and it can be used to construct the following matrix

$$\mathbf{A} = \mathbf{J} + \mathbf{H} = \begin{pmatrix} 0 & 1 & 0 \\ c^2 - u^2 & 2u & gh \\ 0 & 0 & 0 \end{pmatrix}, \quad (54)$$

547 allowing to express the system in quasilinear form. The eigenvalues and  
 548 eigenvectors that diagonalize  $\mathbf{A}$  are given by

$$\lambda^1 = u - c, \quad \lambda^S = 0, \quad \lambda^2 = u + c \quad (55)$$

549 and

$$\mathbf{e}^1 = \begin{pmatrix} 1 \\ \lambda^1 \\ 0 \end{pmatrix}, \quad \mathbf{e}^S = \begin{pmatrix} 1 \\ 0 \\ u^2/gh - 1 \end{pmatrix}, \quad \mathbf{e}^2 = \begin{pmatrix} 1 \\ \lambda^2 \\ 0 \end{pmatrix}. \quad (56)$$

550 For the sake of clarity and consistency throughout the text, the charac-  
 551 teristic field corresponding to the source variable,  $z$ , is denoted by  $S$  while  
 552 the two other fields are denoted by 1 (for the left moving wave) and 2 (for the  
 553 right moving wave). The nature of each characteristic field can be studied  
 554 as pointed out in Section 2.3. Following definition in (29), for this particular  
 555 case we have

$$\begin{aligned} \nabla_u \lambda^1(\mathbf{U}) \cdot \mathbf{e}^1(\mathbf{U}) &= -\frac{\sqrt{g}}{2\sqrt{h}}, \\ \nabla_u \lambda^S(\mathbf{U}) \cdot \mathbf{e}^S(\mathbf{U}) &= 0, \\ \nabla_u \lambda^2(\mathbf{U}) \cdot \mathbf{e}^2(\mathbf{U}) &= \frac{\sqrt{g}}{2\sqrt{h}}, \end{aligned} \quad (57)$$

556 noticing that the  $S$ -characteristic field associated to the bed step is linearly  
 557 degenerate as the eigenvalue  $\lambda^S$  is zero  $\forall \mathbf{U}$  (the step is regarded as a sta-  
 558 tionary discontinuity) while the 1 and 2-characteristic fields are genuinely  
 559 nonlinear.

560 The integral curve for each of the characteristic fields can be derived  
 561 from equation (28). The integral curve associated to the 1-characteristic  
 562 field, parametrized by  $\xi$  and starting at  $(h, hu, z) = (h^*, (hu)^*, z^*)$ , reads

$$\mathbf{U}^1(\xi) = \begin{pmatrix} h(\xi) \\ hu(\xi) \\ z(\xi) \end{pmatrix} = \begin{pmatrix} h^* + \xi \\ (h^* + \xi) \left[ u^* - 2(\sqrt{g(h^* + \xi)} - \sqrt{gh^*}) \right] \\ z^* \end{pmatrix}. \quad (58)$$

563 Similarly, the integral curve for the 2-characteristic field can be calculated,  
 564 obtaining the conjugated of (58). It is more interesting to analyze the result  
 565 for the  $S$ -characteristic field, that reads

$$\mathbf{U}^S(\xi) = \begin{pmatrix} h(\xi) \\ hu(\xi) \\ z(\xi) \end{pmatrix} = \begin{pmatrix} h^* + \xi \\ (hu)^* \\ \frac{u^{*2}}{2g} + z^* - \frac{(hu)^{*2}}{2g(h^* + \xi)^2} - \xi \end{pmatrix}, \quad (59)$$

566 as it can be given a physical meaning. One can realize that the third equation  
 567 in vector (59), in combination with the first and second equations, stands for  
 568 the conservation of the specific mechanical energy across the contact wave.  
 569 Such an idea can be more generally conveyed by saying that the Riemann  
 570 invariants of the  $S$ -characteristic field are the discharge and the mechanical  
 571 energy. In Table 1, the Riemann invariants for all waves are presented.

Characteristic field	1-Riemann invariant	2-Riemann invariant
1	$u + 2\sqrt{gh}$	$z$
$S$	$hu$	$\frac{u^2}{2g} + h + z$
2	$u - 2\sqrt{gh}$	$z$

Table 1: Summary of Riemann invariants for the non-homogeneous SWE.

#### 572 4.2. Conservation of energy across the bed-step contact wave

573 As outlined in the previous section, the  $S$ -characteristic field in the non-  
 574 conservative SWE in (52) is a linearly degenerate field. This kind of field  
 575 arises from the presence of the bed step and is characterized by a contact  
 576 wave of zero celerity,  $\lambda^S = 0$ , since the bed elevation does not vary in time.

577 Discontinuous solutions describing a contact wave are generally expressed  
 578 by (33). For this particular case, the right state will be denoted by  $\mathbf{U}_R$ , hence  
 579 (33) is rewritten as

$$\mathbf{U}(x, t) = \begin{cases} \mathbf{U}_L & x < 0 \\ \mathbf{U}_R & x > 0 \end{cases} \quad (60)$$

580 where  $\mathbf{U}_L = (h_L, (hu)_L, z_L)^T$  and  $\mathbf{U}_R = (h_R, (hu)_R, z_R)^T$  are the left and  
 581 right states respectively. Notice that we may write  $(hu)_L = h_L u_L$  for the  
 582 sake of clarity and recall that this quantity represents the first Riemann  
 583 invariant of the  $S$ -characteristic field, hence  $h_L u_L = h_R u_R$ . The second  
 584 Riemann invariant is the specific mechanical energy, hence  $u_L^2/2 + g(h+z)_L =$   
 585  $u_R^2/2 + g(h+z)_R$ .

586 Across the contact wave in (60), the Generalized Rankine-Hugoniot (GRH)  
 587 condition in (24) must hold for all variables. For this particular case, it reads



$$\begin{aligned}
& h_R u_R - h_L u_L = 0, \\
& \left( g \frac{h_R^2}{2} + h_R u_R^2 \right) - \left( g \frac{h_L^2}{2} + h_L u_L^2 \right) = D,
\end{aligned} \tag{61}$$

588 with  $D$  a suitable approximation of the integral of the source term across the  
589 bed step

$$D = - \int_{z_L}^{z_R} gh dz, \tag{62}$$

590 that can be rewritten as

$$D = - \int_{x_L}^{x_R} gh \frac{dz}{dx} dx. \tag{63}$$

591 As outlined before, GRH condition in (61) must be ensured so that  
592 (60) is a weak solution of the problem, hence the right state  $(h_R, h_R u_R, z_R)$   
593 must lie on the Generalized Hugoniot Locus (GHL) for a given left state  
594  $(h_L, h_L u_L, z_L)$ . However, this condition does not ensure the conservation of  
595 Riemann invariants across the contact wave. Only when condition in (39)  
596 holds, Riemann invariants are conserved and the IC coincide with the GHL.  
597 In other words, we can state that the Integral Curve (IC) coincide with the  
598 GHL if (61) holds and the Riemann invariants of the  $S$ -field in Table 1 are  
599 conserved.

600 It seems clear that the election of a suitable discretization of the integral of  
601 the source term in (63) is crucial. In [14], a particular STD based on physical  
602 considerations that accounts for the dissipation of energy across the step was  
603 chosen. Under this assumption, they showed that equation (39) is not always  
604 satisfied and proved that the Riemann invariant associated to the specific  
605 mechanical energy was not anymore conserved across the step. In this way,  
606 they provided a coherent mathematical framework for the physically-based  
607 dissipative discretization of the bed step and they constructed a Riemann  
608 solver based on such ideas.

609 Unlike [14], in the present work the authors do not include any additional  
610 energy dissipation mechanism. Here, an energy-conservative STD is sought,  
611 hence both the GRH condition and Equation (39) must hold, as Riemann  
612 invariants have to be conserved across the contact wave. Following [14],  
613 equation (39) is rewritten as

$$-\int_0^{\hat{\xi}} \mathbf{H} \cdot \mathbf{e}^S d\xi = \mathbf{D}, \quad (64)$$

614 where  $\hat{\xi} = h_R - h_L$  is the value of  $\xi$  on the right state. We define

$$h(\hat{\xi}) = h_R \quad u(\hat{\xi}) = u_R \quad z(\hat{\xi}) = z_R. \quad (65)$$

615 Our goal here is to find the expression for  $\mathbf{D}$  satisfying (64) and to this end,  
 616 we have to manipulate (64) using extra relations among left and right states.  
 617 It is worth recalling that for the derivation of condition (64) (originally (39)),  
 618  $\mathbf{U}(\xi)$  was imposed to lie on the IC, given by Equation (59). Here we will  
 619 work under the same assumption, hence  $\mathbf{U}(\hat{\xi}) = \mathbf{U}_R = (h_R, h_R u_R, z_R)$  lies on  
 620 the IC for a given left state. Water depth along the IC can be expressed as

$$h(\hat{\xi}) = h_L + \hat{\xi} = h_R \quad (66)$$

621 and in the same way, the velocity along the IC is

$$u(\hat{\xi}) = \frac{h_L u_L}{h_L + \hat{\xi}} = \frac{h_R u_R}{h_R} = u_R, \quad (67)$$

622 with a constant discharge

$$q = hu(\hat{\xi}) = h_L u_L = h_R u_R, \quad (68)$$

623 also denoted by  $q$ , and a variable bed elevation along the IC

$$z(\hat{\xi}) \equiv z_R = z_L + h_L - h_R + \frac{u_L^2}{2g} - \frac{u_R^2}{2g}. \quad (69)$$

624 In the following derivation, condition (64) will be combined with the rela-  
 625 tions between left and right states in (66)-(69), allowing to find the expression  
 626 of  $\mathbf{D}$  satisfying the RI and the GRH conditions. The product  $\mathbf{H} \cdot \mathbf{e}^S$  reads

$$\mathbf{H} \cdot \mathbf{e}^S = \begin{pmatrix} 0 \\ u^2(\xi) - gh(\xi) \\ 0 \end{pmatrix} \quad (70)$$

627 and using (67) in (70), the latter yields

$$-\int_0^{\hat{\xi}} \begin{pmatrix} 0 \\ \left(\frac{h_L u_L}{h_L + \xi}\right)^2 - g(h_L + \xi) \\ 0 \end{pmatrix} d\xi = \begin{pmatrix} 0 \\ D \\ 0 \end{pmatrix}. \quad (71)$$

628 From (71), only the second component will be considered

$$-\int_0^{\hat{\xi}} \left(\frac{h_L u_L}{h_L + \xi}\right)^2 d\xi + \int_0^{\hat{\xi}} g(h_L + \xi) d\xi = D. \quad (72)$$

629 Integrating (72) and using the relation  $h_L u_L = h_R u_R$  in (68) when required,  
630 it yields

$$\left(g \frac{h_R^2}{2} + h_R u_R^2\right) - \left(g \frac{h_L^2}{2} + h_L u_L^2\right) = D, \quad (73)$$

631 with the right state laying on the IC in (59). It can be noticed that equation  
632 (73) coincides with the GRH condition for the conservation of momentum.

633 Now, combination of equation (73) with (69) allows to derive the particu-  
634 lar STD,  $D$ , that under the assumed hypotheses will ensure the conservation  
635 of the Riemann invariants and lead to an energy-conservative scheme. For  
636 the sake of clarity, equation (73) is rewritten as

$$\delta \left( g \frac{h^2}{2} + h u^2 \right)_{L,R} = D \quad (74)$$

637 and so is (69), the equation for the conservation of energy

$$\delta \left( \frac{u^2}{2} + g(h + z) \right)_{L,R} = 0 \quad (75)$$

638 where  $\delta(\cdot)_{L,R} = (\cdot)_R - (\cdot)_L$  is a difference operator. From (74), it is straight-  
639 forward to obtain

$$(g\bar{h}\delta h + \bar{u}\delta(hu) + \bar{h}\bar{u}\delta u)_{L,R} = D, \quad (76)$$

640 where

$$\bar{(\cdot)}_{L,R} = \frac{(\cdot)_L + (\cdot)_R}{2} \quad (77)$$

641 is an average operator. For the sake of simplicity, subscript  $(\cdot)_{L,R}$  is dropped  
 642 in Equations (78)-(82) as they always refer to the left and right states of the  
 643 contact wave in this derivation. Noticing that  $\delta(hu)_{L,R} = h_R u_R - h_L u_L = 0$ ,  
 644 Equation (76) yields

$$g\bar{h}\delta h + \bar{h}\bar{u}\delta u = D. \quad (78)$$

645 The equation for the conservation of energy in (75) is multiplied by  $\bar{h}$  and  
 646 rewritten as

$$\bar{h}\bar{u}\delta u + g\bar{h}\delta h + g\bar{h}\delta z = 0, \quad (79)$$

647 from where the term  $g\bar{h}\delta h$  can be expressed as

$$g\bar{h}\delta h = -\bar{h}\bar{u}\delta u - g\bar{h}\delta z \quad (80)$$

648 and can be inserted in (78), leading to

$$D = -g\bar{h}\delta z + (\bar{h}\bar{u} - \bar{h}\bar{u})\delta u. \quad (81)$$

649 It is straightforward to show that (81) can be rewritten as

$$D = -g\bar{h}\delta z + \delta(hu^2) - \bar{u}\delta(hu) - \bar{h}\delta\left(\frac{1}{2}u^2\right), \quad (82)$$

650 with  $\delta(hu) = 0$  according to the GRH conditions, hence

$$D = -g\bar{h}\delta z + \delta(hu^2) - \bar{h}\delta\left(\frac{1}{2}u^2\right). \quad (83)$$

651 As outlined before, weak solutions for the bed step contact wave are  
 652 always required to satisfy the GRH condition. That is to say, for a given  
 653 left state, the right state is calculated using (61). When the discretization  
 654 of the source term in (63),  $D$ , is undefined, there are infinite solutions for  
 655 the right state and only when choosing a particular discretization, the right  
 656 state can be determined. Unlike the approach proposed in [14] where the  
 657 authors impose a particular STD based on energy dissipation hypothesis, here  
 658 the expression for the discretization of the source term is derived imposing  
 659 the equivalence between GHL and IC. To this end, apart from the GRH  
 660 condition, we require an extra condition given by (39) in order to ensure  
 661 the constancy of Riemann invariants across the wave. Notice that such a  
 662 condition consists of the equation for the conservation of energy provided by  
 663 the IC.

664 *4.3. Numerical discretization of the source term at cell interfaces for aug-*  
 665 *mented solvers*

666 When using augmented solvers, such as the HLLS and ARoe solvers, numerical approximations over the integral of the source term at cell interfaces  
 667 are required. The approximation of the spatial integral of the source term at  
 668 cell interface  $i + 1/2$ , that is inside  $[x_i, x_{i+1}]$ , will be referred to as  
 669

$$\int_{x_i}^{x_{i+1}} -g h \frac{dz}{dx} dx \approx \bar{S}_{i+1/2}. \quad (84)$$

670 We can find in the literature different numerical approaches for Equation  
 671 (84), however, this choice is not trivial since most of such approaches  
 672 are not able to ensure a numerical solution that converges to a physically  
 673 based solution with mesh refinement, even when using high order schemes.  
 674 This problem is put into evidence when looking, for instance, at the discrete  
 675 energy level or at the shock positioning given by the numerical scheme. In  
 676 this section, four different source term discretizations are described. Two of  
 677 them, the differential formulation (DF) and the integral formulation (IF),  
 678 are traditional approaches, which are easy to program and exhibit an over-  
 679 all acceptable performance but they are not able to ensure conservation of  
 680 energy. Moreover, the IF does not allow the numerical scheme to converge  
 681 to the exact shock position, for steady shocks, with mesh refinement. The  
 682 other two STDs described here, in contrast, are energy balanced discretiza-  
 683 tions, that is to say, they allow the numerical scheme to preserve the discrete  
 684 level of energy (when required) and to dissipate the exact amount of energy  
 685 in presence of hydraulic jumps. Such techniques are called weighted energy  
 686 balanced formulation (WEBF) and the selective energy balanced method  
 687 (SEBF) and whereas the former is still not able to make the scheme con-  
 688 verge to the exact position of the hydraulic jump under steady regime, the  
 689 latter does, as it will be shown in the following section. Therefore, among  
 690 the four techniques described here, only the SEBF which is presented here  
 691 for the first time, is well suited for both energy conservation and accurate  
 692 shock capturing.

693 One possibility is to compute it considering a smooth variation of the  
 694 variables across the interface, as

$$\bar{S}_{i+1/2}^{DF} = -g\bar{h}\delta z, \quad (85)$$

695 which will be referred to as differential formulation (DF) and where

$$\bar{h} = \frac{1}{2}(h_{i+1} + h_i), \quad \delta z = z_{i+1} - z_i. \quad (86)$$

696 The second possibility is the so-called integral formulation (IF), derived from  
 697 the integration of the pressure along the bottom step for a piecewise constant  
 698 data reconstruction of the bed elevation,  $z$ . If assuming that the pressure  
 699 distribution is hydrostatic over the step and depends only on the free-surface  
 700 level on the side of the discontinuity where the bottom elevation is lower, the  
 701 source term is evaluated explicitly at  $t = 0$  as [11]

$$\bar{S}_{i+1/2}^{IF} = -g \left( h_j - \frac{|\delta z'|}{2} \right)_{i+\frac{1}{2}} \delta z'_{i+\frac{1}{2}}, \quad (87)$$

702 where  $z$  is the bed level surface, and  $j$  and  $\delta z'$  are given by

$$j = \begin{cases} i & \text{if } \delta z_{i+\frac{1}{2}} \geq 0 \\ i+1 & \text{if } \delta z_{i+\frac{1}{2}} < 0 \end{cases} \quad \delta z' = \begin{cases} h_i & \text{if } \delta z_{i+\frac{1}{2}} \geq 0 \text{ and } d_i < z_{i+1} \\ -h_{i+1} & \text{if } \delta z_{i+\frac{1}{2}} < 0 \text{ and } d_{i+1} < z_i \\ \delta z & \text{otherwise} \end{cases} \quad (88)$$

703 and  $d = h + z$  is the water level surface.

704 In cases of still water with a continuous water level surface, both (85)  
 705 and (87) do ensure quiescent equilibrium. In this particular case hydrostatic  
 706 forces are exactly balanced.

707 In order to extend the well-balanced property for static equilibrium to  
 708 the energy-balanced property, that ensures the exact conservation of energy  
 709 in steady cases with moving water, it is necessary to impose extra conditions  
 710 in the discretization of the source term. Generally, under the assumption  
 711 of conservation of energy across the bed step contact wave, the best choice  
 712 for the discretization of the bed source term seems to be Equation (81).  
 713 However, such a discretization does not allow to construct an explicit scheme  
 714 as it depends upon the intermediate states at both sides of the bed step,  $\mathbf{U}_i^-$   
 715 and  $\mathbf{U}_{i+1}^+$ .

716 Under steady conditions and considering no change in flow regime across  
 717 the RP, it is straightforward to prove that  $\mathbf{U}_i = \mathbf{U}_i^-$  and  $\mathbf{U}_{i+1} = \mathbf{U}_{i+1}^+$ , hence  
 718 (81) can be rewritten in terms of the initial data as

$$\begin{aligned}
D = & -g \left( \frac{h_{i+1} + h_i}{2} \right) (z_{i+1} - z_i) + \\
& \left[ \left( \frac{(hu)_{i+1} + (hu)_i}{2} \right) - \left( \frac{h_{i+1} + h_i}{2} \right) \left( \frac{u_{i+1} + u_i}{2} \right) \right] (u_{i+1} - u_i).
\end{aligned} \tag{89}$$

719 For the sake of clarity, notation for Equation (89) is simplified, considering  
720 variations and averages across the interface  $i + 1/2$ , that is, the left and right  
721 states of the RP. By doing this, (89) is rewritten as

$$D = \{ -g\bar{h}\delta z + (\overline{hu} - \bar{h}\bar{u})\delta u \}_{i+1/2}. \tag{90}$$

722 In shallow flows, there are physically feasible situations where energy is  
723 dissipated, such as hydraulic jumps. Ideally, such a shock would be consid-  
724 ered as a pure discontinuity where energy is suddenly dissipated, however,  
725 when using a finite volume formulation, the shock width is of the size of a  
726 cell, since the discretization considers constant values in each cell and the  
727 discontinuity cannot be kept anymore as a discontinuity inside a cell. As  
728 a consequence, energy dissipation must be imposed at the interfaces of the  
729 cell containing the shock, as it is not possible to explicitly carry out the  
730 dissipation of energy inside the cell.

731 Murillo [25] proposed a novel approach for the discretization of the source  
732 term that allows to construct an exactly energy balanced scheme. This ap-  
733 proximation is based on the principle of conservation of mechanical energy  
734 and is only applied to the leading term, since higher order terms become nil  
735 in steady state when energy is conserved, as mentioned above.

736 Considering the IF and DF approaches for the discretization of the source  
737 term, it is possible to evaluate  $\bar{S}_{i+1/2}$  as a combination of them as

$$\bar{S}_{i+1/2} = (1 - \mathcal{A})S_{i+1/2}^{DF} + \mathcal{A}S_{i+1/2}^{IF}, \tag{91}$$

738 where  $0 \leq \mathcal{A} \leq 1$ . This formulation will be referred to as weighted energy  
739 balanced formulation (WEBF). In order to satisfy both energy and momen-  
740 tum conservation under steady conditions, a value  $\mathcal{A}_E$  is defined as

$$\mathcal{A}_E = \frac{\delta(hu^2) - \bar{h}\delta\left(\frac{u^2}{2}\right)}{S_{i+1/2}^{IF} - S_{i+1/2}^{DF}}, \tag{92}$$

741 according to [25]. Coefficient  $\mathcal{A}_E$  can be used in (91) to ensure the conser-  
 742 vation of energy for smooth solutions. On the other hand, when considering  
 743 transcritical jumps, energy must be dissipated, hence the value of weight  
 744 coefficient  $\mathcal{A}$  in (91) is set to 1. Considering these situations, the complete  
 745 algorithm for the calculation of  $\mathcal{A}$  reads [25]

$$\mathcal{A} = \begin{cases} 1 & \text{if } u_{i+1}u_i > 0 \text{ and } u_i > 0 \text{ and } |Fr_{i+1}| < 1 \text{ and } |Fr_i| > 1 \\ 1 & \text{if } u_{i+1}u_i > 0 \text{ and } u_i < 0 \text{ and } |Fr_{i+1}| > 1 \text{ and } |Fr_i| < 1 \\ \mathcal{A}_E & \text{otherwise} \end{cases} \quad (93)$$

746 where  $Fr_i$  and  $Fr_{i+1}$  are the Froude numbers on the left and right sides of the  
 747 interface. It is worth pointing out that  $\mathcal{A}_E$  can be straightforwardly obtained  
 748 from Equation (90).

749 On the other hand, instead of imposing the exact amount of dissipation  
 750 of energy across the shock by means of a tailored STD at that point, in  
 751 this work we propose to add an additional degree of freedom to the equa-  
 752 tions by means of using a traditional discretization of the source term at the  
 753 interfaces surrounding the hydraulic jump while maintaining the energy con-  
 754 servative formulation in (90) for the rest. The differential discretization of  
 755 the source term is chosen at those interfaces. This technique allows the nu-  
 756 merical scheme to converge to the exact position of the shock while recovering  
 757 the exact solution in both the subcritical and supercritical regions connected  
 758 by the transcritical shock, with independence of the grid refinement.

759 The proposed approach is next explained. We propose to use Roe celerities,  
 760  $\tilde{\lambda}^m$  to identify the cell containing the hydraulic jump, since it is known  
 761 that both celerities at the left interface are positive (supercritical flow enter-  
 762 ing the cell) while the celerities at the right interface correspond to subcriti-  
 763 cal conditions (one negative and the other one positive). Let us consider the  
 764 cells,  $\Omega_i$ , as single cells contained in the computational domain  $\Omega$  such that  
 765  $\Omega = \{\Omega_i \mid i \in [1, \dots, N]\}$ . Considering the possibility of multiple hydraulic  
 766 jumps within the domain, we denote the set of cells containing a positive-  
 767 flow hydraulic jump as

$$\mathcal{D}^+ = \left\{ \Omega_i \mid \Omega_i \in \Omega \wedge \tilde{\lambda}_{i-1/2}^1 \cdot \tilde{\lambda}_{i+1/2}^1 < 0 \wedge h_{i-1} < h_{i+1} \right\} \quad (94)$$

768 and the set of cells containing a negative-flow hydraulic jump as

$$\mathcal{D}^- = \left\{ \Omega_i \mid \Omega_i \in \Omega \wedge \tilde{\lambda}_{i-1/2}^2 \cdot \tilde{\lambda}_{i+1/2}^2 < 0 \wedge h_{i-1} > h_{i+1} \right\} \quad (95)$$



769 and the set of Riemann Problems at the left and right interfaces of cells  
 770  $\Omega_i \in \mathcal{D}^+ \cup \mathcal{D}^-$

$$771 \quad \mathcal{R}_1 = \{\text{RP}_{i+1/2} \mid i \in \mathbb{N} \wedge \Omega_i \in \mathcal{D}^+ \cup \mathcal{D}^-\} \quad (96)$$

$$\mathcal{R}_2 = \{\text{RP}_{i-1/2} \mid i \in \mathbb{N} \wedge \Omega_i \in \mathcal{D}^+ \cup \mathcal{D}^-\} \quad (97)$$

772 respectively, where  $\text{RP}_{i-1/2}$  stands for the Riemann Problem at left interface  
 773 and  $\text{RP}_{i+1/2}$  at right interface. The whole set of RPs is given by

$$\mathcal{R} = \mathcal{R}_1 \cup \mathcal{R}_2. \quad (98)$$

774 By using the previous definitions, the approximation of the integral of the  
 775 source term at any interface is defined as follows

$$\bar{S}_{i+1/2} = \begin{cases} -g\bar{h}\delta z + (\bar{h}u - \bar{h}\bar{u})\delta u & \text{if } \text{RP}_{i+1/2} \notin \mathcal{R} \\ -g\bar{h}\delta z & \text{if } \text{RP}_{i+1/2} \in \mathcal{R} \end{cases} \quad (99)$$

776 and the method will be hereafter referred to as selective energy balanced  
 777 formulation (SEBF).

#### 778 4.4. The ARoe scheme for the SWE

779 When applied to the Shallow Water Equations, the Augmented Roe solver  
 780 provides a linearized solution that can be straightforward expanded from the  
 781 homogeneous case. The approximate Jacobian  $\tilde{\mathbf{J}}$  of the homogeneous part is  
 782 given by [8]

$$\tilde{\mathbf{J}}_{i+1/2} = \begin{pmatrix} 0 & 1 \\ \tilde{c}^2 - \tilde{u}^2 & 2\tilde{u} \end{pmatrix}_{i+1/2}, \quad \delta \mathbf{F}_{i+1/2} = \tilde{\mathbf{J}}_{i+1/2} \delta \mathbf{U}_{i+1/2}, \quad (100)$$

783 where

$$\begin{aligned} \tilde{\lambda}^1 &= \tilde{u} - \tilde{c}, & \tilde{\lambda}^2 &= \tilde{u} + \tilde{c} \\ \tilde{\mathbf{e}}^1 &= \begin{pmatrix} 1 \\ \tilde{u} - \tilde{c} \end{pmatrix}, & \tilde{\mathbf{e}}^2 &= \begin{pmatrix} 1 \\ \tilde{u} + \tilde{c} \end{pmatrix} \end{aligned} \quad (101)$$

784 with

$$\tilde{c} = \sqrt{g \frac{h_i + h_{i+1}}{2}}, \quad \tilde{u} = \frac{u_{i+1} \sqrt{h_{i+1}} + u_i \sqrt{h_i}}{\sqrt{h_{i+1}} + \sqrt{h_i}}. \quad (102)$$

785 *4.5. Test case 1: steady shock capturing for the SWE with bed topography*

786 In this test case, steady solutions for the flow over the following bed  
787 elevation profile

$$z(x) = \begin{cases} 0 & \text{if } x < 8 \\ 0.05(x - 8) & \text{if } 8 \leq x \leq 12 \\ 0.2 - 0.05(x - 12)^2 & \text{if } 12 \leq x \leq 14 \\ 0 & \text{if } x > 14 \end{cases} \quad (103)$$

788 are computed using the ARoe solver in combination with the different dis-  
789 cretization techniques for the source term outlined before. The computa-  
790 tional domain is  $[0, 20]$  and the solution is computed for  $t = 600$  s. CFL  
791 number is set to 0.45 for all cases. The discharge is imposed to  $0.6 \text{ m}^2/\text{s}$   
792 upstream to obtain the critical point at the cell with maximum bed eleva-  
793 tion, that is  $z_{max} = 0.2$ . Downstream, the water depth is also imposed to  
794  $h = 0.621$  m in order to generate a hydraulic jump downstream. Different  
795 computational grids, composed of 100, 200, 400, 800 and 1600 cells respec-  
796 tively, are used to compute the numerical solution.

797 Numerical solutions provided by the ARoe solver when using the different  
798 approximations of the source term presented before, namely the differential  
799 formulation (DF), the integral formulation (IF), the weighted energy bal-  
800 anced formulation (WEBF) and the novel selective energy balanced method  
801 (SEBF), are presented and compared with the exact solution in Figures 4, 5.  
802 In Figure 4, the numerical solutions for  $h + z$  and  $q$  computed by the ARoe  
803 solver in combination with all the previous techniques on two grids of 100  
804 and 400 cells are plotted together and compared with the exact solution. To  
805 study the effect of mesh refinement in the accuracy of the numerical solution  
806 and convergence to the exact position of the shock, a detailed plot of the so-  
807 lution provided by each one of the methods is presented in Figure 5 for three  
808 different grids composed of 200, 400 and 800 cells respectively. Numerical  
809 results evidence that those approximations based on the integral discretiza-  
810 tion of the source term, such as the energy balanced approach from [25] and  
811 the integral discretization itself, do not accurately capture the position of  
812 the shock, with independence of the grid. In any case, the former strategy  
813 provides much better results than the latter, as it is energy-conservative. On  
814 the other hand, it is evidenced that both the differential formulation and  
815 the selective energy balanced formulation do accurately capture the shock  
816 position for any grid.

817 It is also noticed that a spurious spike in the numerical discharge appears  
 818 for all methods and what is of utmost relevance, that the amplitude of this  
 819 spike is not reduced with mesh refinement, as observed in Figure 4.

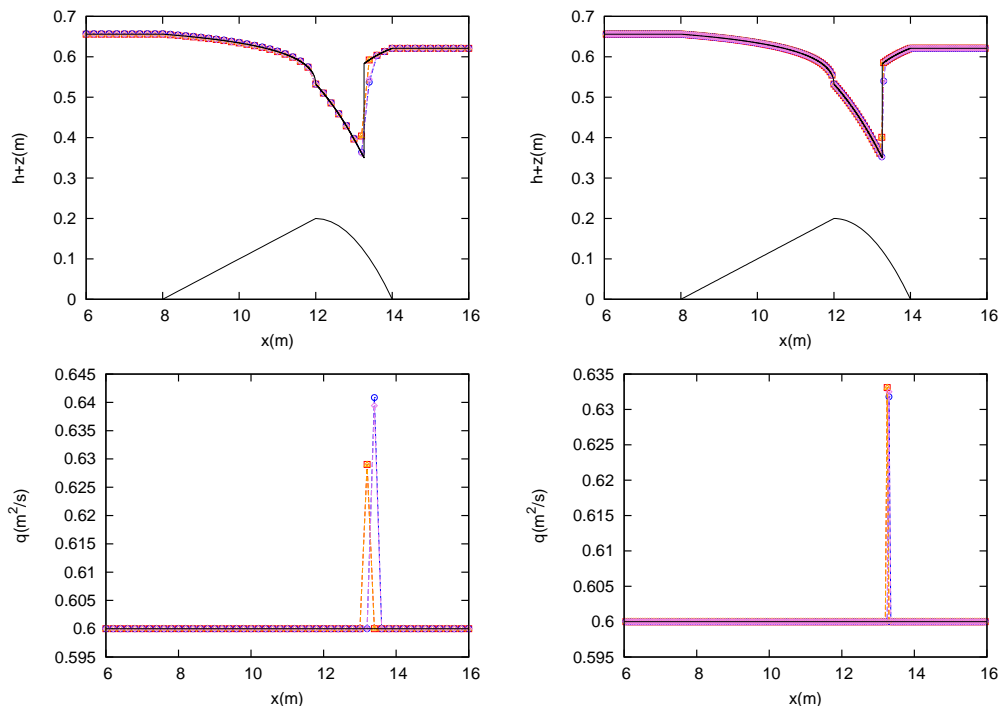


Figure 4: Test case 1. Exact (—) and numerical solution for  $h + z$  (top) and  $q$  (bottom) computed by the ARoe solver in combination with the DF (— $\triangle$ —), IF (— $\circ$ —), SEBF (— $\square$ —) and WEBF (— $\diamond$ —), using 100 (left) and 400 cells (right).

820 The numerical solution for the specific mechanical energy, computed using  
 821 the aforementioned techniques in the grids of 100 and 400 cells, is presented  
 822 in Figure 6 left and right respectively. It is observed that only when using an  
 823 energy-balanced STD (E-scheme), such as the ARoe solver in combination  
 824 with the SEBF or WEBF formulations, energy is conserved. On the other  
 825 hand, when using the DF and IF formulations of the source term, energy  
 826 is not conserved though it converges with mesh refinement. Among the  
 827 assessed methods, the SEBF is the one providing the best performance, as  
 828 it ensures the conservation of energy when required and accurately captures  
 829 the position of the hydraulic jump. This method provides the exact solutions  
 830 in all cells but the one containing the shock, with independence of the grid.

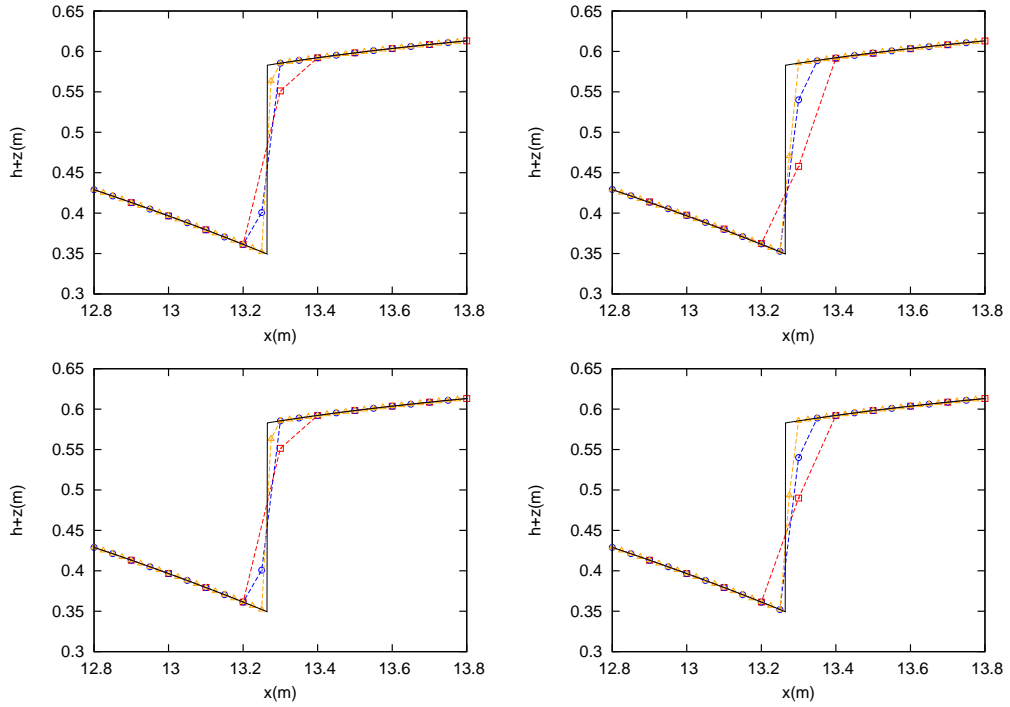


Figure 5: Test case 1. Exact (—) and numerical solution for  $h + z$  computed by the ARoe solver in combination with the DF (top left), IF (top right), SEBF (bottom left) and WEBF (bottom right) using 200 (—□—), 400 (—○—) and 800 (—△—) cells.

831 **5. Numerical shockwave anomalies in the SWE: computation of**  
 832 **the hydraulic jump**

833 It has been widely reported in the literature that significant numerical  
 834 anomalies arise in presence of shock waves. An example of such problems are  
 835 the Carbuncle, the slowly-moving shock and the wall-heating phenomenon,  
 836 all of them leading to spurious numerical solutions. The aforementioned  
 837 problems have been deeply studied in the framework of Euler equations and  
 838 some authors have proposed different numerical techniques to address them.  
 839 Here, we will focus on the numerical anomalies present when computing  
 840 steady and moving hydraulic jumps, which are a particular type of shock  
 841 waves in the framework of the Shallow Water Equations (SWE). Specifically,  
 842 our interest lies in the reduction of the spike in the discharge, reported in  
 843 the previous section.

844 The hydraulic jump occurs when a supercritical flow suddenly changes to

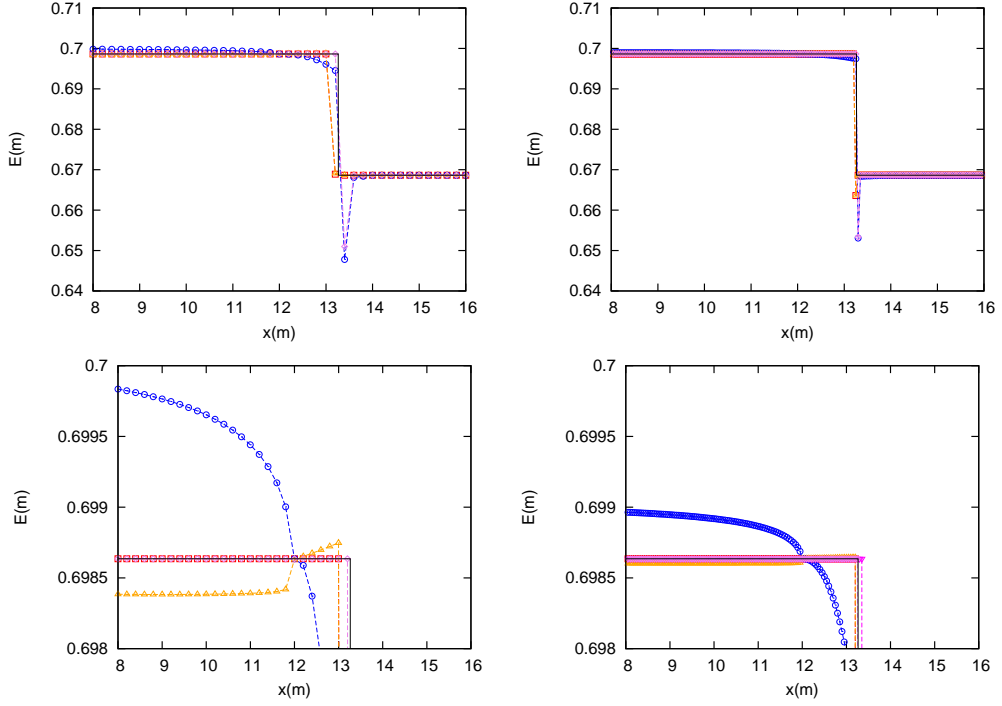


Figure 6: Test case 1. Numerical solution for the specific mechanical energy computed by the ARoe solver in combination with the DF ( $-\triangle-$ ), IF ( $-\circ-$ ), SEBF ( $-\square-$ ) and WEBF ( $-\diamond-$ ) (top) and detail of the solution (bottom), using 100 (left) and 400 (right) cells.

845 subcritical conditions, generating a steep free surface elevation where intense  
 846 mixing takes place and a large amount of mechanical energy is dissipated.  
 847 Mathematically, hydraulic jumps are modelled by a discontinuity correspond-  
 848 ing to a shock wave and the relation between the states at each side of the  
 849 discontinuity is provided by the RH conditions.

### 850 5.1. Hugoniot locus of the hydraulic jump

851 To understand the mathematical treatment of the hydraulic jump and  
 852 the numerical anomalies arising from such a wave, it is worth studying first  
 853 the analytical solution of this type of wave under the simplest conditions,  
 854 that is over flat bed. From Rankine-Hugoniot (RH) conditions, all possible  
 855 values connecting the left and right states can be determined and represented  
 856 in phase space as  $(h(\xi), hu(\xi))$  by means of the so-called Hugoniot locus

$$\mathbf{U}(\xi) = \begin{pmatrix} h(\xi) \\ hu(\xi) \end{pmatrix} = \begin{pmatrix} h_L + \xi \\ (hu)_L + \xi \left( u_L \pm \sqrt{gh_L + \frac{1}{2}g\xi \left( 3 + \frac{\xi}{h_L} \right)} \right) \end{pmatrix}, \quad (104)$$

857 where  $\xi = h - h_L$ , with  $h$  the independent variable used for the parametriza-  
858 tion. From (104), we notice that two families of curves are possible, denoted  
859 by  $\Psi^1$  and  $\Psi^2$ , which are associated to the 1-wave and 2-wave respectively.  
860 Such curves are defined by

$$\Psi^1(\xi) = \begin{pmatrix} \psi_1^1(\xi) \\ \psi_2^1(\xi) \end{pmatrix} = \begin{pmatrix} h_L + \xi \\ (hu)_L + \xi \left( u_L - \sqrt{gh_L + \frac{1}{2}g\xi \left( 3 + \frac{\xi}{h_L} \right)} \right) \end{pmatrix}, \quad (105)$$

$$\Psi^2(\xi) = \begin{pmatrix} \psi_1^2(\xi) \\ \psi_2^2(\xi) \end{pmatrix} = \begin{pmatrix} h_L + \xi \\ (hu)_L + \xi \left( u_L + \sqrt{gh_L + \frac{1}{2}g\xi \left( 3 + \frac{\xi}{h_L} \right)} \right) \end{pmatrix}. \quad (106)$$

861 Figure 7 depicts different curves obtained for different left-reference states  
862 using (105) in red and (106) in blue, for  $\Psi^1, \Psi^2 \in \mathbb{R}^+ \times \mathbb{R}^+$ . Also the curve  
863  $hu(h) = \sqrt{gh^3}$  that represents the transition between supercritical (white  
864 background) and subcritical region (green background) is depicted in the  
865 figure. For any given set of two points laying on a curve, a weak solution of  
866 the PDEs in the form of a shock wave is mathematically possible. It is worth  
867 pointing out that further representations of the aforementioned curves will  
868 be carried out by the parametrization of  $\psi_2^m$ , which is the discharge  $hu$ , in  
869 terms of  $\psi_1^m$ , which is  $h$ , so that their representation in the phase space  $h, hu$   
870 is straightforward.

871 It must be borne in mind that not every choice of subcritical state that  
872 is connected to a given supercritical state represents a hydraulic jump. For  
873 instance, let us consider a left supercritical state given by  $h_L = 0.85$  and  
874  $hu_L = 3.411764705882353$  and let us find two possible right states connected  
875 to it, each of them laying on each branch of the Hugoniot locus. This is  
876 depicted in Figure 8, where the original left state is denoted by F, the right  
877 state lying on the 1-curve,  $\Psi^1$ , is denoted by G and the right state lying on

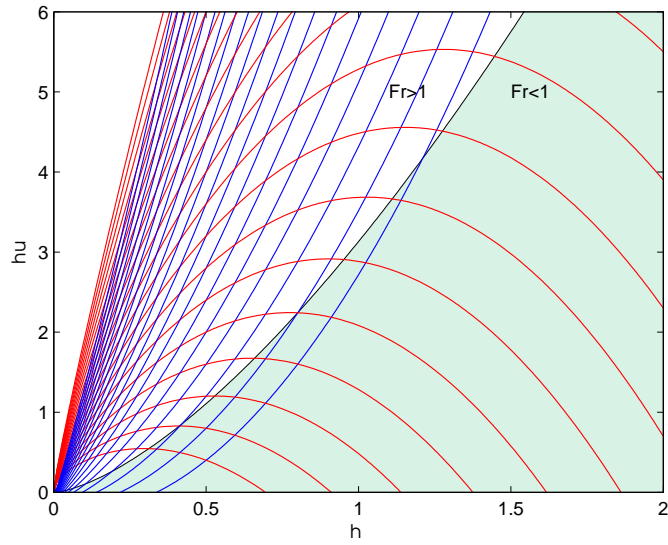


Figure 7: Phase space  $(h, hu) \in \mathbb{R}^+ \times \mathbb{R}^+$  with the subcritical region depicted in green background and the supercritical region in white background, showing the Hugoniot locus  $\Psi^1$  in red and  $\Psi^2$  in blue, obtained for different left-reference states using (105) and (106) respectively.

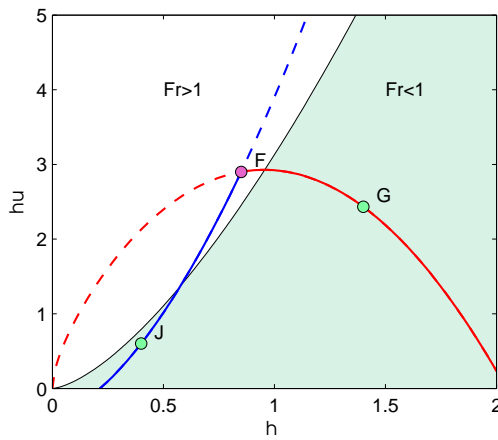


Figure 8: Phase space  $(h, hu) \in \mathbb{R}^+ \times \mathbb{R}^+$  with the subcritical region depicted in green background and the supercritical region in white background, showing the Hugoniot locus  $\Psi^1$  in red and  $\Psi^2$  in blue.

878 the 2-curve,  $\Psi^2$ , is denoted by J. The entropically inadmissible region of the  
 879 curves has been represented by dashed line. It is observed that both G and J

880 lie on the subcritical region of the phase plane and they are both entropically  
 881 admissible, however, only the combination of states F–G leads to a hydraulic  
 882 jump, because G, unlike J, has a higher water depth than F and, what is  
 883 decisive in this case, wave celerities of F and G have opposite sign. More  
 884 generally, we can define an hydraulic jump as:

885 **Definition 3.** (*Hydraulic jump*). *Let the following discontinuous solution*

$$\mathbf{U}(x, t) = \begin{cases} (h, hu)_L & x < 0 \\ (h, hu)_R & x > 0 \end{cases} \quad (107)$$

886 *be a weak solution of the SWE system, where  $(h, hu)_L$  and  $(h, hu)_R$  are two*  
 887 *different states laying on  $\Psi^m$  and satisfying the entropy condition  $\lambda^m(\mathbf{U}_L) >$*   
 888  *$\mathcal{S}^m > \lambda^m(\mathbf{U}_R)$ , with  $\mathcal{S}^m$  the speed of the jump, that undergoes a flow transi-*  
 889 *tion as  $Fr_L < 1 < Fr_R$  or  $Fr_R < 1 < Fr_L$ . Solution in (107) is termed as*  
 890 *hydraulic jump if and only if  $\lambda^m(\mathbf{U}_L) > 0 > \lambda^m(\mathbf{U}_R)$ .*

891 Notice that, according to the previous definition, hydraulic jumps admit  
 892 that  $\mathcal{S}^m$  be nil, hence they are the only shock-type solution for the SWE that  
 893 can be stationary at a fixed position.

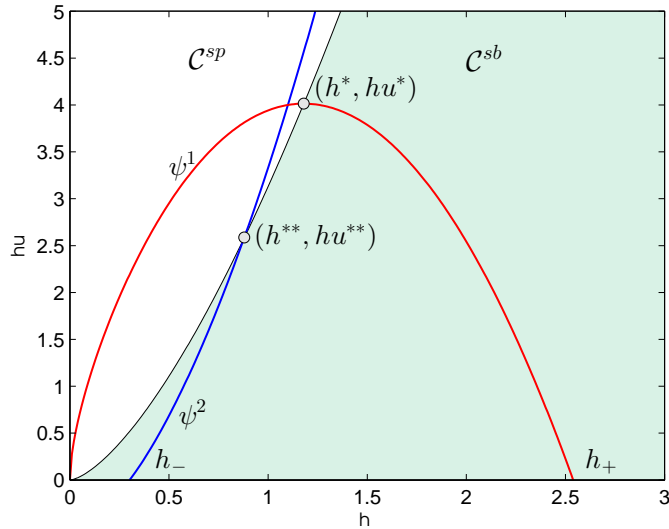


Figure 9: Phase space  $(h, hu) \in \mathbb{R}^+ \times \mathbb{R}^+$  with the subcritical region,  $\mathcal{C}^{sb}$ , depicted in green and the supercritical,  $\mathcal{C}^{sp}$ , region in white, showing the Hugoniot locus  $\Psi^1$  in red and  $\Psi^2$  in blue and the corresponding intersection.



894 From the analysis of the Hugoniot locus considering  $h, hu > 0$  and departing  
 895 from a left reference point located in the supercritical region, we notice the  
 896 following points:

- 897 • Curve  $\sqrt{gh^3}$  is monotonically increasing and divides the space  $\mathbb{R}^+ \times \mathbb{R}^+$   
 898 in two sets,  $\mathcal{C}^{sp}$  and  $\mathcal{C}^{sb}$ , as follows

$$\mathcal{C}^{sp} = \left\{ (h, hu) \in \mathbb{R}^2 \mid hu > \sqrt{gh^3} \wedge h > 0 \right\}, \quad (108)$$

$$\mathcal{C}^{sb} = \left\{ (h, hu) \in \mathbb{R}^2 \mid hu < \sqrt{gh^3} \wedge h > 0 \right\}, \quad (109)$$

899 such that  $\mathcal{C}^{sp} \cup \mathcal{C}^{sb} \cup \mathcal{C}^{cr} = \mathbb{R}^+ \times \mathbb{R}^+$ , where

$$\mathcal{C}^{cr} = \left\{ (h, hu) \in \mathbb{R}^2 \mid hu = \sqrt{gh^3} \wedge h > 0 \right\}. \quad (110)$$

- 900 • Curve  $\sqrt{gh^3}$  is monotonically increasing.
- 901 • Curve  $\psi_2^1$  has a global maximum at  $h_{max}$  such that  $(h_{max}, hu_{max}) \in$   
 902  $\mathbb{R}^+ \times \mathbb{R}^+$ .
- 903 • Curve  $\psi_2^2$  is monotonically increasing in  $\mathbb{R}^+ \times \mathbb{R}^+$ .
- 904 • Curves  $\sqrt{gh^3}$  and  $\psi_2^1$  intersect at a single point denoted by  $(h^*, hu^*) \in$   
 905  $\mathbb{R}^+ \times \mathbb{R}^+$ , with  $hu^* < hu_{max}$ .
- 906 • Curves  $\sqrt{gh^3}$  and  $\psi_2^2$  intersect at a single point denoted by  $(h^{**}, hu^{**}) \in$   
 907  $\mathbb{R}^+ \times \mathbb{R}^+$ .
- 908 • We can define two sets of  $h$  states,  $\mathcal{H}^{sp,1} = (0, h^*)$  and  $\mathcal{H}^{sb,1} = (h^*, h_+)$ ,  
 909 with  $h_+$  the value of  $h$  for which  $\Psi^1 = (h_+, 0)$ , such that  $\Psi^1 \in \mathcal{C}^{sp} \forall h \in$   
 910  $\mathcal{H}^{sp,1}$  and  $\Psi^1 \in \mathcal{C}^{sb} \forall h \in \mathcal{H}^{sb,1}$ .
- 911 • We can define two set of  $h$  states,  $\mathcal{H}^{sp,2} = (h_-, h^{**})$  and  $\mathcal{H}^{sb,2} =$   
 912  $(h^{**}, \infty)$ , with  $h_-$  the value of  $h$  for which  $\Psi^2 = (h_-, 0)$ , such that  
 913  $\Psi^2 \in \mathcal{C}^{sb} \forall h \in \mathcal{H}^{sp,2}$  and  $\Psi^2 \in \mathcal{C}^{sp} \forall h \in \mathcal{H}^{sb,2}$ .

914 Definitions introduced in the previous statements are depicted in the top-left  
 915 plot in Figure 10. From the previous points, the following observations are  
 916 worth being mentioned:

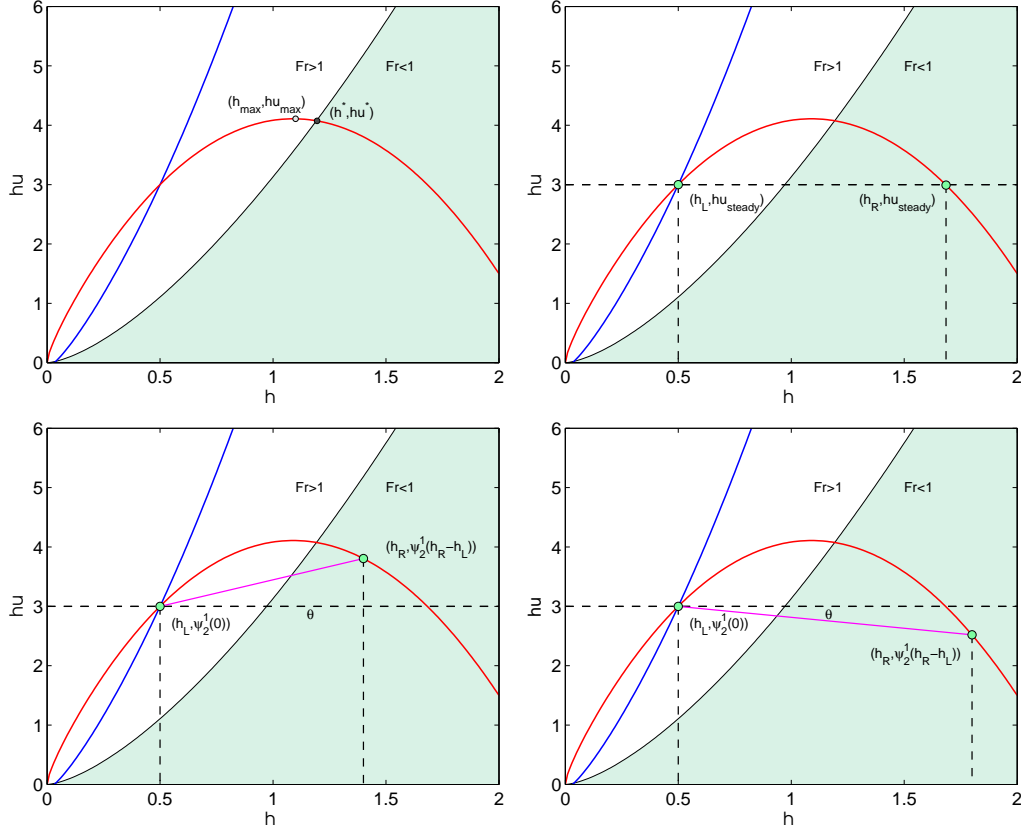


Figure 10: Hugoniot locus  $\Psi^1$  in red and  $\Psi^2$  in blue for the left state  $(h, hu) = (0.5, 3)$ , showing three possible solutions in the form of a hydraulic jump: a steady jump (top-right), a right-moving jump (bottom-left) and a left-moving jump (bottom-right).

- 917 • According to the two last points stated before, hydraulic jumps with  
918  $hu > 0$  only take place when  $\Psi^m \in \mathcal{C}^{sp} \forall h \in \mathcal{H}^{sp,m}$  and  $\Psi^m \in$   
919  $\mathcal{C}^{sb} \forall h \in \mathcal{H}^{sb,m}$ , which is only possible for  $\Psi^1$ . Hence, any solution  
920 for  $\text{RP}(\mathbf{U}_L, \mathbf{U}_R)$ , with  $\mathbf{U}_L = \Psi^1(0)$  and  $\mathbf{U}_R = \Psi^1(h - h_L) \forall h \in$   
921  $\mathcal{H}^{sb,1}, \forall h_L \in \mathcal{H}^{sp,1}$ , evolves as a hydraulic jump.
- 922 • There exist two points  $h_L \in \mathcal{H}^{sp,1}$  and  $h_R \in \mathcal{H}^{sb,1}$  such that  $\psi_2^1(0) =$   
923  $\psi_2^1(h_R - h_L) \equiv (hu)_{steady}$  and  $\psi_2^1(0), \psi_2^1(h_R - h_L) \in (0, hu^*) \subset \mathbb{R}^+$ .  
924 Such points correspond to the left and right states of the hydraulic  
925 jump under steady conditions with a constant discharge of  $(hu)_{steady}$ .  
926 This case is depicted in Figure 10 (top-right plot)

- 927 • There exist two other points  $h_L \in \mathcal{H}^{sp,1}$  and  $h_R \in \mathcal{H}^{sb,1}$  such that  
 928  $\psi_2^1(0) \in (0, hu_{max}) \subset \mathbb{R}^+$  and  $\psi_2^1(h_R - h_L) \in (0, hu^*) \subset \mathbb{R}^+$ . If  $\psi_2^1(0) <$   
 929  $\psi_2^1(h_R - h_L)$  a right-moving transient shock will appear as depicted  
 930 in Figure 10 (bottom-left plot). If  $\psi_2^1(0) > \psi_2^1(h_R - h_L)$ , a left-moving  
 931 transient shock will appear as depicted in Figure 10 (bottom-right plot).
- 932 • Shock speed is equal to the slope of the magenta straight line in Figure  
 933 10, that is  $\mathcal{S} = \tan \theta$ .
- 934 • The previous statements apply to  $\psi_2^2$  in the region  $\mathbb{R}^+ \times \mathbb{R}^-$  when  
 935 considering  $hu < 0$ .

936 *5.2. Analytical study and comparison of the exact solution for 2 and 3-states*  
 937 *hydraulic jumps.*

938 Prior to analyzing the numerical solutions of Godunov's scheme to the  
 939 hydraulic jump, it is worth studying the analytical solutions to this problem,  
 940 which will help to understand the nature and characteristics of the numerical  
 941 (discrete) solution to it. It is well known that an intermediate state appears in  
 942 the numerical solution provided by Godunov's scheme, with independence of  
 943 the solver [42]. The presence of this intermediate state, hereafter denoted by  
 944  $\mathbf{U}_M$ , is not of any physical relevance as it provides an unrealistic estimation  
 945 of the average discharge in the intermediate cell (spike) which does not match  
 946 the constant value of discharge. However, when using conservative schemes  
 947 the intermediate value may be useful to compute a rough estimate of the  
 948 shock position. The position of the shock inside the cell can be computed  
 949 imposing conservation of mass as

$$x_S = \frac{h_M - h_R}{h_L - h_R}, \quad (111)$$

950 where  $x_S \in [0, 1]$  represents the normalized position of the shock (where  
 951  $0 \equiv$ left interface,  $0.5 \equiv$ middle position and  $1 \equiv$ right interface) [42].

952 As a first approach and before getting into the numerical issues concern-  
 953 ing hydraulic jumps, let us compare analytically the solution for the ideal  
 954 steady hydraulic jump (pure discontinuity) with another solution for the  
 955 steady hydraulic jump that includes an intermediate state, which resembles  
 956 the discrete solution provided by Godunov's scheme. Both solutions are weak  
 957 solutions of the equations and they are both valid. Whereas the former is  
 958 characterized by two states, namely  $\mathbf{U}_L$  and  $\mathbf{U}_R$ , the latter is given by  $\mathbf{U}_L$ ,

959  $\mathbf{U}_M$  and  $\mathbf{U}_R$ . Moreover, the latter does not experience a sudden transition of  
 960 flow regime, hence it cannot be considered a pure, or ideal, hydraulic jump.

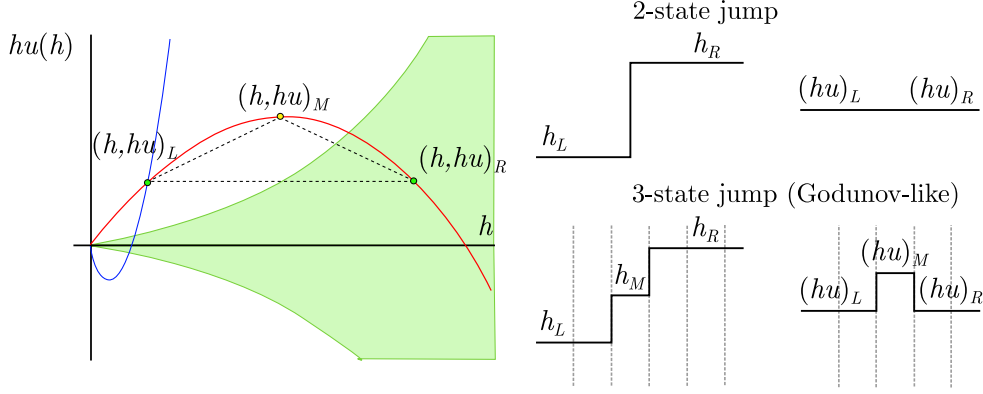


Figure 11: Hugoniot Locus and sketch of the analytical solutions for a 2-state and 3-state hydraulic jumps.

961 Let us consider first the ideal hydraulic jump composed of two states. This  
 962 solution consists of a supercritical right-moving steady flow that suddenly  
 963 decelerates through a pure discontinuity to subcritical conditions, as depicted  
 964 schematically in Figure 11 (top-right). The Hugoniot locus that connects the  
 965 left and right states of the jump,  $\Psi^1$ , is depicted in Figure 11 (left), showing  
 966 that such states are located at the intersection of the Hugoniot Locus with  
 967 the curve of constant discharge  $(hu)_L = (hu)_R$ , ensuring the steady regime.

968 On the other hand, when seeking a weak solution of the equations that  
 969 includes an intermediate state,  $\mathbf{U}_M$ , as depicted in Figure 11 (bottom-right),  
 970 we need to look for this additional state on the Hugoniot curve. According to  
 971 Figure 11 (left), the intermediate state  $(h_M, (hu)_M)$  (yellow point) will lie on  
 972 Hugoniot Locus and is connected to the left and right states (green points)  
 973 through this curve. From the previous observations, we realize that only a  
 974 linear Hugoniot Locus would ensure a constant discharge in the intermediate  
 975 state [42].

976 If a curve of the family of

$$\check{\Psi}(\xi) = \begin{pmatrix} h(\xi) \\ (hu)_{steady} \end{pmatrix} \quad (112)$$

977 was considered in state space, with  $(hu)_{steady} \in \mathbb{R}^+$  for a right-moving flow,  
 978 a constant discharge for the intermediate state would be possible. Only if

979  $\Psi^1$  was of the type of  $\check{\Psi}$ , constant discharge would be ensured across the  
 980 intermediate cell. This means that we would have a linear Hugoniot [42].  
 981 This concept can be extended to moving hydraulic jumps by examination  
 982 of Figure 10 (bottom left). Let us redefine the states denoted in the plot  
 983 by  $(h', hu')$  and  $(h'', hu'')$  as left state  $(h_L, hu_L)$  and right state  $(h_R, hu_R)$ ,  
 984 respectively. The linear Hugoniot must lie on the line depicted in magenta,  
 985 with slope  $\theta = (h_R - h_L)/(hu_R - hu_L)$  and can be parametrized in terms of  
 986  $x_S$  in (111). Hence, it can be expressed as

$$\check{\Psi}(x_S) = \begin{pmatrix} h(x_S) \\ hu(x_S) \end{pmatrix}, \quad (113)$$

987 where  $h(x_S) = x_S(h_R - h_L) + h_L$ ,

$$hu(x_S) = hu_L + \theta h(x_S) \quad (114)$$

988 and  $x_S \in [0, 1]$ . Note that parametrization  $\check{\Psi}(\xi)$  is straightforward as  $\xi =$   
 989  $(h_R - h_L)x_S$ .

990 Considering again the steady case described above and depicted in Figure  
 991 11, we can observe that the exact Hugoniot is neither linear nor monotone  
 992 and  $\psi_2^1$  has a global maxima  $hu_{max}$  at  $h_{max} \in [h_L, h_R] \subset \mathbb{R}^+$  therefore,  
 993 for any  $h_M \in [h_L, h_R] \subset \mathbb{R}^+$ , we have that  $(hu)_M \geq (hu)_L = (hu)_R \equiv$   
 994  $(hu)_{steady}$ . This can be observed in Figure 11 (bottom-right), where a spike  
 995 in the discharge appears.

### 996 5.3. Properties of the intermediate state in discrete Godunov-type solutions

997 Up to this point throughout this section, we have only considered exact  
 998 solutions to the hydraulic jump. Theoretically, when considering the exact  
 999 solution, the presence of an intermediate constant state  $\mathbf{U}_M = (h_M, (hu)_M)$   
 1000 is not stable, that is, it cannot be kept under steady conditions. The reason  
 1001 for this is that both jumps (left to middle and middle to right) have non-zero  
 1002 wave velocities of opposite sign, hence both jumps would converge to form  
 1003 a unique jump. This behavior, shown in Figure 12, is only present in the  
 1004 exact solution. On the other hand, when considering a discrete solution in  
 1005 a computational grid, both waves could be kept at a stationary position (at  
 1006 the cell interfaces of the intermediate cell) and the intermediate cell could  
 1007 keep the intermediate value in the steady regime. The reason for this is that  
 1008 the numerical fluxes at the interfaces of such a cell would coincide, that is

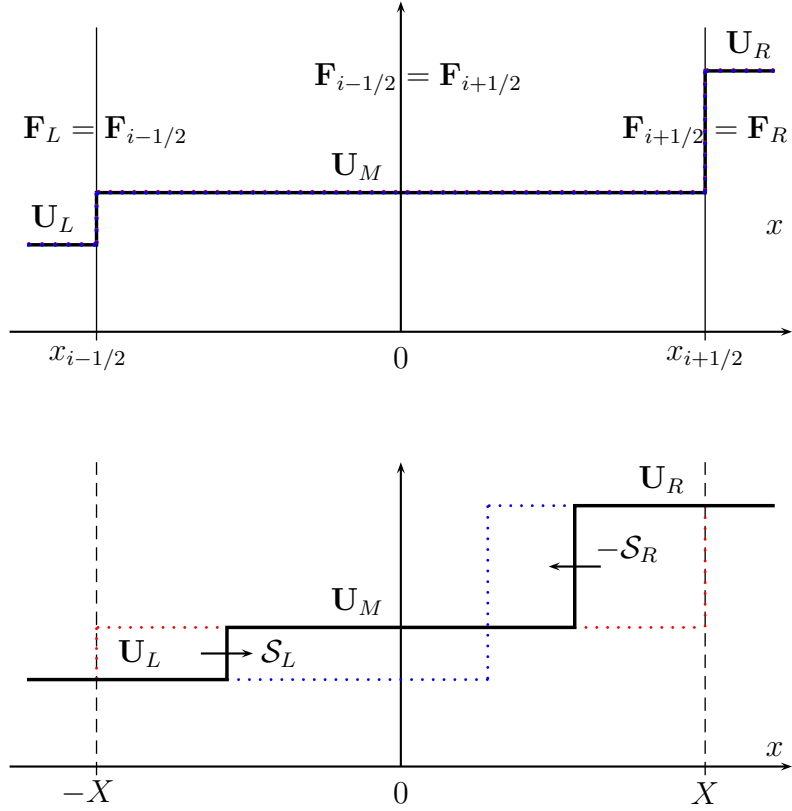


Figure 12: Initial condition considering an intermediate state (red), transient evolution of the discontinuities  $U_L-U_M$  and  $U_M-U_R$  (black) and final steady solution (blue).

$$\mathbf{F}_{i+1/2}^- = \mathbf{F}_{i-1/2}^+, \quad (115)$$

1009 when considering the numerical resolution of the problem by means of FV  
 1010 Godunov's scheme in (44).

1011 Figure 12 depicts the contrasting behavior of the 3-state hydraulic jump  
 1012 when considering the discrete (top) and exact (bottom) solution. The initial  
 1013 condition is represented by red dotted line, the final solution (when steady-  
 1014 ness is achieved) is represented by blue dotted line and the solution at an  
 1015 arbitrary time before reaching the steady state is represented by black solid  
 1016 line. It can be observed that the initial condition is maintained in the dis-  
 1017 crete solution, where the intermediate state,  $U_M$ , has been defined inside the

1018 cell  $[x_{i-1/2}, x_{i+1/2}]$ .

1019 There is another important issue worth being mentioned. Only when the  
1020 intermediate state coincides with the left or right states, the approximate  
1021 solver would provide the exact solution. Hence, only when the shock position  
1022 is located exactly at the interface, the approximate solver provides the exact  
1023 solution [53, 54]. Moreover, it must be borne in mind that the intermediate  
1024 state,  $\mathbf{U}_M$ , does depend on the Riemann solver used for the computation of  
1025 the fluxes, and will only coincide with the value of  $\mathbf{U}_M$  provided by the ana-  
1026 lytical Hugoniot locus when using an exact solver. A exhaustive comparison  
1027 of the numerical performance in shock-capturing of different flux functions  
1028 in the framework of Euler equations can be found in [55].

## 1029 6. Flux fixes for the computation of the hydraulic jump

1030 In this section, some spike-reduction numerical techniques based on flux  
1031 interpolation are recalled and applied to the Shallow Water Equations (SWE).  
1032 This idea of flux interpolation was first presented by Zaide and Roe [42],  
1033 who proposed to find the fluxes in the untrustworthy intermediate cells by  
1034 extrapolation from trustworthy neighbors and presented two new flux func-  
1035 tions. The first one, named by the authors flux function A, was constructed  
1036 based on the flux-wave approach, by computing the fluctuations in the inter-  
1037 polated fluxes across each wave. The second one, called flux function B, is  
1038 based on the classical Roe solver and relies on conserved variables to deter-  
1039 mine the jumps across each wave and the contribution of each wave to the  
1040 numerical flux. The authors claim that, by enforcing a linear shock structure  
1041 and unambiguous sub-cell shock position, numerical shockwave anomalies  
1042 are dramatically reduced.

1043 Zaide and Roe [42] proposed to compute the fluxes in the intermediate  
1044 cells by extrapolation from neighboring cells, hence a more general idea of  
1045 a homogeneous flux function of the type  $\mathbf{F}_{i+1/2}^* = \mathbf{F}_{i+1/2}^*(\mathbf{U}_{i-m}, \dots, \mathbf{U}_{i-n})$   
1046 was introduced, rather than a Riemann solver that computes the numerical  
1047 flux as  $\mathbf{F}_{i+1/2}^* = \mathbf{F}_{i+1/2}^*(\mathbf{U}_i, \mathbf{U}_{i+1})$ , with  $m$  and  $n$  two integer numbers. The  
1048 authors in [42] outline that the conserved variables must be trusted since this  
1049 is the only way to ensure conservation, however, the flux values should not  
1050 be trusted.

1051 Prior to the construction of the novel numerical fluxes  $\mathbf{F}_{i+1/2}^*$ , physical  
1052 fluxes (which are the cell centered fluxes,  $\mathbf{F}_i$ ) are used to construct a novel

1053 approximation of the fluxes in every cell. Cell-centered fluxes,  $\mathbf{F}_i$ , are re-  
 1054 computed by means of extrapolation from neighboring cells. At every cell,  
 1055 the new flux is calculated as

$$\check{\mathbf{F}}_i = \frac{1}{2}(\mathbf{F}_{i+1} + \mathbf{F}_{i-1}) - \frac{1}{2}\tilde{\mathbf{J}}_{i-1,i+1}(\mathbf{U}_{i+1} - 2\mathbf{U}_i + \mathbf{U}_{i-1}), \quad (116)$$

1056 with  $\tilde{\mathbf{J}}_{i-1,i+1} = \tilde{\mathbf{J}}_{i-1,i+1}(\mathbf{U}_{i+1}, \mathbf{U}_{i-1})$  a Jacobian Roe's matrix,

$$\mathbf{F}_{i+1} - \mathbf{F}_{i-1} = \tilde{\mathbf{J}}_{i-1,i+1}(\mathbf{U}_{i+1} - \mathbf{U}_{i-1}). \quad (117)$$

1057 To construct those more general numerical fluxes, two alternatives, named  
 1058 flux function A and flux function B, are proposed in [42]. Such alternatives,  
 1059 as well as the traditional Roe flux, are detailed below:

- 1060 • Traditional Roe homogeneous flux:

1061 The traditional Roe homogeneous flux (B.8) in Appendix B is used. It  
 1062 is constructed using Roe's matrix  $\tilde{\mathbf{J}}_{i+\frac{1}{2}}$ ,

$$\mathbf{F}_{i+1/2}^{*,Roe} = \frac{1}{2}(\mathbf{F}_i + \mathbf{F}_{i+1}) - \frac{1}{2}|\tilde{\mathbf{J}}_{i+1/2}| \delta\mathbf{U}_{i+1/2}, \quad (118)$$

1063 evaluated conventionally as  $\tilde{\mathbf{J}}_{i+\frac{1}{2}} = \tilde{\mathbf{J}}_{i+\frac{1}{2}}(\mathbf{U}_i, \mathbf{U}_{i+1})$ .

- 1064 • Flux function A:

1065 The extrapolated fluxes,  $\check{\mathbf{F}}_i$ , computed by (116), can be directly pro-  
 1066 jected onto the Jacobian's eigenvectors basis and upwinded according  
 1067 to the propagation velocities of the Jacobian. The resulting numerical  
 1068 flux is constructed using (B.8), yielding [42]

$$\mathbf{F}_{i+1/2}^{*,A} = \frac{1}{2}(\check{\mathbf{F}}_i + \check{\mathbf{F}}_{i+1}) - \frac{1}{2}\text{sgn}\left(\tilde{\mathbf{J}}_{i+\frac{1}{2}}\right) \delta\check{\mathbf{F}}_{i+1/2}. \quad (119)$$

- 1069 • Flux function B:

1070 This new flux function is computed by means of a novel Roe's matrix  
 1071 that spans a wider set of cells, instead of just the two cells at each side  
 1072 of the discontinuity. It reads [42]

$$\mathbf{F}_{i+1/2}^{*,B} = \frac{1}{2}(\check{\mathbf{F}}_i + \check{\mathbf{F}}_{i+1}) - \frac{1}{2}|\bar{\mathbf{J}}_{i+1/2}| \delta\mathbf{U}_{i+1/2}, \quad (120)$$



1073 with  $\bar{\mathbf{J}}_{i+1/2} = \bar{\mathbf{J}}_{i+1/2}(\mathbf{U}_{i-1}, \mathbf{U}_{i+2})$  Roe's matrix computed with cells  
1074  $i - 1$  and  $i + 2$ .

### 1075 6.1. Test case 2: assessment of flux functions A and B for the SWE

1076 In order to test flux functions A and B in the framework of the SWE  
1077 and compare their performance with the traditional homogeneous Roe flux,  
1078 the following numerical experiment is proposed. It consists of a RP with  
1079 initial data  $h_L = 0.5$ ,  $(hu)_L = 3$ ,  $h_R = 1.6$  and  $(hu)_R = 3.28787832816$ , that  
1080 generates a moving shock wave with speed  $\mathcal{S} = 0.26171$ . The computational  
1081 domain is set to  $[0, 450]$ , with the discontinuity located at  $x = 225$ . Regarding  
1082 the numerical discretization, the computational domain is divided in 900 cells  
1083 of size  $\Delta x = 0.5$  and the CFL number is set to 0.8. The simulation time is  
1084 25 s.

1085 This test case is computed using the traditional Roe flux in (118) as well  
1086 as the flux functions A and B in (119) and (120) respectively. The numerical  
1087 solution for the discharge provided by such methods is plotted in space and  
1088 time in Figure 13. Complementary results for the study of the spike in the  
1089 numerical solution are presented in Figure 14, where the evolution in time  
1090 of cell average values are depicted for the 8 leftmost cells on the right hand  
1091 side of the RP (e.g. the first cell on the right of the initial discontinuity is  
1092 depicted in blue, the second one in cyan and so on).

1093 From figures 13 and 14, it is clearly evidenced that whereas the tradi-  
1094 tional Roe solver leads to a high spike in the discharge, which generates a  
1095 shedding of spurious waves, when using the novel flux functions the spike is  
1096 dramatically reduced and hence the shedding of such waves. A closer exam-  
1097 ination of the numerical results evidences that flux function A provides the  
1098 best performance concerning the reduction of the spike, on the other hand,  
1099 flux function B does also reduce this anomalous behavior at the cell where  
1100 the shock is contained but still leaves a small spike behind it. This particu-  
1101 larity of flux function B is clearly noticed in Figure 14 (bottom) where the  
1102 spikes appear to be shifted to the left, which means that it occurs on the  
1103 right side of the wavefront, as observed in Figure 13 (bottom).

1104 In Figure 15 (left), the numerical solutions provided by the traditional  
1105 Roe solver, the solver using flux function A and the solver using flux function  
1106 B is depicted at  $t = 25$  s in purple, green and magenta, respectively. It is  
1107 observed that both the Roe flux and the flux A capture the exact position of  
1108 the shock whereas the flux B underestimates the shock speed, hence providing  
1109 a slightly shifted, though convergent, shock position.

1110 The analysis of the properties of the novel flux functions from [42] can be  
 1111 completed by plotting the numerical results in the phase space. Figure 15  
 1112 (right) shows the exact and approximate Hugoniot locus for the intermediate  
 1113 states between the left and right states of the RP. The exact Hugoniot locus is  
 1114 represented by a red continuous line, the approximate locus for the traditional  
 1115 Roe flux by purple dots, the approximate locus for flux function A by green  
 1116 dots and that for flux function B by magenta dots. As outlined in [42], the  
 1117 optimal locus that prevents the numerical solution from exhibiting any spike  
 1118 and spurious waves is the straight line between the left and right state. It  
 1119 can be observed in Figure 15 (right) that only flux function A achieves this  
 1120 requirement and therefore it is the preferred technique for the reduction of  
 1121 the spike in the SWE.

### 1122 *6.2. Extension of the flux function A to the SWE with source term*

1123 It is evidenced that flux function A is a better choice than B for the  
 1124 resolution of moving hydraulic jumps as it provides a better estimate of the  
 1125 shock speed. Previous numerical experiments do not include the presence of  
 1126 source terms, but most realistic cases are dominated by the action of those  
 1127 sources. In this section, the extension of flux function A to non-homogeneous  
 1128 equations is carried out by means of a suitable correction of the interpola-  
 1129 tion technique that ensures a virtually exact equilibrium between fluxes and  
 1130 source term. In addition to this, the numerical fluxes at the interfaces must  
 1131 be rewritten to account for the source term.

1132 First, it is time to find out which is the most suitable correction of the flux  
 1133 extrapolation to reduce the spike of discharge in both transient and steady  
 1134 cases. Following a similar procedure than in [42], the idea is to find an ap-  
 1135 proximation of such fluxes that ensures the exact equilibrium between fluxes  
 1136 and source term across cell interfaces under steady conditions, while keeping  
 1137 the idea of having an interpolated flux in the cell containing the shock in  
 1138 order to prevent the scheme from using the equilibrium flux, which leads to  
 1139 the spike. To this end, it is first required to find the cell where the shock is  
 1140 contained. We propose to use Roe celerities,  $\tilde{\lambda}^m$  to unequivocally locate such  
 1141 a cell, since it is known that both celerities at the left interface are positive  
 1142 (supercritical flow entering the cell) while a combination of celerities corre-  
 1143 sponding to subcritical conditions (one negative and the other one positive)  
 1144 is identified at the right interface.

1145 Let us consider the cells,  $\Omega_i$ , as single items contained in the domain  $\Omega$   
 1146 such that  $\Omega = \{\Omega_i \mid i \in [1, \dots, N]\}$ . Considering the possibility of multiple

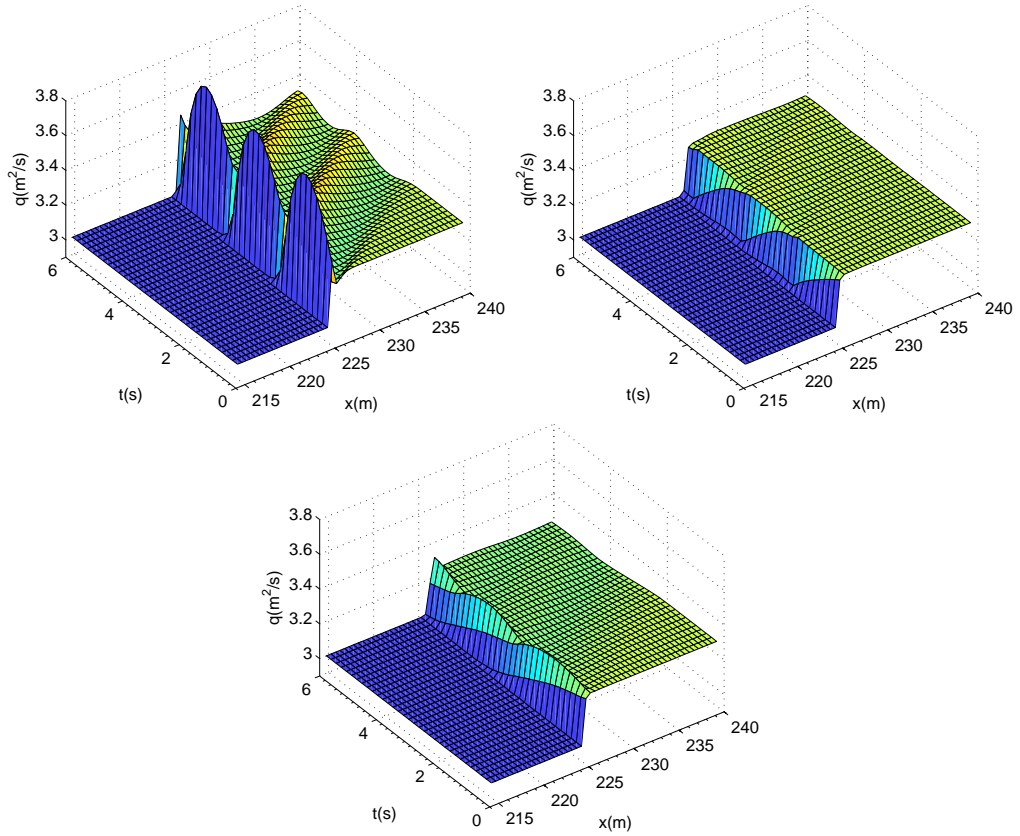


Figure 13: Test case 2. Numerical solution provided by the traditional Roe solver (top-left) as well as the flux functions A (top-right) and B (bottom) proposed in [42] within the time interval  $[0, 6]$  s.

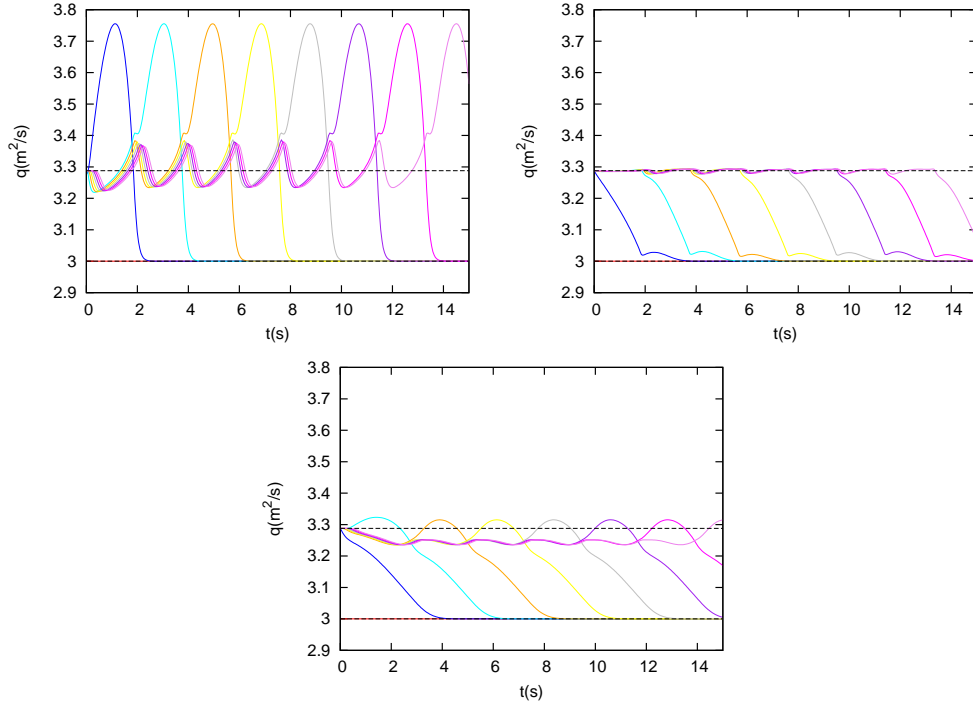


Figure 14: Test case 2. Evolution in time of cell average values for the 8 leftmost cells on the right hand side of the RP using the Roe flux (top-left), flux function A (top-right) and flux function B (bottom).

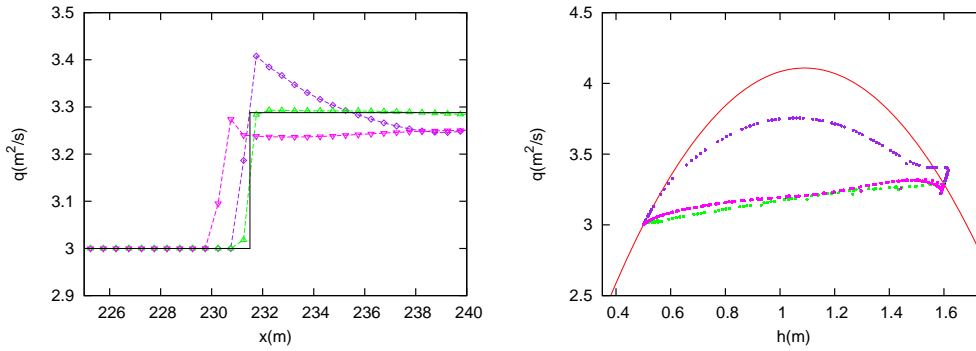


Figure 15: Test case 2. Left: numerical solution using the Roe flux ( $-\diamond-$ ), flux function A ( $-\triangle-$ ) and flux function B ( $-\nabla-$ ) at  $t = 25$  s. Right: exact Hugoniot locus and approximate locus for the Roe flux, flux function A and flux function B.

1147 hydraulic jumps within the domain, we denote the set of cells containing a  
 1148 positive-flow hydraulic jump as

$$\mathcal{D}^+ = \left\{ \Omega_i \mid \Omega_i \in \Omega \wedge \tilde{\lambda}_{i-1/2}^1 \cdot \tilde{\lambda}_{i+1/2}^1 < 0 \wedge h_{i-1} < h_{i+1} \right\} \quad (121)$$

1149 and the set of cells containing a negative-flow hydraulic jump as

$$\mathcal{D}^- = \left\{ \Omega_i \mid \Omega_i \in \Omega \wedge \tilde{\lambda}_{i-1/2}^2 \cdot \tilde{\lambda}_{i+1/2}^2 < 0 \wedge h_{i-1} > h_{i+1} \right\}. \quad (122)$$

1150 Once the hydraulic jumps are found, the following cell-centered fluxes are  
 1151 proposed in order to generate an spike fix

$$\hat{\mathbf{F}}_i = \begin{cases} \mathbf{F}_i & \text{if } \Omega_i \notin \mathcal{D}^+ \cup \mathcal{D}^- \\ \check{\mathbf{F}}_i - (1 - x_{\mathcal{S},i})\bar{\mathbf{S}}_{i-1,i+1} + \bar{\mathbf{S}}_{i-1/2} & \text{if } \Omega_i \in \mathcal{D}^+ \cup \mathcal{D}^- \end{cases} \quad (123)$$

1152 with  $\check{\mathbf{F}}_i$  the interpolated flux in (116),  $\bar{\mathbf{S}}_{i-1,i+1}$  a centered integral of the  
 1153 source term, that can be computed computed as

$$\bar{\mathbf{S}}_{i-1,i+1} = \begin{pmatrix} 0 \\ -g \frac{h_{i-1} + h_{i+1}}{2} (z_{i+1} - z_{i-1}) \end{pmatrix}, \quad (124)$$

1154  $\bar{\mathbf{S}}_{i-1/2}$  the integral of the source term across the left interface, that can be  
 1155 computed as

$$\bar{\mathbf{S}}_{i-1/2} = \begin{pmatrix} 0 \\ -g \frac{h_{i-1} + h_i}{2} (z_i - z_{i-1}) \end{pmatrix}. \quad (125)$$

1156 Parameter  $x_{\mathcal{S},i}$  accounts for the normalized position of the shock inside the  
 1157 cell, here approximated by

$$x_{\mathcal{S},i} = \frac{h_i - h_{i+1}}{h_{i-1} - h_{i+1}}, \quad (126)$$

1158 if considering that the intermediate state is a linear combination of the left  
 1159 and right states (linear Hugoniot)

$$\mathbf{U}_i = x_{\mathcal{S},i} \mathbf{U}_{i-1} + (1 - x_{\mathcal{S},i}) \mathbf{U}_{i+1}, \quad (127)$$

1160 where  $\mathbf{U}_{i-1}$ ,  $\mathbf{U}_i$  and  $\mathbf{U}_{i+1}$  are any arbitrary left, middle and right states  
 1161 defining a hydraulic jump as depicted in Figure 12.

1162 It is worth pointing out that the corrected flux in (123) provides an ap-  
 1163 proximation of the cell-centered flux in the shock cell that converges to the  
 1164 exact steady flux, unlike traditional methods, that only converge to an equi-  
 1165 librium flux (different to the exact flux) that allows the steadiness of the  
 1166 solution. The reason why the proposed technique does not always ensure  
 1167 the exact flux with independence of the grid is due to the assumption we  
 1168 make for the definition of (123): the intermediate state (at cell  $\Omega_i$  where the  
 1169 shock is located) lies on a linear Hugoniot between the left and right states,  
 1170 according to (127), which is not completely true under the presence of a bed  
 1171 step source term. The exact linear Hugoniot would be expressed instead as

$$\mathbf{U}_i = x_{\mathcal{S},i} \mathbf{U}_i^- + (1 - x_{\mathcal{S},i}) \mathbf{U}_i^+, \quad (128)$$

1172 where  $\mathbf{U}_i^-$  and  $\mathbf{U}_i^+$  are the left and right intermediate states at the interfaces  
 1173 of cell  $\Omega_i$ . In spite of this, the approximation in (127) provides a trustworthy  
 1174 approximation of the shock position when solving for  $x_{\mathcal{S},i}$  and what is of most  
 1175 importance, it converges to the exact position as the grid is refined, when  
 1176 dealing with a smooth bed topography.

1177 It is straightforward to show that (123) provides the exact flux under  
 1178 steady conditions by considering the shock located at cell  $\Omega_M$  and applying  
 1179 steady state conditions to the second equation of (123), as follows

$$\hat{\mathbf{F}}_i = \frac{1}{2}(\mathbf{F}_{i-1} + \mathbf{F}_{i+1}) - \frac{1}{2} \tilde{\mathbf{J}}_{i-1,i+1} (\mathbf{U}_{i+1} - 2\mathbf{U}_i + \mathbf{U}_{i-1}) - (1 - x_{\mathcal{S},i}) \bar{\mathbf{S}}_{i-1,i+1} + \bar{\mathbf{S}}_{i-1/2}, \quad (129)$$

1180 where substitution of  $\mathbf{U}_i$  using (127) yields

$$\hat{\mathbf{F}}_i = \frac{1}{2}(\mathbf{F}_{i-1} + \mathbf{F}_{i+1}) + \frac{1}{2}(1 - 2x_{\mathcal{S},i}) \tilde{\mathbf{J}}_{i-1,i+1} (\mathbf{U}_{i+1} - \mathbf{U}_{i-1}) - (1 - x_{\mathcal{S},i}) \bar{\mathbf{S}}_{i-1,i+1} + \bar{\mathbf{S}}_{i-1/2}. \quad (130)$$

1181 From the definition of Roe's Jacobian matrix, we know that  $\tilde{\mathbf{J}}_{i-1,i+1} (\mathbf{U}_{i+1} -$   
 1182  $\mathbf{U}_{i-1}) = \mathbf{F}_{i+1} - \mathbf{F}_{i-1}$  and under steady conditions  $\mathbf{F}_{i+1} - \mathbf{F}_{i-1} = \bar{\mathbf{S}}_{i-1,i+1}$ .  
 1183 Substitution of this term into (130) reads

$$\hat{\mathbf{F}}_i = \frac{1}{2}(\mathbf{F}_{i-1} + \mathbf{F}_{i+1}) + \frac{1}{2}(1 - 2x_{\mathcal{S},i}) \bar{\mathbf{S}}_{i-1,i+1} - (1 - x_{\mathcal{S},i}) \bar{\mathbf{S}}_{i-1,i+1} + \bar{\mathbf{S}}_{i-1/2}, \quad (131)$$

1184 Now, making use of  $\mathbf{F}_{i+1} - \mathbf{F}_{i-1} = \bar{\mathbf{S}}_{i-1,i+1}$  again, it does lead to

$$\hat{\mathbf{F}}_i - \mathbf{F}_{i-1} = \bar{\mathbf{S}}_{i-1/2}, \quad (132)$$

1185 the GRH condition.

1186 Finally, the expression for the numerical fluxes at cell interfaces is pre-  
 1187 sented. Using definitions in Section Appendix A, we can write the non-  
 1188 homogeneous version of the numerical flux in (119) to account for the con-  
 1189 tribution of the source term as

$$\begin{aligned} \mathbf{F}_{i+1/2}^- &= \hat{\mathbf{F}}_i + \sum_{m=1}^I [(\hat{\gamma} - \beta)\tilde{\mathbf{e}}]_{i+\frac{1}{2}}^m, \\ \mathbf{F}_{i+1/2}^+ &= \hat{\mathbf{F}}_{i+1} - \sum_{m=I+1}^{N_\lambda} [(\hat{\gamma} - \beta)\tilde{\mathbf{e}}]_{i+\frac{1}{2}}^m. \end{aligned} \quad (133)$$

1190 where  $\hat{\gamma}$  are the components of  $\hat{\mathbf{\Gamma}}_{i+1/2} = \tilde{\mathbf{P}}_{i+1/2}^{-1} \delta \hat{\mathbf{F}}_{i+1/2}$ , the projection of the  
 1191 jump in the extrapolated fluxes across cell interfaces,  $\hat{\mathbf{F}}_{i+1/2} = \hat{\mathbf{F}}_{i+1} - \hat{\mathbf{F}}_i$ .

### 1192 6.3. Test case 3: Steady jump over smoothly varying bed profile

1193 In this test case, steady solutions for the flow over the following bed  
 1194 elevation profile

$$z(x) = \begin{cases} 0 & \text{if } x < 8 \\ 0.05(x - 8) & \text{if } 8 \leq x \leq 12 \\ 0.2 - 0.05(x - 12)^2 & \text{if } 12 \leq x \leq 14 \\ 0 & \text{if } x > 14 \end{cases} \quad (134)$$

1195 are computed using the proposed technique. The computational domain is  
 1196  $[0, 20]$  and the solution is computed for  $t = 400$  s. CFL number is set to 0.45  
 1197 for all cases and the computational domain is discretized in 100 cells. The  
 1198 discharge is imposed to  $0.6 \text{ m}^2/\text{s}$  upstream to obtain the sonic point at the  
 1199 cell with the maximum bed elevation, that is  $z_{max} = 0.2$ . Downstream, the  
 1200 water depth is also imposed in order to generate the hydraulic jump. Dif-  
 1201 ferent values for  $h$  downstream, are chosen to generate the jump at different  
 1202 locations and assess the performance of the proposed scheme. The complete  
 1203 configuration of boundary conditions is presented in Table 2.

1204 Numerical results provided by the novel scheme are presented for test  
 1205 case 1.A in Figure 16 (top) and compared with the results provided by the  
 1206 traditional Roe solver, depicted in Figure 16 (bottom). No differences can be

1207 noticed when considering the solution for the water surface elevation, but it  
 1208 is clearly evidenced that the spike in the solution for the discharge at the cell  
 1209 where the shock is located is strongly reduced when using the novel numerical  
 1210 technique.

Case	$q_{BC:left}(m^2/s)$	$h_{BC:right}(m)$	Shock position (m)	$x_S$
1.A	0.6	0.6185	13.298	0.01
1.B	0.6	0.6200	13.278	0.11
1.C	0.6	0.6220	13.252	0.24
1.D	0.6	0.6256	13.201	0.495
1.E	0.6	0.6280	13.166	0.67
1.F	0.6	0.6300	13.135	0.825
1.G	0.6	0.6320	13.102	0.99

Table 2: Different boundary condition configurations for Test case 3.

1211 To study the behavior of this spike, the solution for the discharge in the  
 1212 shock cell is depicted for tests cases 1.A-1.G in Figure (17) (left). In this  
 1213 plot, the value of discharge against the normalized shock position has been  
 1214 depicted for the results provided by the traditional Roe solver as well as the  
 1215 modified solver using flux interpolation in [42] and the proposed technique.  
 1216 It can be observed that the method in [42] already helps decreasing the spike  
 1217 of discharge but only when including the correction term, as done in the  
 1218 novel method, the spike is virtually reduced to zero.

1219 As outlined before, the proposed scheme does not always provide the  
 1220 exact discharge in the shock cell, however, the numerical estimate of the  
 1221 discharge in this cell converges to the exact value as the grid is refined. This  
 1222 property is of utmost importance, as the novel scheme can be considered  $L_1$ ,  
 1223  $L_2$  and  $L_\infty$  convergent, while previous schemes were not able to converge  
 1224 when regarding  $L_\infty$  error norm. Convergence rate results for  $L_\infty$  error norm  
 1225 are presented in Figure 17 (right) for the traditional Roe solver and for the  
 1226 proposed scheme. The convergence rate test has been carried out for case  
 1227 1.D using four different grids, composed of 100, 200, 400 and 800 cells. It  
 1228 is worth mentioning that the grid is shifted in order to keep a constant  
 1229 distance between the exact position of the jump and the right cell interface.  
 1230 It is clearly evidenced that the proposed technique allows the scheme to  
 1231 converge to the exact solution as the grid is refined, unlike the traditional Roe  
 1232 solver that does not exhibit any convergence with grid refinement because



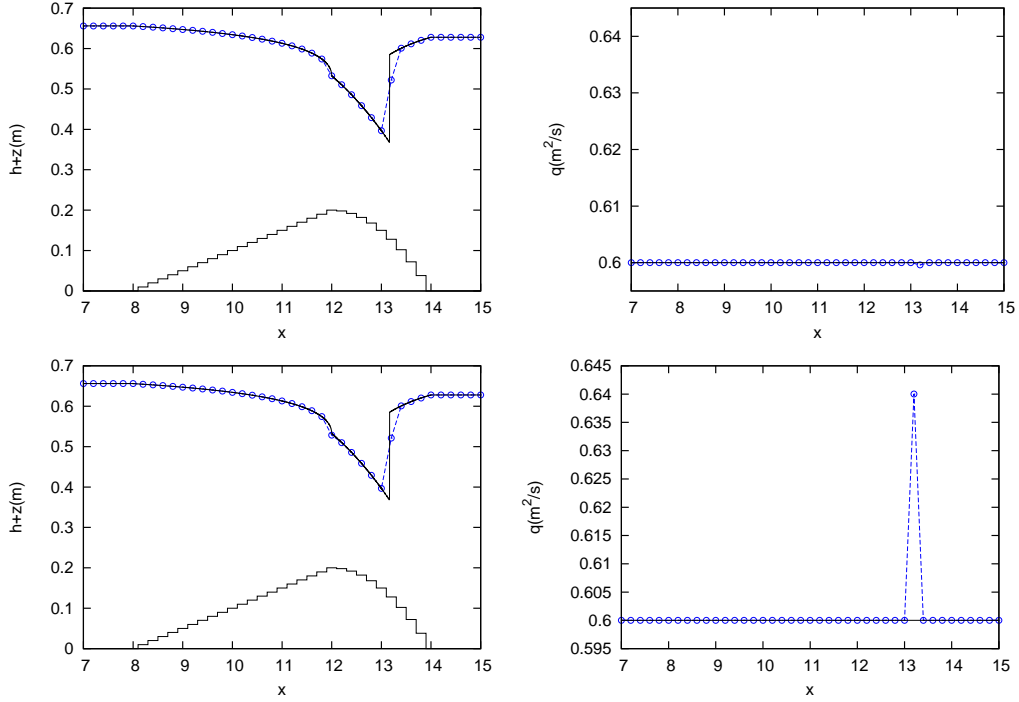


Figure 16: Test case 3. Numerical results for  $h + z$  (left) and  $q$  (right) provided by the proposed spike-reducing method (top) and by the traditional Roe solver (bottom), compared to the exact solution, using 100 cells and CFL=0.45.

1233 the equilibrium discharge at the shock cell is always different than the exact  
 1234 discharge when the shock is not located at cell interfaces.

#### 1235 6.4. Test case 4: Traveling jump over different bed profiles

1236 In this test case, traveling shock waves over different bed elevation pro-  
 1237 files  $z(x)$  are computed. For all bed profiles, the maximum bed elevation  
 1238 is  $z_{max} = 0.2$  m and the bed elevation at the boundaries is zero. To con-  
 1239 struct a solution consisting of a single jump traveling across the domain,  
 1240 we first compute a steady transcritical solution over the bed profile by im-  
 1241 posing a constant discharge upstream of  $q = 0.6$  m<sup>2</sup>/s. When the steady  
 1242 regime is reached, the boundary condition upstream is redefined, imposing  
 1243 now  $q = 0.556749458405104$  m<sup>2</sup>/s and  $h = 0.12$  m, which generates a super-  
 1244 critical state that is connected with the original subcritical state by means of  
 1245 a traveling hydraulic jump, according to the Hugoniot locus. The computa-  
 1246 tional domain is  $[0, 560]$  and the solution is computed at  $t = 610$  s. The CFL

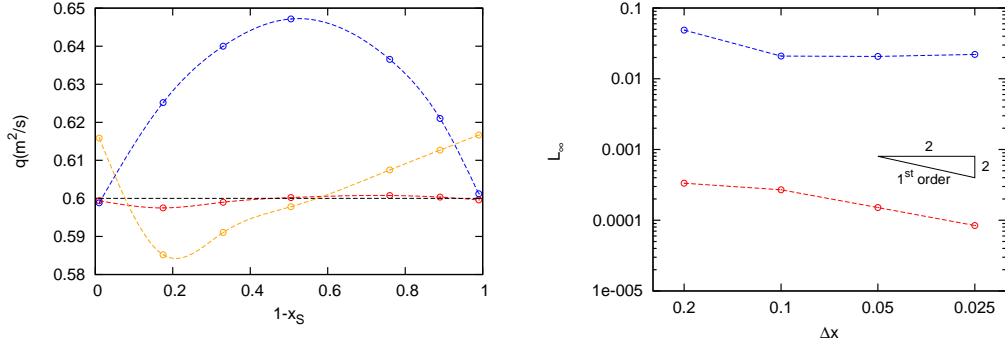


Figure 17: Test case 3. Left: representation of the spike of discharge against the position of the shock within the cell for the traditional Roe flux ( $- \circ -$ ), for the method using the interpolated flux in [42] ( $- \circ -$ ) and for the proposed spike-reducing method ( $- \circ -$ ), using 100 cells and CFL=0.45. Right: convergence rate test for the traditional Roe method ( $- \circ -$ ) and for the proposed method ( $- \circ -$ ), using CFL=0.45.

1247 number is set to 0.45 and the domain is discretized in 140 computational  
 1248 cells.

1249 The bed profile will be constructed as

$$z(x) = \begin{cases} \frac{0.2}{276}(x - 4) + g(x) & \text{if } 4 \leq x < 280 \\ 0.2 - \frac{0.2}{276}(x - 280) & \text{if } 280 \leq x \leq 556 \\ 0 & \text{otherwise} \end{cases} \quad (135)$$

1250 where  $g(x)$  is an additional geometric function that allows to make variations  
 1251 in the basic constant slope profile (when  $g(x) = 0$ ). Three different bed slopes  
 1252 are defined:

- 1253 • *Constant slope (Test 4.1)*: The first test is carried out over a constant  
 1254 slope profile, setting  $g(x) = 0$  in (135).
- 1255 • *Sinusoidal variations in a constant slope (Test 4.2)*: Now, a sinusoidal  
 1256 variation is added to (135) by means of

$$g(x) = \begin{cases} 0.02 \sin(0.04\pi(x - 12)) & \text{if } 12 \leq x < 212 \\ 0 & \text{otherwise} \end{cases} \quad (136)$$

- 1257 • *Discontinuities in the constant slope (Test 4.3)*: Here, some disconti-  
 1258 nuities are added to (135) by means of

$$g(x) = \begin{cases} 0.02 & \text{if } 12 \leq x < 32 \\ -0.02 & \text{if } 32 \leq x < 52 \\ 0.04 & \text{if } 52 \leq x < 72 \\ -0.04 & \text{if } 72 \leq x < 92 \\ 0 & \text{otherwise} \end{cases} \quad (137)$$

1259 Numerical results for tests 4.1, 4.2 and 4.3 are presented in Figures 18,  
 1260 19, 20 and 21. Figure 18 shows the numerical solution at  $t = 610$  s for  
 1261 the water surface elevation and discharge provided by the ARoe scheme and  
 1262 by the proposed spike-reducing method in Section 6.2. For all the test, the  
 1263 SEBF discretization of the source term is chosen. In the figures mentioned  
 1264 above, major differences are observed in the solution of the discharge, which  
 1265 is much more oscillatory when computed by the ARoe method. On the other  
 1266 hand, differences on the water surface elevation are less sensitive to the spike.  
 1267 A space-time representation of the numerical discharge is presented in Fig-  
 1268 ure 19, where the elimination of post-shock oscillations can be observed. In  
 1269 Figure 20, the numerical solution for the water surface elevation and dis-  
 1270 charge inside the cell with maximum bed elevation (cell 71) is plotted in  
 1271 time, showing that the proposed spike-reducing scheme performs adequately  
 1272 with independence of the bed profile, as it prevents the solution from gener-  
 1273 ating oscillations. On the other hand, the numerical solution computed by  
 1274 means of the traditional ARoe scheme shows the oscillations produced by the  
 1275 spike, which travel downwards at a higher speed than the hydraulic jump.  
 1276 In order to carry out an exhaustive analysis on the spike reducing effect of  
 1277 the proposed method, the evolution in time of the numerical solution for the  
 1278 discharge in cells 2 to 11, computed by means of the aforementioned schemes,  
 1279 is plotted in Figure 21. It is evidenced that the numerical solution provided  
 1280 by the proposed scheme completely reduces the spike and only leaves very  
 1281 small peaks that are virtually bounded by the values of the discharge at each  
 1282 side of the shock, hence they are not of any relevance.

### 1283 6.5. Test case 5: Interaction of two jumps over a smooth bed profile

1284 In this case, two hydraulic jumps moving in opposite directions are in-  
 1285 troduced in a steady transcritical flow over the bed profile in (134), inside  
 1286 the domain  $[0, 20]$ . The initial condition corresponds to the steady solution  
 1287 generated when setting  $q = 0.6$  m<sup>2</sup>/s upstream in most part of the domain,  
 1288 and also includes the two jumps as

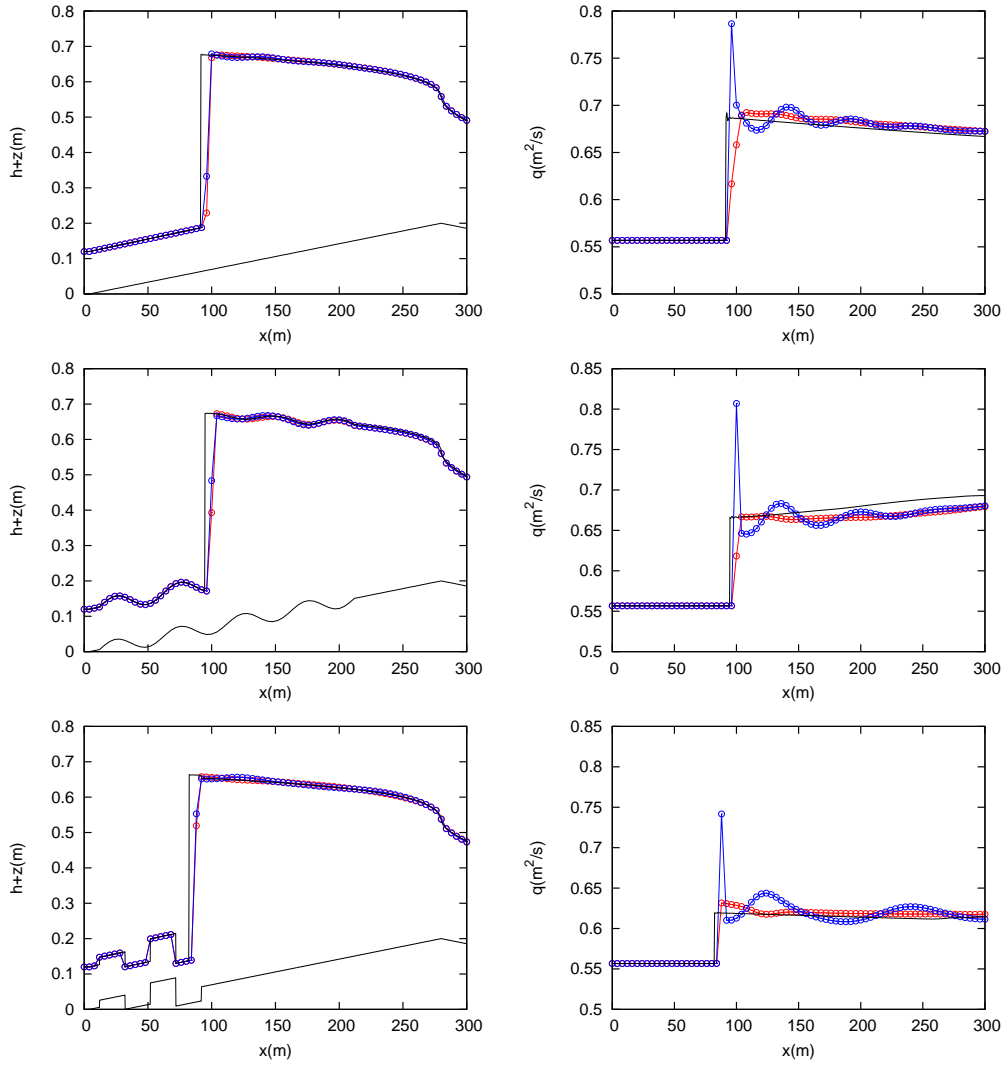


Figure 18: Test case 4. Numerical solution at  $t = 610$  s for the water surface elevation (left) and discharge (right) provided by the traditional Roe flux ( $- \circ -$ ) and by the proposed spike-reducing method ( $- \circ -$ ), using 140 cells and CFL=0.45.

$$\mathbf{U}(x) = \begin{cases} \mathbf{U}_{in} & \text{if } 0 \leq x \leq 1 \\ \mathbf{U}_s & \text{if } 1 < x < 17 \\ \mathbf{U}_{out} & \text{if } 17 \leq x \leq 20 \end{cases} \quad (138)$$

1289 where  $\mathbf{U}_s$  is the steady energy-conservative solution with  $q = 0.6 \text{ m}^2/\text{s}$ ,  $\mathbf{U}_{in} =$

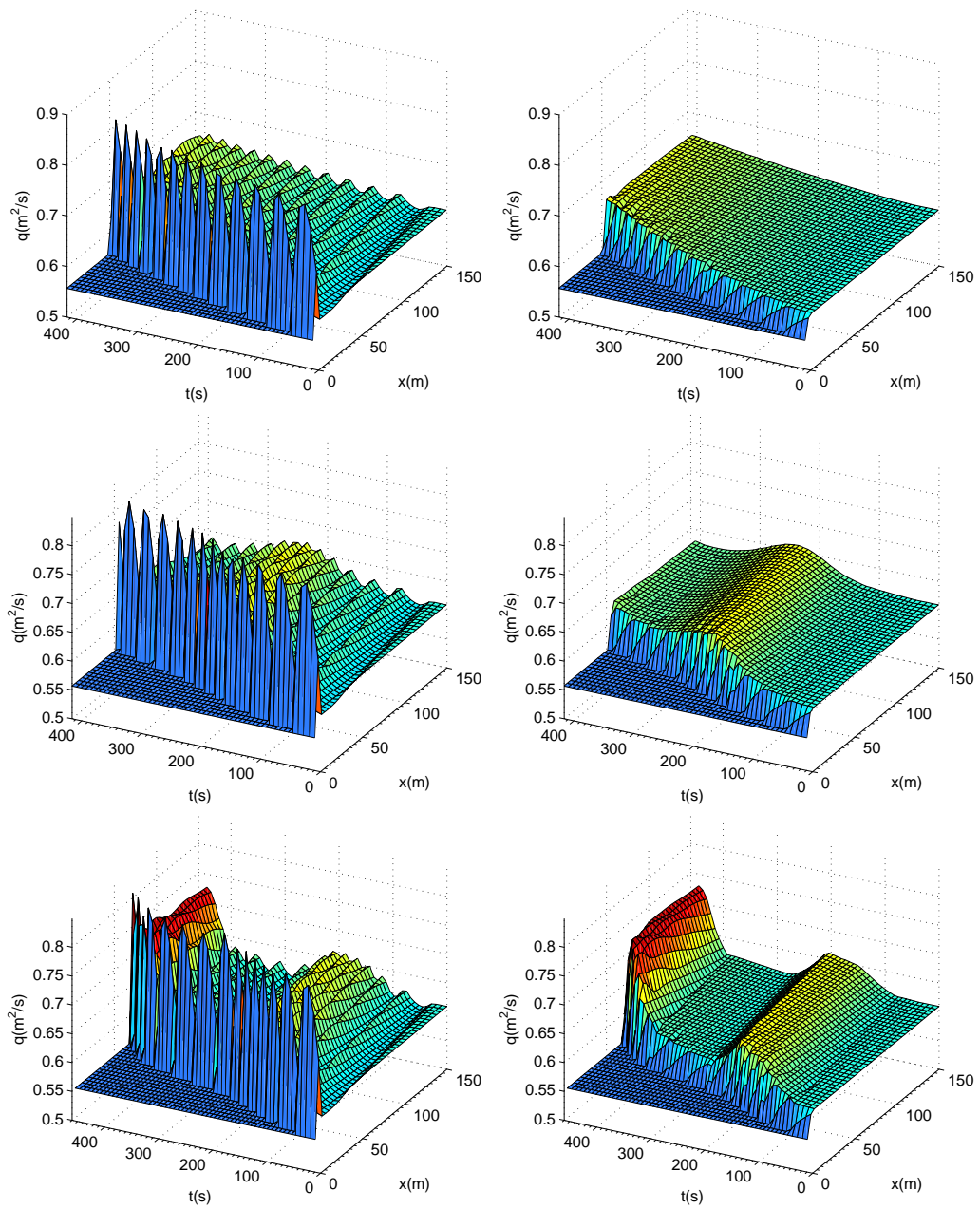


Figure 19: Test case 4. Space-time representation of the numerical discharge provided by the traditional Roe flux (left) and by the proposed spike-reducing method (right), using 140 cells and CFL=0.45.

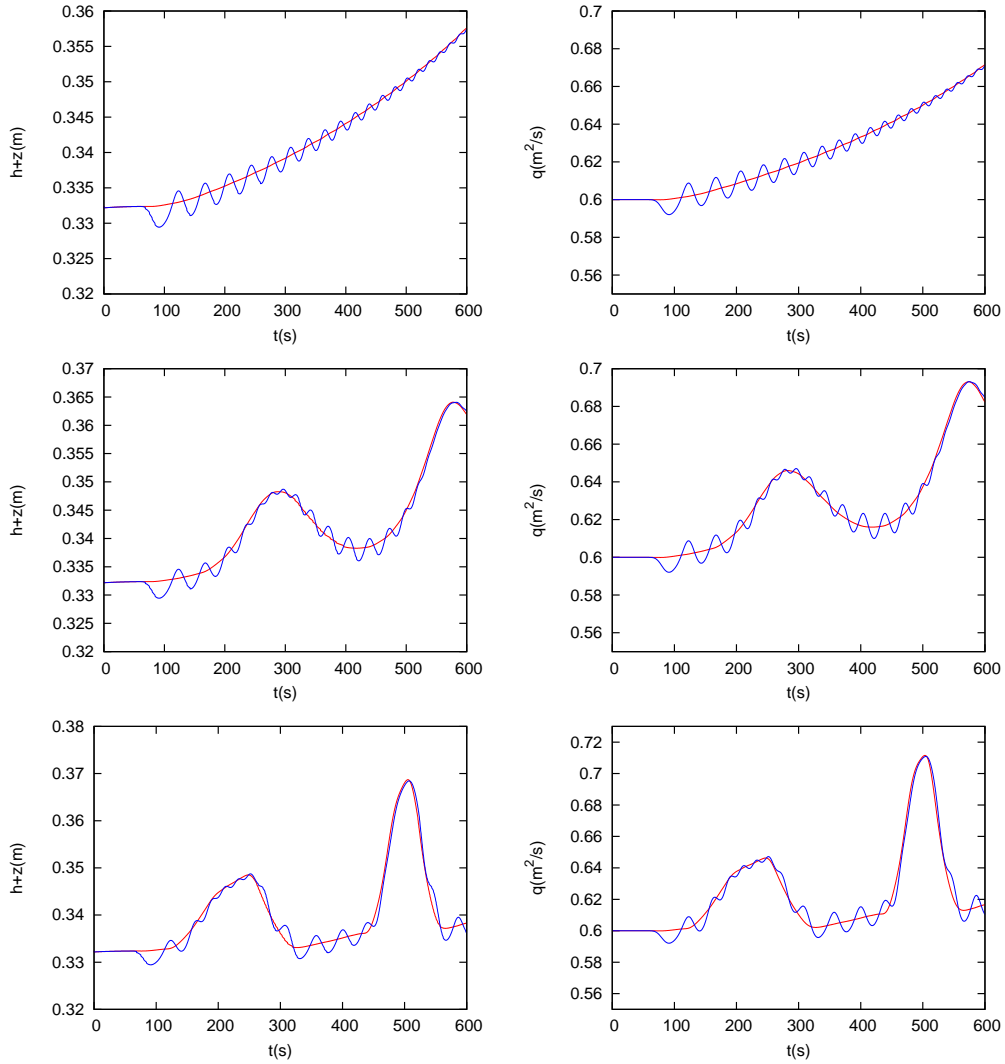


Figure 20: Test case 4. Evolution in time of the numerical solution for the water surface elevation (left) and discharge (right) in the cell with initial  $Fr = 1$  (cell 71) provided by the traditional Roe flux (—) and by the proposed spike-reducing method (—), using 140 cells and  $CFL=0.45$ .

1290  $(h_{in}, q_{in})$  and  $\mathbf{U}_{out} = (h_{out}, q_{out})$ , with  $h_{in} = 0.12$  m,  $q_{in} = 0.556749458405104$   
 1291  $\text{m}^2/\text{s}$ ,  $h_{out} = 0.62$  m and  $q_{out} = 0.410276289759429$   $\text{m}^2/\text{s}$

1292 In order to maintain the hydraulic jumps, the boundary conditions are set  
 1293 supercritical upstream and subcritical downstream, hence we impose  $h = h_{in}$

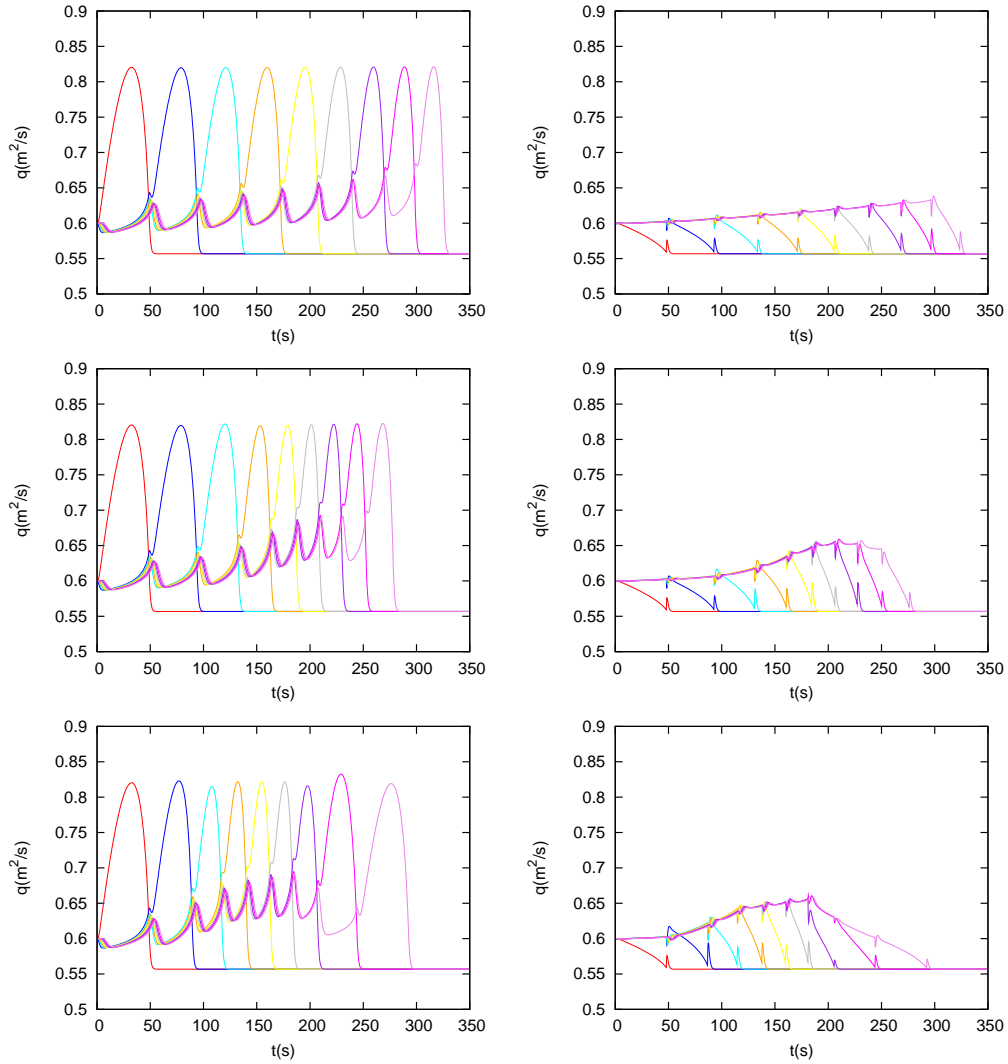


Figure 21: Test case 4. Evolution in time of the numerical solution for the discharge inside cells 2 to 11 provided by the traditional Roe flux (left plot) and by the proposed spike-reducing method (right plot), using 140 cells and CFL=0.45.

1294 and  $q = q_{in}$  upstream and  $h = h_{out}$  downstream. For this test case, we set  
 1295  $\Delta x = 0.2$  and  $\Delta x = 0.1$  m and CFL=0.45. As time goes forward, the left-  
 1296 moving shock on the right decelerates and eventually stops, as the thrust  
 1297 exerted by the bed slope is sufficiently large for it. On the other hand, the  
 1298 right-moving shock on the left does not stop and continuously moves along

1299 the domain. In most part of this simulation, the aforementioned shock moves  
1300 over a flat bottom.

1301 The numerical solution computed by the ARoe scheme and the proposed  
1302 spike-reducing method are presented in Figures 22 and 23, for grid sizes  
1303  $\Delta x = 0.2$  and  $\Delta x = 0.1$  m respectively. The top plots show the solution  
1304 for the water surface elevation and discharge at  $t = 70$  s and the bottom  
1305 plots show the evolution in time of such quantities inside the cell where the  
1306 right jump stops and remains steady. It is observed that the spike-reducing  
1307 method provides a numerical solution much closer to the reference solution as  
1308 no shedding of spurious oscillation occurs, unlike the traditional Roe scheme  
1309 that is unable to avoid those oscillations. It is also observed that oscillations  
1310 are barely reduced with mesh refinement. This is because the spike is still  
1311 present, as the approximate Hugoniot locus of the Roe solver does not depend  
1312 on the discretization (the hydraulic jump is still produced between the same  
1313 left and right states). This means that only the spike-reducing method can  
1314 ensure convergence with mesh refinement.

## 1315 7. Conclusions

1316 This work focuses on the study and design of efficient and robust numeri-  
1317 cal schemes for the computation of hyperbolic conservation laws with source  
1318 terms, with application to the SWE. The goal of the methods proposed here  
1319 is to overcome some present difficulties that have been well documented in  
1320 previous literature, such as the exact conservation of the discrete energy  
1321 (when necessary), the accurate positioning of steady shockwaves and the re-  
1322 duction of the numerical shockwave anomalies arising from slowly-moving  
1323 shocks, among others.

1324 Regarding the conservation of energy in the numerical solution of the  
1325 Shallow Water Equations (SWE), we carry out a theoretical study on the  
1326 relations among variables across the bed step contact wave, showing that  
1327 the conservation of energy can be ensured by imposing conservation of the  
1328 Riemann invariants associated to this wave, or in other words, making the  
1329 Generalized Hugoniot locus (GHL) and the Integral Curve (IC) coincide.  
1330 We consider then the design of a suitable source term discretization (STD)  
1331 that ensures the conservation of energy, showing that the WEBF [25] can be  
1332 derived from these assumptions under the conditions of steady state. The  
1333 WEBF has proven a good performance in a variety of situations, however,



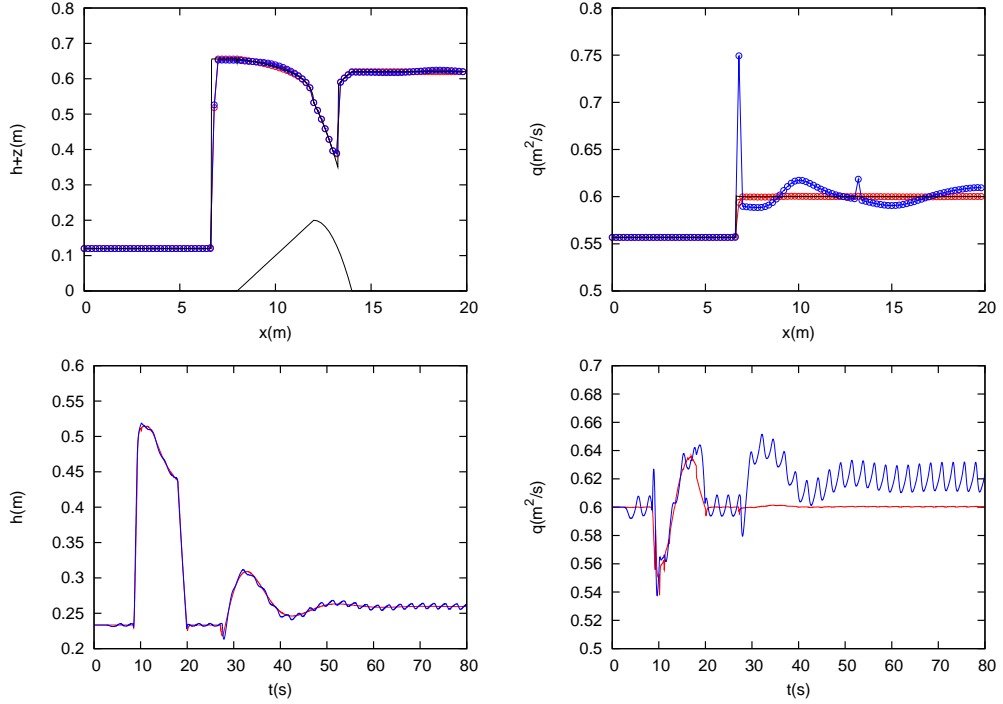


Figure 22: Test case 5. Top: Numerical solution at  $t = 70$  s for the water surface elevation (left) and discharge (right) provided by the traditional Roe flux ( $-o-$ ) and by the proposed spike-reducing method ( $-o-$ ). Bottom: Numerical solution inside cell containing the right jump for the water depth (left) and discharge (right), provided by the traditional Roe flux ( $-$ ) and by the proposed spike-reducing method ( $-$ ). Grid size is set to  $\Delta x = 0.2$ .

1334 when using it for the computation of hydraulic jumps, it is not able to provide  
 1335 an accurate positioning of the discontinuity.

1336 To address the aforementioned issues of shock positioning, a novel dis-  
 1337 cretization of the source term that ensures the exact conservation of the  
 1338 discrete energy while capturing the exact position of the hydraulic jump is  
 1339 proposed. This technique allows to unequivocally identify the position of  
 1340 hydraulic jumps and dissipate the exact amount of energy across them. It is  
 1341 referred to as selective energy balanced formulation (SEBF) of the integral  
 1342 of the source term and can be applied to the ARoe and HLLS solvers, and  
 1343 their high order versions.

1344 Numerical shockwave anomalies in the framework of the SWE, particu-  
 1345 larly the so-called slowly-moving shock anomalies, are also considered in  
 1346 this work. Following the approach in [42], we propose a novel spike-reducing

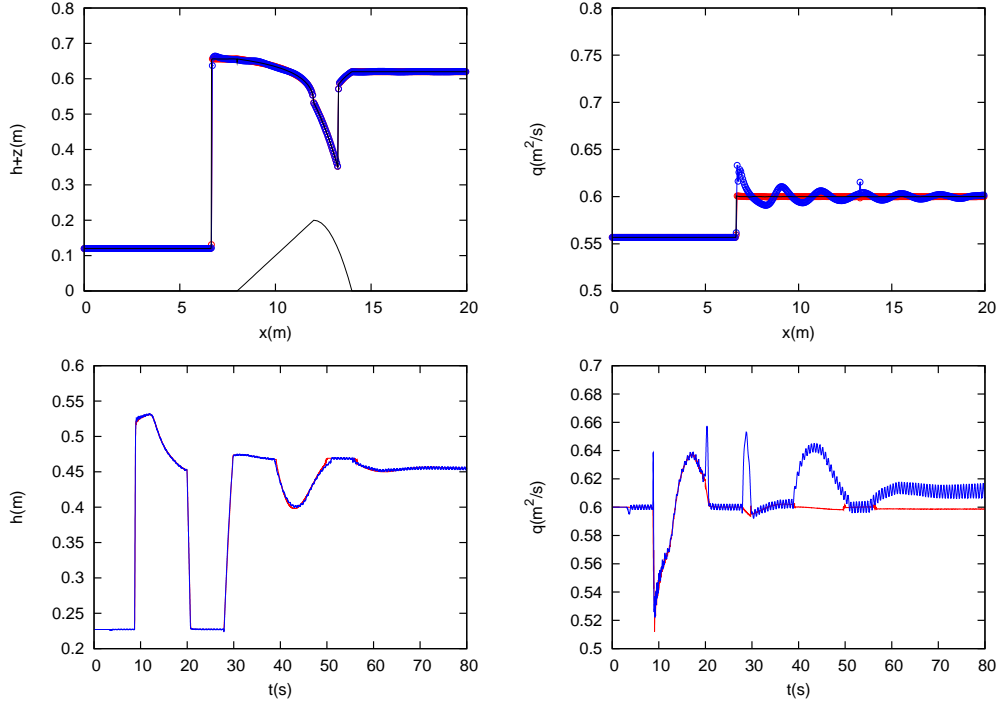


Figure 23: Test case 5. Top: Numerical solution at  $t = 70$  s for the water surface elevation (left) and discharge (right) provided by the traditional Roe flux ( $-o-$ ) and by the proposed spike-reducing method ( $-o-$ ). Bottom: Numerical solution inside cell containing the right jump for the water depth (left) and discharge (right), provided by the traditional Roe flux ( $-$ ) and by the proposed spike-reducing method ( $-$ ). Grid size is set to  $\Delta x = 0.1$ .

1347 flux function for the SWE with varying bed. To this end, we first study the  
 1348 problem of slowly-moving shocks in the SWE and notice that they are only  
 1349 produced when dealing with hydraulic jumps. A complete description of such  
 1350 kind of waves is provided and a thorough study on the shock structure, com-  
 1351 paring exact and Godunov type solutions, is carried out by using the phase  
 1352 space representation. Moreover, prior to the presentation of the proposed  
 1353 technique, flux functions A and B in [42] are assessed for the computation of  
 1354 moving hydraulic jumps over flat bed, evidencing a strong reduction of the  
 1355 spike when using such methods.

1356 The novel spike-reducing flux proposed in this work is computed in the  
 1357 same way than function A [42], but with two main differences. First, a  
 1358 modified flux interpolation technique is carried out in order to account for  
 1359 the contribution of the source. Second, the novel flux function includes the

1360 source strengths across each wave as done in the ARoe solver in [25]. Here  
 1361 we propose to modify the interpolation in [42] by means of a correction term  
 1362 that leads to the exact balance between sources and fluxes in the steady state.  
 1363 This spike fix is based on the hypothesis that the intermediate state should  
 1364 lie on a linear Hugoniot that connects the left and right states, which is not  
 1365 completely general, specially for large discontinuities in the bed elevation,  
 1366 but still leads to satisfactory numerical results for any practical purpose.

1367 The proposed technique is assessed in a variety of situations, including  
 1368 steady and transient cases, over continuous and discontinuous bed. Numerical  
 1369 results evidence that the spike is dramatically reduced to a point where  
 1370 the shedding of spurious waves is virtually not noticeable and also that the  
 1371 proposed scheme leads to a convergent numerical solution because the size  
 1372 of the spike can now be reduced with mesh refinement. For the numerical  
 1373 tests presented in this work, the new scheme does not impose additional sta-  
 1374 bility restrictions and the numerical solution is stable for any CFL number  
 1375 below the traditional bound of 1.0. Numerical results for steady cases with  
 1376 hydraulic jumps are presented, proving that the proposed scheme leads to a  
 1377 convergent solution, even when measured with  $L_\infty$  error norm.

## 1378 **Appendix A. The ARoe solver for systems of $N_\lambda$ waves**

1379 Depending on the nature of the source term, a centered integration of  
 1380 this term may prevent the numerical scheme from preserving the exact bal-  
 1381 ance between fluxes and sources under steady state. This is the case of the  
 1382 so-called geometric source terms, described in (3). In this case, the so-called  
 1383 augmented Riemann solvers are of application for the resolution of the RP,  
 1384 providing an approximation of the numerical fluxes that includes the contri-  
 1385 bution of the source term. Numerical fluxes can be generally expressed as  
 1386  $\mathbf{F}_{i+\frac{1}{2}}^- = \mathbf{F}_{i+\frac{1}{2}}^-(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n; \bar{\mathbf{S}}_{i+1/2})$ ,  $\mathbf{F}_{i-\frac{1}{2}}^+ = \mathbf{F}_{i-\frac{1}{2}}^+(\mathbf{U}_{i-1}^n, \mathbf{U}_i^n; \bar{\mathbf{S}}_{i-1/2})$ , where  $\bar{\mathbf{S}}_{i+1/2}$   
 1387 is a suitable approximation of the integral of the source term across the cell  
 1388 edge.

1389 Riemann Problems are defined at each interface, as depicted in Figure  
 1390 A.24, as

$$\text{RP}(\mathbf{U}_i, \mathbf{U}_{i+1}) : \begin{cases} \frac{\partial \mathbf{U}}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{U})}{\partial x} = \mathbf{S} \\ \mathbf{U}(x, 0) = \begin{cases} \mathbf{U}_i & x < 0 \\ \mathbf{U}_{i+1} & x > 0 \end{cases} \end{cases} \quad (\text{A.1})$$

1391 It is worth mentioning that, for each RP, spatial and temporal variables  
 1392 are redefined setting the reference for the spatial coordinate at  $x_{i+\frac{1}{2}}$  to  $x = 0$   
 1393 and for the time  $t^n$  to  $t = 0$ . Superscript  $n$  is also dropped. As mentioned  
 1394 before, the contribution of the source term is included in the solution of the  
 1395 Riemann Problems as a pointwise quantity at the interface.

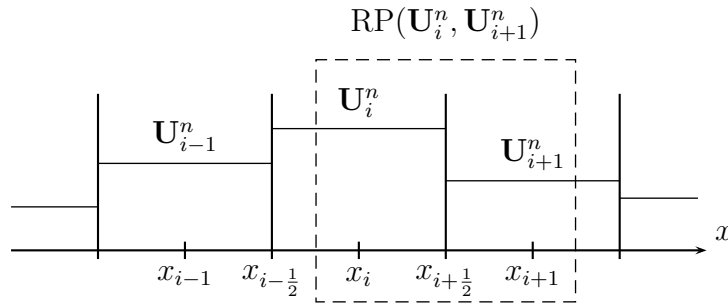


Figure A.24: Neighbouring region of cell  $\Omega_i$  and representation of piecewise defined data, showing RP at  $x_{i+\frac{1}{2}}$  that will be referred to as  $\text{RP}(\mathbf{U}_i^n, \mathbf{U}_{i+1}^n)$ .

1396 RP in (A.1) can be approximated by exactly solving the following con-  
 1397 stant coefficient linear RP [13]

$$\begin{cases} \frac{\partial \hat{\mathbf{U}}}{\partial t} + \tilde{\mathbf{J}}_{i+\frac{1}{2}} \frac{\partial \hat{\mathbf{U}}}{\partial x} = \mathbf{S} \\ \hat{\mathbf{U}}(x, 0) = \begin{cases} \mathbf{U}_i & x < 0 \\ \mathbf{U}_{i+1} & x > 0 \end{cases} \end{cases} \quad (\text{A.2})$$

1398 where  $\hat{\mathbf{U}}(x, t)$  is the approximate solution of (A.1) and  $\tilde{\mathbf{J}}_{i+\frac{1}{2}} = \tilde{\mathbf{J}}_{i+\frac{1}{2}}(\mathbf{U}_i, \mathbf{U}_{i+1})$   
 1399 is a constant matrix defined as a function of left and right states that rep-  
 1400 resents an approximation of the Jacobian at  $x_{i+\frac{1}{2}}$ . This matrix is chosen so  
 1401 that

$$\delta \mathbf{F}_{i+\frac{1}{2}} = \tilde{\mathbf{J}}_{i+\frac{1}{2}} \delta \mathbf{U}_{i+\frac{1}{2}} \quad (\text{A.3})$$

1402 holds [8]. Matrix  $\tilde{\mathbf{J}}_{i+\frac{1}{2}}$  is considered to be diagonalizable with  $N_\lambda$  approximate  
 1403 real eigenvalues

$$\tilde{\lambda}_{i+\frac{1}{2}}^1 < \dots < \tilde{\lambda}_{i+\frac{1}{2}}^I < 0 < \tilde{\lambda}_{i+\frac{1}{2}}^{I+1} < \dots < \tilde{\lambda}_{i+\frac{1}{2}}^{N_\lambda} \quad (\text{A.4})$$

1404 and  $N_\lambda$  eigenvectors  $\tilde{\mathbf{e}}^1, \dots, \tilde{\mathbf{e}}^{N_\lambda}$ . With them, two approximate matrices,  
 1405  $\tilde{\mathbf{P}}_{i+\frac{1}{2}} = (\tilde{\mathbf{e}}^1, \dots, \tilde{\mathbf{e}}^{N_\lambda})_{i+\frac{1}{2}}$  and  $\tilde{\mathbf{P}}_{i+\frac{1}{2}}^{-1}$  are built with the following property

$$\tilde{\mathbf{J}}_{i+\frac{1}{2}} = (\tilde{\mathbf{P}}\tilde{\mathbf{\Lambda}}\tilde{\mathbf{P}}^{-1})_{i+\frac{1}{2}}, \quad \tilde{\mathbf{\Lambda}}_{i+\frac{1}{2}} = \begin{pmatrix} \tilde{\lambda}^1 & & 0 \\ & \ddots & \\ 0 & & \tilde{\lambda}^{N_\lambda} \end{pmatrix}_{i+\frac{1}{2}} \quad (\text{A.5})$$

1406 where  $\tilde{\mathbf{\Lambda}}_{i+\frac{1}{2}}$  is a diagonal matrix with approximate eigenvalues in the main  
 1407 diagonal. System in (A.2) can be transformed using  $\tilde{\mathbf{P}}^{-1}$  matrix as follows

$$\frac{\partial \hat{\mathbf{W}}}{\partial t} + \tilde{\mathbf{\Lambda}}_{i+\frac{1}{2}} \frac{\partial \hat{\mathbf{W}}}{\partial x} = \mathbf{B}_{i+\frac{1}{2}} \quad (\text{A.6})$$

1408 expressing (A.2) in terms of the characteristic variables  $\hat{\mathbf{W}} = \tilde{\mathbf{P}}_{i+\frac{1}{2}}^{-1} \hat{\mathbf{U}}$ , with  
 1409  $\hat{\mathbf{W}} = (\hat{w}^1, \dots, \hat{w}^{N_\lambda})$  and  $\mathbf{B}_{i+\frac{1}{2}} = (\tilde{\mathbf{P}}^{-1} \mathbf{S})_{i+\frac{1}{2}}$

1410 Approximate fluxes on the left and right side of the  $t$  axis,  $\mathbf{F}_i^-$  and  $\mathbf{F}_{i+1}^+$ ,  
 1411 can be derived using the results for the scalar equation. Combination of the  
 1412 solutions for the characteristic variables,  $\hat{w}^m(x, t)$ , allows to construct the  
 1413 numerical fluxes at the interface as [13]

$$\begin{aligned} \mathbf{F}_i^- &= \mathbf{F}_i + \sum_{m=1}^I \left[ (\tilde{\lambda}\alpha - \tilde{\beta}) \tilde{\mathbf{e}} \right]_{i+\frac{1}{2}}^m, \\ \mathbf{F}_{i+1}^+ &= \mathbf{F}_{i+1} - \sum_{m=I+1}^{N_\lambda} \left[ (\tilde{\lambda}\alpha - \tilde{\beta}) \tilde{\mathbf{e}} \right]_{i+\frac{1}{2}}^m, \end{aligned} \quad (\text{A.7})$$

1414 where the set of wave strengths is defined as

$$\mathbf{A}_{i+\frac{1}{2}} = (\alpha^1, \dots, \alpha^{N_\lambda})_{i+\frac{1}{2}}^T = (\tilde{\mathbf{P}}^{-1} \delta \mathbf{U})_{i+\frac{1}{2}}, \quad (\text{A.8})$$

1415 and the set of source strengths as

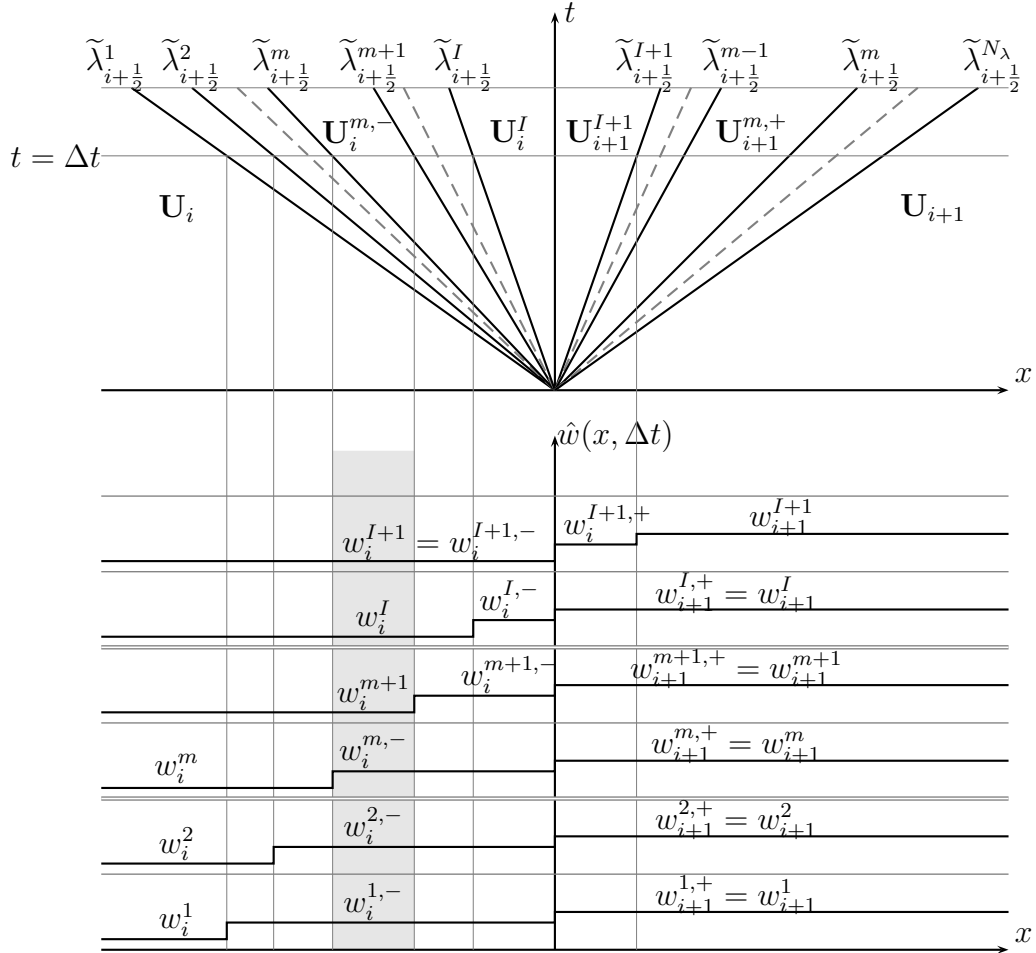


Figure A.25: Upper: Approximate solution  $\hat{\mathbf{U}}(x, t)$ . The solution consists of  $N_\lambda$  inner constant states separated by a stationary contact discontinuity, with celerity  $S = 0$  at  $x = 0$ . Lower: The solution for characteristic variables  $\hat{w}^m(x, t)$  for  $m = 1, \dots, I + 1$  is depicted at  $t = \Delta t$ .

$$\bar{\mathbf{B}}_{i+\frac{1}{2}} = (\bar{\beta}^1, \dots, \bar{\beta}^{N_\lambda})_{i+\frac{1}{2}}^T = \left( \tilde{\mathbf{P}}^{-1} \bar{\mathbf{S}} \right)_{i+\frac{1}{2}}. \quad (\text{A.9})$$

1416 It is worth recalling that  $\delta w_{i+\frac{1}{2}}^m = \alpha_{i+\frac{1}{2}}^m$ . Analogously, if defining  $\delta \mathbf{F}_{i+1/2} =$

1417  $\tilde{\mathbf{P}}_{i+1/2}\mathbf{\Gamma}_{i+1/2}$ , it is straightforward to obtain the following relation

$$\mathbf{\Gamma}_{i+1/2} = \tilde{\mathbf{\Lambda}}_{i+1/2}\tilde{\mathbf{A}}_{i+1/2} \quad (\text{A.10})$$

1418

1419 with  $\mathbf{\Gamma}_{i+1/2} = (\gamma^1, \dots, \gamma^{N_\lambda})_{i+1/2}$ , that allows to rewrite (A.7) as

$$\begin{aligned} \mathbf{F}_{i+1/2}^- &= \hat{\mathbf{F}}_i + \sum_{m=1}^I [(\gamma - \bar{\beta})\tilde{\mathbf{e}}]_{i+\frac{1}{2}}^m, \\ \mathbf{F}_{i+1/2}^+ &= \hat{\mathbf{F}}_{i+1} - \sum_{m=I+1}^{N_\lambda} [(\gamma - \bar{\beta})\tilde{\mathbf{e}}]_{i+\frac{1}{2}}^m. \end{aligned} \quad (\text{A.11})$$

1420

1421 For the sake of simplicity, the term  $(\gamma - \bar{\beta})_{i+\frac{1}{2}}^m$ , or  $(\tilde{\lambda}\alpha - \bar{\beta})_{i+\frac{1}{2}}^m$  analogously,  
 1422 can be expressed as  $(\tilde{\lambda}\theta\alpha)_{i+\frac{1}{2}}^m$ , where  $\theta_{i+\frac{1}{2}}^m = 1 - \bar{\beta}/\tilde{\lambda}\alpha$ . Using this compact  
 1423 form, the difference between left and right states across the interface can be  
 1424 expressed as

$$\mathbf{U}_{i+1}^+ - \mathbf{U}_i^- = \mathbf{U}_{i+1} - \mathbf{U}_i - \sum_{m_1=1}^{N_\lambda} (\theta\alpha\tilde{\mathbf{e}})_{i+\frac{1}{2}}^{m_1} \quad (\text{A.12})$$

1425 where wave contributions can be written in their matrix form as

$$\sum_{m_1=1}^{N_\lambda} (\theta\alpha\tilde{\mathbf{e}})_{i+\frac{1}{2}}^{m_1} = \left( \tilde{\mathbf{P}}\Theta\mathbf{A} \right)_{i+\frac{1}{2}} \quad (\text{A.13})$$

1426 with  $\Theta_{i+\frac{1}{2}} = \text{diag}(\theta_{i+\frac{1}{2}}^1, \theta_{i+\frac{1}{2}}^2, \dots, \theta_{i+\frac{1}{2}}^{N_\lambda})$  a diagonal matrix that allows to rewrite  
 1427  $\tilde{\mathbf{P}}\Theta\mathbf{A} = \tilde{\mathbf{P}}\mathbf{A} - \tilde{\mathbf{P}}\tilde{\mathbf{\Lambda}}^{-1}\bar{\mathbf{B}}$ . Substituting the previous results in (A.12) and  
 1428 noticing that  $\tilde{\mathbf{P}}\mathbf{A}_{i+\frac{1}{2}} = \mathbf{U}_{i+1} - \mathbf{U}_i$ , it becomes

$$\mathbf{U}_{i+1}^+ - \mathbf{U}_i^- = \left( \tilde{\mathbf{P}}\tilde{\mathbf{\Lambda}}^{-1}\bar{\mathbf{B}} \right)_{i+\frac{1}{2}} \quad (\text{A.14})$$

1429 from which it can be observed that the difference between left and right  
 1430 states is only due to the presence of the source term. Expressing  $\bar{\mathbf{B}}_{i+\frac{1}{2}} =$   
 1431  $\left( \tilde{\mathbf{P}}^{-1}\bar{\mathbf{S}} \right)_{i+\frac{1}{2}}$ , the following relation is noticed

$$\bar{\mathbf{S}}_{i+\frac{1}{2}} = \left( \tilde{\mathbf{J}}^{-1} \right)_{i+\frac{1}{2}} (\mathbf{U}_{i+1}^+ - \mathbf{U}_i^-). \quad (\text{A.15})$$

1432 This relation is worth keeping in mind, as it will come along with other  
 1433 derivations within the text.

1434 When using the ARoe numerical fluxes, the first order Godunov scheme  
 1435 in (44) reads

$$\mathbf{U}_i^{n+1} = \mathbf{U}_i^n - \frac{\Delta t}{\Delta x} [\mathbf{F}_i^- - \mathbf{F}_i^+]. \quad (\text{A.16})$$

## 1436 Appendix B. The traditional Roe solver

1437 When considering a homogeneous RP, that is, the contribution of the  
 1438 source term is nil, RH condition across the interface yields  $\mathbf{F}_i^- = \mathbf{F}_{i+1}^+$ , ac-  
 1439 cording to the notation used in this work. Such fluxes are now a unique value  
 1440 and are denoted by  $\mathbf{F}_{i+1/2}^*$ , which can be expressed in terms of the left or  
 1441 right contributions according to (A.7) as follows

$$\begin{aligned} \mathbf{F}_{i+1/2}^* &= \mathbf{F}_i + \sum_{m_1=1}^I (\tilde{\lambda} \alpha \tilde{\mathbf{e}})_{i+\frac{1}{2}}^{m_1} \\ \mathbf{F}_{i+1/2}^* &= \mathbf{F}_{i+1} - \sum_{m_1=I+1}^{N_\lambda} (\tilde{\lambda} \alpha \tilde{\mathbf{e}})_{i+\frac{1}{2}}^{m_1}. \end{aligned} \quad (\text{B.1})$$

1442 Combination of the expressions in (B.1) leads to

$$\mathbf{F}_{i+1/2}^* = \frac{\mathbf{F}_i + \mathbf{F}_{i+1}}{2} - \frac{1}{2} \sum_{m_1=1}^{N_\lambda} (|\tilde{\lambda}| \alpha \tilde{\mathbf{e}})_{i+\frac{1}{2}}^{m_1} \quad (\text{B.2})$$

1443 that can be rewritten in matrix form as

$$\mathbf{F}_{i+1/2}^* = \frac{\mathbf{F}_i + \mathbf{F}_{i+1}}{2} - \frac{1}{2} \left( \tilde{\mathbf{P}} \mid \tilde{\mathbf{\Lambda}} \mid \tilde{\mathbf{A}} \right)_{i+\frac{1}{2}} \quad (\text{B.3})$$

1444 where

$$\mid \tilde{\mathbf{\Lambda}} \mid_{i+\frac{1}{2}} = \begin{pmatrix} |\tilde{\lambda}^1| & & 0 \\ & \ddots & \\ 0 & & |\tilde{\lambda}^{N_\lambda}| \end{pmatrix}_{i+\frac{1}{2}} \quad (\text{B.4})$$

1445 If defining  $\mid \tilde{\mathbf{J}} \mid_{i+\frac{1}{2}} = \left( \tilde{\mathbf{P}} \mid \tilde{\mathbf{\Lambda}} \mid \tilde{\mathbf{P}}^{-1} \right)_{i+\frac{1}{2}}$ , the last term in Equation (B.3) can be  
 1446 rewritten as



$$\left(\tilde{\mathbf{P}} \mid \tilde{\mathbf{\Lambda}} \mid \tilde{\mathbf{A}}\right)_{i+\frac{1}{2}} = \left(\tilde{\mathbf{P}} \mid \tilde{\mathbf{\Lambda}} \mid \tilde{\mathbf{P}}^{-1}\delta\mathbf{U}\right)_{i+\frac{1}{2}} = \left(\mid \tilde{\mathbf{J}} \mid \delta\mathbf{U}\right)_{i+\frac{1}{2}} \quad (\text{B.5})$$

1447 leading to the following intercell homogeneous flux

$$\mathbf{F}_{i+1/2}^* = \frac{\mathbf{F}_i + \mathbf{F}_{i+1}}{2} - \frac{1}{2} \left(\mid \tilde{\mathbf{J}} \mid \delta\mathbf{U}\right)_{i+\frac{1}{2}} \quad (\text{B.6})$$

1448 Analogously, if defining  $\delta\mathbf{F}_{i+1/2} = \tilde{\mathbf{P}}_{i+1/2}\mathbf{\Gamma}_{i+1/2}$ , it is straightforward to  
1449 obtain the following relation

$$\mathbf{\Gamma}_{i+1/2} = \tilde{\mathbf{\Lambda}}_{i+1/2}\tilde{\mathbf{A}}_{i+1/2} \quad (\text{B.7})$$

1450 with  $\mathbf{\Gamma}_{i+1/2} = (\gamma^1, \dots, \gamma^{N_\lambda})_{i+1/2}$ , that can be introduced in (B.3) to obtain

$$\mathbf{F}_{i+1/2}^* = \frac{\mathbf{F}_i + \mathbf{F}_{i+1}}{2} - \frac{1}{2} \text{sgn}(\tilde{\mathbf{J}}_{i+\frac{1}{2}})\delta\mathbf{F}_{i+1/2} \quad (\text{B.8})$$

1451 where  $\text{sgn}(\tilde{\mathbf{J}}_{i+\frac{1}{2}}) = \left(\tilde{\mathbf{P}} \mid \tilde{\mathbf{\Lambda}} \mid \tilde{\mathbf{\Lambda}}^{-1}\tilde{\mathbf{P}}^{-1}\right)_{i+\frac{1}{2}}$  is the upwinding matrix. The pre-  
1452 vious equation can be rewritten as follows

$$\mathbf{F}_{i+1/2}^* = \frac{\mathbf{F}_i + \mathbf{F}_{i+1}}{2} - \frac{1}{2} \sum_{m_1=1}^{N_\lambda} \left(\text{sgn}(\tilde{\lambda})\gamma\tilde{\mathbf{e}}\right)_{i+\frac{1}{2}}^{m_1} \quad (\text{B.9})$$

1453 or, analogously to equation (B.1)

$$\begin{aligned} \mathbf{F}_{i+1/2}^* &= \mathbf{F}_i + \sum_{m_1=1}^I (\gamma\tilde{\mathbf{e}})_{i+\frac{1}{2}}^{m_1} \\ \mathbf{F}_{i+1/2}^* &= \mathbf{F}_{i+1} - \sum_{m_1=I+1}^{N_\lambda} (\gamma\tilde{\mathbf{e}})_{i+\frac{1}{2}}^{m_1}. \end{aligned} \quad (\text{B.10})$$

1454 When using the homogeneous Roe fluxes, the first order Godunov scheme  
1455 in (44) reads

$$\mathbf{U}_i^{n+1} = \mathbf{U}_i^n - \frac{\Delta t}{\Delta x} [\mathbf{F}_{i+1/2}^* - \mathbf{F}_{i-1/2}^*] \quad (\text{B.11})$$

1456 and can be used to solve a homogeneous PDE.

1457 **Acknowledgment**

1458 This work has been funded by the Spanish Ministerio de Economía y  
1459 Competitividad under research project CGL2015-66114-R.

- 1460 [1] S. Sahnima, F. Benkhaldounb, F. Alcrudo, A sign matrix based scheme  
1461 for non-homogeneous PDE's with an analysis of the convergence stag-  
1462 nation phenomenon, *J. Comput. Phys.* 148 (2007) 1753-1783.
- 1463 [2] L. O. Mller and E. F. Toro, A global multiscale mathematical model  
1464 for the human circulation with emphasis on the venous system, *Int. J.*  
1465 *Numer. Meth. Biomed. Engng.* 30 (2014)) 681-725 .
- 1466 [3] S.K. Godunov, Finite difference methods for the computation of discon-  
1467 tinuous solutions of the equations of fluid dynamics, *Mat. Sb.* 47 (1959)  
1468 271-306.
- 1469 [4] E. Godlewski, P.-A. Raviart *Numerical Approximation of Hyperbolic*  
1470 *Systems of Conservation Laws.* Springer Science and Business Media,  
1471 Berlin, 2013.
- 1472 [5] D.L. George. Augmented Riemann solvers for the shallow water equa-  
1473 tions over variable topography with steady states and inundation, *J.*  
1474 *Comput. Phys.* 227 (2008) 3089-3113.
- 1475 [6] J. Murillo, J. Burguete, P. Brufau, P. García-Navarro. The influence of  
1476 source terms on stability, accuracy and conservation in two-dimensional  
1477 shallow flow simulation using triangular finite volumes, *Int. J. Numer.*  
1478 *Meth. Fluids* (2007) 54 543-590.
- 1479 [7] J. Murillo, P. García-Navarro, Augmented Roe's approaches for Rie-  
1480 mann problems including source terms: definition of stability region with  
1481 application to the shallow water equations with rigid and deformable  
1482 bed. In M. E. Vázquez-Cendón and A. Hidalgo and P. García-Navarro  
1483 and L. Cea, eds., *Numerical Methods for Hyperbolic Equations. Theory*  
1484 *and Applications*, pages 149-154. Taylor-Francis Group, 2013.
- 1485 [8] Roe, Approximate Riemann solvers, parameter vectors, and difference  
1486 schemes, *J. Comput. Phys.* 43 (1981) 357-372.

- 1487 [9] A. Harten, P. Lax, B. van Leer, On upstream differencing and Godunov  
1488 type methods for hyperbolic conservation laws, *SIAM review.* 25 (1983)  
1489 35–61.
- 1490 [10] E.F. Toro, M. Spruce, W. Spears, Restoration of the contact surface in  
1491 the HLL Riemann solver, *Shock Waves.* 4 (1994) 25–34.
- 1492 [11] J. Murillo, P. García-Navarro, Weak solutions for partial differential  
1493 equations with source terms: application to the shallow water equations,  
1494 *J. Comput. Phys.* 229 (2010) 4327–4368.
- 1495 [12] J. Murillo, P. García-Navarro, Augmented versions of the HLL and  
1496 HLLC Riemann Solvers including source terms in one and two dimen-  
1497 sions for shallow flow applications, *J. Comput. Phys.* 231 (2012) 6861–  
1498 6906.
- 1499 [13] J. Murillo and A. Navas-Montilla, A comprehensive explanation and  
1500 exercise of the source terms in hyperbolic systems using Roe type solu-  
1501 tions. Application to the 1D-2D shallow water equations, *Advances in*  
1502 *Water Resources* 98 (2016) 70–96.
- 1503 [14] G. Rosatti, L. Begnudelli, The Riemann Problem for the one-  
1504 dimensional, free-surface Shallow Water Equations with a bed step:  
1505 theoretical analysis and numerical simulations, *J. Comput. Phys.* 229  
1506 (2010) 760-787.
- 1507 [15] A. Bermudez and M.E. Vázquez-Cendón, Upwind methods for hyper-  
1508 bolic conservation laws with source terms, *Comput. Fluids.* 23 (1994)  
1509 1049–1071.
- 1510 [16] M.E. Vázquez-Cendón. Improved treatment of source terms in upwind  
1511 schemes for the shallow water equations in channels with irregular ge-  
1512 ometry, *J. Comput. Phys.* 148 (1999) 497–498.
- 1513 [17] J.M. Greenberg, A.Y. Leroux, A well-balanced scheme for the numerical  
1514 processing of source terms in hyperbolic equations, *SIAM J. Numer.*  
1515 *Anal.* 33 (1996) 1–16.
- 1516 [18] P. García-Navarro, M.E. Vázquez-Cendón. On numerical treatment of  
1517 the source terms in the shallow water equations, *Comput. and Fluids.*  
1518 29 (2000) 951–979.

- 1519 [19] A. Chinnayya, A.-Y. LeRoux, N. Seguin, A well-balanced numerical  
1520 scheme for the approximation of the shallow water equations with to-  
1521 pography: the resonance phenomenon, *Int. J. Finite Vol.* 1 (2004) 1–33.
- 1522 [20] M. E. Hubbard, P. García-Navarro, Flux difference splitting and the  
1523 balancing of source terms and flux gradients. *J. Comp. Phys.* 165 (2000)  
1524 89–125.
- 1525 [21] Noelle, S., Xing, Y., Shu, C., High-order well-balanced finite volume  
1526 WENO schemes for shallow water equation with moving water, *J. Com-  
1527 put. Phys.* 226 (2007) 29–58.
- 1528 [22] U.S. Fjordholm, S. Mishra, E. Tadmor, Well-balanced and energy stable  
1529 schemes for the shallow water equations with discontinuous topography,  
1530 *J. Comput. Phys.* 230 (2011) 5587–5609.
- 1531 [23] M.J. Castro Díaz , J.A. López-García, Carlos Parés, High order exactly  
1532 well-balanced numerical methods for shallow water systems, *J. Comput.  
1533 Phys.* 246 (2013) 242–264.
- 1534 [24] Y. Xing, Exactly well-balanced discontinuous Galerkin methods for the  
1535 shallow water equations with moving water equilibrium, *J. Comput.  
1536 Phys.* 257 (2014) 536–553.
- 1537 [25] J. Murillo, P. García-Navarro, Energy balance numerical schemes for  
1538 shallow water equations with discontinuous topography, *J. Comput.  
1539 Phys.* 236 (2012) 119–142.
- 1540 [26] J. Murillo, P. García-Navarro, Accurate numerical modeling of 1D flow  
1541 in channels with arbitrary shape. Application of the energy balanced  
1542 property, *J. Comput. Phys.* 260 (2014) 222–248.
- 1543 [27] A. Navas-Montilla, J. Murillo, Energy balanced numerical schemes with  
1544 very high order. The Augmented Roe Flux ADER scheme. Application  
1545 to the shallow water equations, *J. Comput. Phys.* 290 (2015) 188–218.
- 1546 [28] A. Navas-Montilla, J. Murillo, Asymptotically and exactly energy bal-  
1547 anced augmented flux-ADER schemes with application to hyperbolic  
1548 conservation laws with geometric source terms, *J. Comput. Phys.* 317  
1549 (2016) 108–147.

- 1550 [29] A. Chinnayya, A. Y. LeRoux, N. Seguin, A well-balanced numerical  
1551 scheme for the approximation of the shallow-water equations with to-  
1552 pography: the resonance phenomenon, *Int. J. Finite Volumes* 1 (2004)  
1553 1-33.
- 1554 [30] T. Galloet, J.M. Herard, N. Seguin, Some approximate Godunov  
1555 schemes to compute shallow-water equations with topography, *Com-  
1556 puters and Fluids* 32 (2003) 479-513.
- 1557 [31] F. Alcrudo, F. Benkhaldoun, Exact solutions to the Riemann problem  
1558 of the shallow water equations with a bottom step, *Comput. Fluids* 30  
1559 (2001) 643–671.
- 1560 [32] R. Bernetti, V.A. Titarev, E.F. Toro, Exact solution of the Riemann  
1561 problem for the shallow water equations with discontinuous bottom ge-  
1562 ometry, *J. Comput. Phys.* 227 (2008) 3212–3243.
- 1563 [33] D. S. Balsara, T. Rumpf, M. Dumbser, C.-D. Munz, Efficient, high  
1564 accuracy ADER-WENO schemes for hydrodynamics and divergence-free  
1565 magnetohydrodynamics, *J. Comput. Phys.*, 228 (2009) 2480-2516.
- 1566 [34] F. Franzini, S. Soares-Frazão Efficiency and accuracy of Lateralized  
1567 HLL, HLLS and Augmented Roes scheme with energy balance for river  
1568 flows in irregular channels, *Appl. Math. Model.* 40 (2016) 7427–7446.
- 1569 [35] K.M. Peery and S.T. Imlay, Blunt-body flow simulations, AIAA paper,  
1570 (1988) 88-2924.
- 1571 [36] K. Kitamura, E Shima, and PL Roe, Three-dimensional carbuncles and  
1572 euler fluxes, *Proceedings of the 48th AIAA Aerospace Sciences Meeting*  
1573 (2010).
- 1574 [37] T. W. Roberts, The behavior of flux difference splitting schemes near  
1575 slowly moving shock waves, *J. Comput. Phys.*, 90 (1990) 141–160.
- 1576 [38] D. W. Zaide, Numerical Shockwave Anomalies, PhD thesis, Aerospace  
1577 Engineering and Scientific Computing, University of Michigan, 2012.
- 1578 [39] R. S. Myong, and P. L. Roe, Shock waves and rarefaction waves in  
1579 magnetohydrodynamics. part 2. the mhd system. *J. Plasma Ph.*, 58  
1580 (1997) 21–552.

- 1581 [40] M. Arora, P. L. Roe, On postshock oscillations due to shock capturing  
1582 schemes in unsteady flows, *J. Comput. Phys.*, 130 (1997) 25–40.
- 1583 [41] W.F. Noh, Errors for calculations of strong shocks using an artificial  
1584 viscosity and an artificial heat flux, *J. Comput. Phys.*, 72 (1987) 78-120.
- 1585 [42] D. W. Zaide, P. L. Roe, Flux functions for reducing numerical shockwave  
1586 anomalies. ICCFD7, Big Island, Hawaii, (2012) 9–13.
- 1587 [43] G. Cameron, An analysis of the errors caused by using artificial viscosity  
1588 terms to represent steady-state shock waves. *J. Comput. Phys.* 1 (1966)  
1589 1–20.
- 1590 [44] A. Emery, An evaluation of several differencing methods for inviscid  
1591 fluid flow problems, *J. Comput. Phys.*, 2 (1968) 306–331.
- 1592 [45] S. Karni, S. Canic, Computations of slowly moving shocks, *J. Comput.*  
1593 *Phys.*, 136 (1997) 132–139.
- 1594 [46] S. Jin, J. G. Liu, The Effects of Numerical Viscosities, *J. Comput. Phys.*,  
1595 126 (1996) 373–389.
- 1596 [47] M. H. Carpenter, J. H. Casper, Accuracy of Shock Capturing in Two  
1597 Spatial Dimensions, *AIAA Journal*, 37 (1999) 1072–1079.
- 1598 [48] N. K. Yamaleev, M. H. Carpenter, On accuracy of adaptive grid methods  
1599 for captured shocks, *J. Comput. Phys.*, 181 (2002) 280–316.
- 1600 [49] Y. Stiriba, R. Donat, A numerical study of postshock oscillations in  
1601 slowly moving shock waves, *Comput. Math. with Appl.*, 46 (2003) 719–  
1602 739.
- 1603 [50] E. Johnsen, S. K. Lele, Numerical errors generated in simulations of  
1604 slowly moving shocks, Center for Turbulence Research, Annual Research  
1605 Briefs, (2008) 1–12.
- 1606 [51] R. J. LeVeque, Finite volume methods for hyperbolic problems (Vol.  
1607 31). Cambridge university press, (2002).
- 1608 [52] E.F. Toro, Riemann solvers and numerical methods for fluid dynamics:  
1609 a practical introduction, third ed., Springer-Verlag, Berlin, Heidelberg,  
1610 2009.

- 1611 [53] T.J. Barth, "Some Notes on Shock-Resolving Flux Functions Part 1:  
1612 Stationary Characteristics," NASA TM-101087 (1989)
- 1613 [54] P.L. Roe, Fluctuations and Signals - A Framework for Numerical Evolu-  
1614 tion Problems, Numerical Methods for Fluid Dynamics, edited by K. W.  
1615 Morton, and M. J. Baines, Academic Press, New York, (1982) 219–257.
- 1616 [55] K. Kitamura, P.L. Roe, F. Ismail, Evaluation of Euler fluxes for hyper-  
1617 sonic flow computations, AIAA Journal, 47 (2009) 44–53