



UNIVERSIDAD DE MÁLAGA

PROGRAMA DE DOCTORADO EN MATEMÁTICAS

FACULTAD DE CIENCIAS

DEPARTAMENTO DE ANÁLISIS MATEMÁTICO, ESTADÍSTICA E  
INVESTIGACIÓN OPERATIVA Y MATEMÁTICA APLICADA

# Semi-implicit well-balanced schemes for 1D shallow flows

CELIA CABALLERO CÁRDENAS

PHD THESIS

ADVISORS:

MANUEL JESÚS CASTRO DÍAZ, MARÍA DE LA LUZ MUÑOZ RUIZ

UNIVERSIDAD DE MÁLAGA

2024





UNIVERSIDAD  
DE MÁLAGA

AUTORA: Celia Caballero Cárdenas

 <https://orcid.org/0000-0003-0073-3610>

EDITA: Publicaciones y Divulgación Científica. Universidad de Málaga



Esta obra está bajo una licencia de Creative Commons Reconocimiento-NoComercial-SinObraDerivada 4.0 Internacional:

<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>

Cualquier parte de esta obra se puede reproducir sin autorización pero con el reconocimiento y atribución de los autores.

No se puede hacer uso comercial de la obra y no se puede alterar, transformar o hacer obras derivadas.

Esta Tesis Doctoral está depositada en el Repositorio Institucional de la Universidad de Málaga (RIUMA): [riuma.uma.es](http://riuma.uma.es)





## DECLARACIÓN DE AUTORÍA Y ORIGINALIDAD DE LA TESIS PRESENTADA PARA OBTENER EL TÍTULO DE DOCTOR

D./Dña CELIA CABALLERO CÁRDENAS

Estudiante del programa de doctorado EN MATEMÁTICAS de la Universidad de Málaga, autor/a de la tesis, presentada para la obtención del título de doctor por la Universidad de Málaga, titulada: SEMI-IMPLICIT WELL-BALANCED SCHEMES FOR 1D SHALLOW FLOWS

Realizada bajo la tutorización de MANUEL JESÚS CASTRO DÍAZ y dirección de MANUEL JESÚS CASTRO DÍAZ Y MARÍA DE LA LUZ MUÑOZ RUIZ (si tuviera varios directores deberá hacer constar el nombre de todos)

DECLARO QUE:

La tesis presentada es una obra original que no infringe los derechos de propiedad intelectual ni los derechos de propiedad industrial u otros, conforme al ordenamiento jurídico vigente (Real Decreto Legislativo 1/1996, de 12 de abril, por el que se aprueba el texto refundido de la Ley de Propiedad Intelectual, regularizando, aclarando y armonizando las disposiciones legales vigentes sobre la materia), modificado por la Ley 2/2019, de 1 de marzo.

Igualmente asumo, ante a la Universidad de Málaga y ante cualquier otra instancia, la responsabilidad que pudiera derivarse en caso de plagio de contenidos en la tesis presentada, conforme al ordenamiento jurídico vigente.

En Málaga, a 2 de DICIEMBRE de 2023

Fdo.: CELIA CABALLERO CÁRDENAS Doctorando/a	Fdo.: MANUEL JESÚS CASTRO DÍAZ Tutor/a
Fdo.: MANUEL JESÚS CASTRO DÍAZ Director/es de tesis	
MARÍA DE LA LUZ MUÑOZ RUIZ	





UNIVERSIDAD  
DE MÁLAGA

D. Manuel Jesús Castro Díaz, Catedrático del Departamento de Análisis Matemático, Estadística e Investigación Operativa, y Matemática Aplicada de la Universidad de Málaga, y D<sup>a</sup>. María de la Luz Muñoz Ruiz, Profesora Titular del Departamento de Análisis Matemático, Estadística e Investigación Operativa, y Matemática Aplicada de la Universidad de Málaga,

**CERTIFICAN:**

- Que Celia Caballero Cárdenas, con grado en Matemáticas y máster en Ingeniería Matemática, ha realizado en el Departamento de Análisis Matemático, Estadística e Investigación Operativa, y Matemática Aplicada de la Universidad de Málaga, bajo nuestra dirección, el trabajo de investigación correspondiente a su Tesis Doctoral, titulado:

**Semi-implicit well-balanced schemes for 1D shallow flows**

- Que las publicaciones que avalan la tesis no han sido utilizadas en ninguna otra tesis doctoral ni lo serán en el futuro, puesto que todos los coautores de los trabajos, con la excepción de la doctoranda, son doctores.

Revisado el presente trabajo, estimamos que puede ser presentado al Tribunal que ha de juzgarlo. Y para que constate a efectos de lo establecido en el Reglamento 4/2022, de 24 de octubre, de la Universidad de Málaga sobre los estudios de doctorado, autorizamos la presentación de este trabajo en la Universidad de Málaga.

Málaga, 5 de diciembre de 2023

Dr. Manuel Jesús Castro Díaz

Dra. María de la Luz Muñoz Ruiz



UNIVERSIDAD  
DE MÁLAGA

# Agradecimientos

Hacer una tesis no es tarea fácil, pero si algo he aprendido durante estos años es que rodearte de la gente adecuada es lo más importante, y yo he tenido la suerte de contar con un equipo inmejorable.

En primer lugar, me gustaría agradecer a Carlos Parés, que siendo mi tutor de TFG me transmitió la pasión por la Matemática Aplicada y me animó a adentrarme en el mundo de la investigación, por lo que es uno de los culpables de que acabara optando por este camino.

Por otro lado, no encuentro palabras suficientes para expresar mi agradecimiento a Manolo, Mariluz y Tomás. Vuestra dedicación incansable, generosidad y apoyo constante han sido fundamentales para mí durante esta etapa. Gracias por enseñarme tanto, por ayudarme siempre que lo he necesitado y por vuestra paciencia con mi escaso blablablá.

Gracias también al resto del grupo EDANYA: Jorge Macías, José María Gallardo, José Manuel González, Cipriano Escalante, Marc de la Asunción, María López, Carlos Sánchez, Sergio Ortega; por abrirme sus puertas y hacerme sentir como una más desde el primer día. Y por supuesto, muchísimas gracias a mis compañeros durante estos años: Ernesto Pimentel, Ernesto Guerrero, Juan F. Rodríguez, Alejandro González, León Ávila y Francisco Kieninger.

Evidentemente, gracias también a Irene. Qué suerte haber encontrado a una amiga con la que recorrer este nada fácil camino y con la que compartir las alegrías y los agobios como si fueran propios.

Me gustaría agradecer también a Christophe Chalons, por acogerme tan bien durante mi estancia en Versalles y porque trabajar y aprender con él ha sido un gusto y una suerte para mí.

Por último, tengo que dar las gracias especialmente a mi familia. Mamá, Papá, Jorgito: gracias por ser un apoyo constante en todos los aspectos de mi vida. Todo es mucho más fácil con vosotros a mi lado. Gracias a Bibiki por apoyarme en todas mis decisiones, por estar siempre, en los momentos buenos y en los no tan buenos, y por quererme tan bien. Y gracias también al resto de mi familia y a mis amigos: Rubén, Salva, Laura, Anna, Paula, Óscar, Elena, Belén, Clara, Marta, Blanca, Marina y a los amigos que me regaló el Colesep.

Gracias a todos y a todas.



# Contents

<b>List of figures</b>	<b>iii</b>
<b>Resumen</b>	<b>vii</b>
<b>Introduction</b>	<b>xix</b>
<b>Abstract</b>	<b>xxiii</b>
<b>1 Mathematical settings</b>	<b>1</b>
1.1 Finite volume numerical methods for 1D systems of balance laws . . . . .	1
1.1.1 Systems of balance laws . . . . .	1
1.1.2 Some systems of balance laws models . . . . .	4
1.1.3 Finite volume numerical methods . . . . .	8
1.1.3.1 Reconstruction operators . . . . .	9
1.1.3.2 Relaxation solvers . . . . .	11
1.1.3.3 Time integrators . . . . .	15
1.2 Well-balanced methods . . . . .	17
1.3 Systems in Lagrangian coordinates . . . . .	22
<b>2 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE</b>	<b>25</b>
2.1 The Lagrange-Projection strategy . . . . .	27
2.1.1 The Lagrange-Projection numerical algorithm . . . . .	27
2.2 The Lagrangian step . . . . .	28
2.2.1 Exactly well-balanced reconstruction operators . . . . .	33
2.2.2 Implicit and implicit-explicit exactly well-balanced Lagrangian schemes	35
2.2.2.1 First order schemes . . . . .	36
2.2.2.2 Second order schemes . . . . .	39
2.2.3 Explicit exactly well-balanced Lagrangian schemes . . . . .	41
2.2.3.1 First order explicit Lagrangian scheme . . . . .	41
2.2.3.2 Second order explicit Lagrangian scheme . . . . .	42
2.3 The projection step . . . . .	42
2.3.1 First order projection scheme . . . . .	43
2.3.2 Second order projection scheme . . . . .	44
2.4 Numerical results . . . . .	45
2.4.1 Exactly well-balanced property test . . . . .	46
2.4.2 Computational time vs error . . . . .	46
2.4.3 Order test . . . . .	49
2.4.4 Perturbation of water at rest . . . . .	52
2.4.5 Perturbed water at rest with shock waves . . . . .	54

2.4.6	Generation of subcritical steady state . . . . .	57
<b>3</b>	<b>Implicit LP well-balanced scheme for the Ripa model</b>	<b>61</b>
3.1	Introduction . . . . .	61
3.2	The Lagrangian step . . . . .	63
3.2.1	First order nonlinear implicit well-balanced Lagrangian scheme . . .	65
3.2.2	First order implicit-explicit well-balanced Lagrangian scheme . . . .	66
3.3	The projection step . . . . .	67
3.4	Numerical results . . . . .	67
3.4.1	Isobaric steady state . . . . .	67
3.4.2	Water at rest case . . . . .	69
3.4.3	A general hydrostatic steady state case . . . . .	72
3.4.4	A general steady hydrostatic state with a perturbation . . . . .	76
<b>4</b>	<b>Semi-implicit fully exactly well-balanced schemes for SWE</b>	<b>81</b>
4.1	Introduction . . . . .	81
4.2	Splitting and relaxation techniques . . . . .	82
4.2.1	Well-balanced variable reconstructions . . . . .	88
4.3	First order scheme . . . . .	89
4.3.1	Explicit scheme . . . . .	90
4.3.2	Semi-implicit scheme . . . . .	91
4.4	Second order scheme . . . . .	92
4.4.1	Explicit scheme . . . . .	93
4.4.2	Semi-implicit scheme . . . . .	93
4.5	Numerical experiments . . . . .	93
4.5.1	Fully well-balanced property . . . . .	93
4.5.2	Accuracy test . . . . .	95
4.5.3	Perturbation of water at rest . . . . .	95
4.5.4	Perturbation of water at rest with shock waves . . . . .	98
4.5.5	Perturbation of a subcritical solution . . . . .	98
4.5.6	Perturbation of a transcritical smooth solution . . . . .	100
<b>5</b>	<b>Conclusions and future work</b>	<b>103</b>
	<b>Bibliography</b>	<b>105</b>

# List of Figures

1.1	Sketch of a simple Riemann solver . . . . .	13
2.1	Sketch of the relation between Eulerian and Lagrangian coordinates . . . .	28
2.2	Computational time vs. error for an increasing number of cells using first order schemes . . . . .	48
2.3	Computational time vs. error for an increasing number of cells using second order schemes . . . . .	48
2.4	Free surface for the different first and second order schemes . . . . .	49
2.5	Free surface corresponding to the initial condition for the order test case .	50
2.6	Perturbation of water at rest initial condition . . . . .	53
2.7	Solution for $\eta$ and $q$ at $t = 0.5$ with 200 cells. Explicit: CFL=0.5, Implicit: CFL=2 . . . . .	53
2.8	Solution for $\eta$ and $q$ at $t = 1$ with 200 cells. Explicit: CFL=0.5, Implicit: CFL=2 . . . . .	53
2.9	Solution for $\eta$ and $q$ at $t = 0.5$ with 200 cells. Explicit: CFL=0.5, Implicit: CFL=5 . . . . .	54
2.10	Solution for $\eta$ and $q$ at $t = 1$ with 200 cells. Explicit: CFL=0.5, Implicit: CFL=5 . . . . .	54
2.11	water at rest solution with a discontinuous perturbation on the surface . .	55
2.12	Solution at $t = 0.1$ with 200 cells for a discontinuous perturbation of a water at rest. Explicit: CFL=0.5, Implicit: CFL=2 . . . . .	55
2.13	Solution at $t = 1$ with 200 cells for a discontinuous perturbation of a water at rest. Explicit: CFL=0.5, Implicit: CFL=2 . . . . .	56
2.14	Solution at $t = 0.1$ with 200 cells for a discontinuous perturbation of a water at rest. Explicit: CFL=0.5, Implicit: CFL=5 . . . . .	56
2.15	Solution at $t = 1$ with 200 cells for a discontinuous perturbation of a water at rest. Explicit: CFL=0.5, Implicit: CFL=5. . . . .	56
2.16	Generation of subcritical steady state a time $t = 2$ . Explicit: CFL=0.5, Implicit: CFL=2 . . . . .	57
2.17	Generation of subcritical steady state a time $t = 50$ , Explicit: CFL=0.5, Implicit: CFL=2 . . . . .	58



2.18	Generation of subcritical steady state a time $t = 100$ . Explicit: CFL=0.5, Implicit: CFL=2 . . . . .	58
2.19	Generation of subcritical steady state a time $t = 2$ . Explicit: CFL=0.5, Implicit: CFL=5 . . . . .	58
2.20	Generation of subcritical steady state a time $t = 50$ , Explicit: CFL=0.5, Implicit: CFL=5 . . . . .	59
2.21	Generation of subcritical steady state a time $t = 100$ . Explicit: CFL=0.5, Implicit: CFL=5 . . . . .	59
3.1	Initial condition for an isobaric steady state . . . . .	68
3.2	Free surface computed with well-balanced and non well-balanced schemes for the water at rest steady state . . . . .	70
3.3	Velocity computed with well-balanced and non well-balanced schemes for the water at rest steady state . . . . .	70
3.4	Temperature computed with well-balanced and non well-balanced schemes for the water at rest steady state . . . . .	71
3.5	Exact and discrete solution for $\theta$ for a general hydrostatic steady state case	73
3.6	Error in $h$ for the non well-balanced and exactly well-balanced schemes for a general hydrostatic steady state case . . . . .	74
3.7	Error in $u$ for the non well-balanced and exactly well-balanced schemes for a general hydrostatic steady state case . . . . .	74
3.8	Error in $\theta$ for the non well-balanced and exactly well-balanced schemes for a general hydrostatic steady state case . . . . .	75
3.9	$h$ at different times for the well-balanced scheme for a general hydrostatic steady state case with a perturbation . . . . .	77
3.10	Final steady state for $h$ for the well-balanced scheme for a general hydrostatic steady state case with a perturbation . . . . .	77
3.11	Final steady state for $\theta$ for the well-balanced scheme for a general hydrostatic steady state case with a perturbation . . . . .	78
3.12	Difference for $u$ . . . . .	78
3.13	Difference for $\theta$ . . . . .	79
4.1	Solution at time $t=1$ obtained with the first order schemes using 200 cells .	96
4.2	Solution at time $t=1$ obtained with the second order schemes using 200 cells	97
4.3	Solution for $\eta$ at time $t=1$ obtained using 200 cells when increasing the CFL	97
4.4	Solution at time $t=1$ obtained with the first order schemes using 200 cells .	98
4.5	Solution at time $t=1$ obtained with the second order schemes using 200 cells	99
4.6	Perturbation of a subcritical solution initial condition . . . . .	99
4.7	Difference between the result of the scheme and the steady state at time $t=0.1$ using $N = 200$ cells for the variable $h$ . . . . .	100
4.8	Perturbation of a transcritical smooth solution initial condition . . . . .	101

---

4.9	Difference between the result of the first order schemes and the steady state at time $t = 0.15$ using 200 cells . . . . .	102
4.10	Difference between the result of the second order schemes and the steady state at time $t = 0.15$ using 200 cells . . . . .	102
4.11	Solution for the free surface at time $t = 0.15$ obtained using 200 cells . . .	102





# Resumen

La mecánica de fluidos computacional constituye una valiosa herramienta en la simulación de todo tipo de fenómenos naturales. En efecto, si bien las matemáticas nos permiten describir la evolución de los fluidos mediante sistemas de ecuaciones en derivadas parciales, la resolución exacta de los problemas planteados por lo general no es posible, lo que hace necesaria la utilización de métodos numéricos avanzados para obtener soluciones aproximadas adecuadas.

Los métodos numéricos tienen importantes aplicaciones en campos tan diversos como la oceanografía, la biología, la meteorología, la climatología y la aeronáutica, por citar algunos de ellos. Para abordar el desarrollo de métodos eficaces es esencial contar con un sólido entendimiento, tanto de las características físicas de los fluidos, como de las propiedades matemáticas de los sistemas que los describen.

Una de las líneas de investigación del grupo EDANYA (Ecuaciones Diferenciales, Análisis Numérico Y Aplicaciones) de la Universidad de Málaga, dentro de cuya actividad se ha desarrollado esta tesis doctoral, es la resolución numérica de sistemas de leyes de equilibrio de carácter hiperbólico. Estos sistemas no lineales de ecuaciones en derivadas parciales de evolución, que incluyen un flujo y un término fuente, se utilizan en la simulación de numerosos problemas de la dinámica de fluidos. Es el caso de los modelos de aguas someras, de fluidos multifásicos, de la dinámica de gases, de la magnetohidrodinámica, etc.

En esta tesis se abordan algunos problemas relacionados con la resolución numérica de sistemas hiperbólicos de leyes de equilibrio. En particular, se tratan el sistema constituido por las ecuaciones de aguas someras o aguas poco profundas, también encontradas habitualmente en la literatura como ecuaciones de *shallow water*, y el sistema de Ripa, que se corresponde a una variación del sistema de ecuaciones de aguas someras en que se consideran de forma especial las variaciones de temperatura.

Las ecuaciones de aguas someras se obtienen a partir de las ecuaciones de Navier-Stokes, que se utilizan en mecánica de fluidos para describir el movimiento de un fluido viscoso. Las ecuaciones de Navier-Stokes describen las leyes físicas de conservación de la masa, de la cantidad de movimiento y de la energía, teniendo en cuenta una ecuación de estado que relaciona presión, energía y densidad.

La derivación del sistema unidimensional de ecuaciones de aguas someras que aquí estudiamos se realiza a partir de las de Navier-Stokes mediante un procedimiento de integración vertical en el que se asumen una serie de hipótesis:

- El agua es homogénea e incompresible.
- La presión es hidrostática, lo que implica que aumenta con la profundidad.
- La única fuerza interna que actúa en el fluido es la presión, de modo que se desprecian los efectos viscosos.
- El fondo sobre el que evoluciona el fluido se puede representar mediante una función que depende únicamente de una de las variables horizontales,  $x$ , mientras que la superficie libre lo hace de esa variable horizontal  $x$  y del tiempo  $t$ .
- La velocidad del fluido solo depende igualmente de  $x$  y de  $t$ , despreciándose las variaciones verticales de las componentes horizontales de la velocidad.

El modelo unidimensional así obtenido, que fue deducido por primera vez en 1871 por el ingeniero y matemático francés Adhémar Jean Claude Barré de Saint-Venant (véase [41]), motivo por el que también se conoce bajo el nombre de ecuaciones de Saint-Venant, representa el flujo de una capa delgada de fluido, en que la dimensión horizontal es considerablemente mayor que la vertical.

Consideramos pues el sistema de aguas someras dado por las siguientes ecuaciones:

$$\begin{cases} \partial_t h + \partial_x(hu) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + g\frac{h^2}{2}\right) = -gh\partial_x z, \end{cases}$$

donde  $h(x, t)$  representa el grosor de la capa de agua,  $u(x, t)$  es la velocidad horizontal promediada en la dirección vertical,  $z(x)$  es una función suave conocida, que denota una cierta topografía medida desde un nivel de referencia, y  $g > 0$  es la constante gravitatoria. La superficie libre suele denotarse por  $\eta$  y viene dada por  $\eta = h + z$ . La primera ecuación corresponde a la conservación de masa y la segunda, a la cantidad de movimiento.

En cuanto al modelo de Ripa, su derivación se basa en la consideración de modelos oceánicos multicapa, integrando verticalmente la densidad, el gradiente horizontal de presión y los campos de velocidad en cada capa (véase [79, 80]). Este modelo incorpora el efecto de los gradientes horizontales de temperatura. Las ecuaciones que lo conforman son las siguientes:

$$\begin{cases} \partial_t h + \partial_x(hu) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + \frac{g}{2}h^2\theta\right) = -gh\theta\partial_x z, \\ \partial_t(h\theta) + \partial_x(h\theta u) = 0, \end{cases}$$

siendo  $\theta(x, t)$  el campo de temperatura potencial.

Tanto el sistema de aguas someras como el sistema de Ripa pueden encuadrarse dentro del marco de los sistemas de leyes de equilibrio, que responden a una expresión de la siguiente forma:

$$U_t + f(U)_x = S(U)H_x, \quad (\text{I. 1})$$

en la que  $U$  es una función vectorial que toma valores en un abierto  $\Omega \subset \mathbb{R}^N$ ,  $f : \Omega \rightarrow \mathbb{R}^N$  es una función conocida como la función de flujo y  $S(U)H_x$  es el término fuente. En algunos casos  $H$  puede ser la función identidad. En el caso en el que  $S(U) = 0$  o bien  $H$  sea una función constante, estamos ante lo que se conoce como una ley de conservación:

$$U_t + f(U)_x = 0. \quad (\text{I. 2})$$

Como ya se ha mencionado, en muchos casos resulta imposible resolver los sistemas anteriores de manera exacta, por lo que para realizar simulaciones de flujos gobernados por estas ecuaciones se hace necesario el uso de métodos numéricos adecuados. El objetivo principal de esta tesis es el diseño de métodos numéricos de carácter implícito para estos sistemas. La ventaja de los métodos implícitos con respecto a los explícitos en este caso tiene que ver con la eficiencia computacional en situaciones en las que el número de Froude es bajo. El número de Froude es un número adimensional que relaciona el efecto de las fuerzas de inercia y las fuerzas de gravedad que actúan sobre el fluido. En el caso de aguas someras el número de Froude se define como  $F_r = \frac{|u|}{\sqrt{gh}}$ .

En efecto, los métodos numéricos explícitos presentan mayores restricciones en la elección del paso de tiempo para conseguir estabilidad numérica, mientras que la utilización de métodos implícitos nos permite considerar pasos de tiempo mayores y, con ello, menos iteraciones temporales que cuando se consideran esquemas explícitos.

Asimismo, nos preocuparemos especialmente de que los esquemas desarrollados sean esquemas bien equilibrados o, utilizando la terminología inglesa, esquemas *well-balanced*. Esto es, esquemas que preserven en algún sentido las soluciones de equilibrio, también denominadas estados estacionarios, que son aquellas que no dependen del tiempo, y que por tanto satisfacen

$$f(U)_x = S(U)H_x.$$

La importancia de la propiedad de buen equilibrado radica en el hecho de que, en numerosos escenarios, los flujos que han de simularse surgen como resultado de la perturbación de una solución en estado de equilibrio. Esto ocurre, por ejemplo, en el caso particular de la simulación de tsunamis. Si los errores de discretización del propio método numérico son del mismo orden que la perturbación inicial, sería imposible diferenciar entre las ondas generadas por la perturbación y aquellas provocadas por los errores de discretización. Si bien los errores propios del método podrían reducirse refinando la malla espacial considerada, esto puede llegar a resultar altamente costoso desde el punto de vista computacional. Por este motivo, es crucial que los métodos numéricos considerados sean capaces de preservar de manera precisa las soluciones de equilibrio, lo que permitirá llevar a cabo una simulación adecuada en situaciones como las ya descritas.

En el caso en el que un método numérico preserva una cierta familia de estados estacionarios diremos que es *well-balanced*, mientras que si preserva todos los posibles estados estacionarios lo denominaremos *fully well-balanced*.

En el ámbito particular de los modelos de aguas someras, fue en el trabajo de Bermúdez y Vázquez-Cendón [4] donde se introdujo por primera vez el desarrollo de métodos bien

equilibrados, así como la propiedad C, consistente en preservar las soluciones estacionarias de velocidad nula. Además de este, existen multitud de trabajos en los que se ha estudiado esta propiedad y se han diseñado esquemas de este tipo. Algunos de ellos son, por ejemplo, [61, 77, 51, 50, 7, 1, 38, 2, 9, 45, 49], en los que se preservan los estados estacionarios del agua en reposo. También existen trabajos en los que se presentan esquemas que preservan todos los estados estacionarios, como son [57, 20, 93, 92, 5, 6, 24, 74, 81, 8].

Para lograr que los métodos numéricos desarrollados en este trabajo sean bien equilibrados se utiliza la estrategia descrita en [25], que ha sido empleada en distintos modelos para desarrollar esquemas que preservan todas las soluciones estacionarias de sistemas de leyes de equilibrio, cuando se cuenta con una representación explícita o implícita de estas soluciones. Algunos ejemplos son el modelo de aguas someras (ver [24]), el modelo de flujo de sangre en los vasos sanguíneos (ver [73]), el modelo de Ripa (ver [82]) o el modelo de Euler con gravedad (ver [48, 65]).

El diseño de métodos numéricos implícitos bien equilibrados se aborda en este trabajo utilizando dos estrategias diferentes. Por un lado, en los Capítulos 2 y 3 aplicaremos la estrategia Lagrangiano-Proyectado, ya utilizada en otros trabajos como [23, 72, 33], que consiste, en cada iteración temporal, en resolver primero el sistema en coordenadas Lagrangianas, para proyectar a continuación la solución así obtenida en coordenadas Eulerianas. Por otro lado, en el Capítulo 4 aplicaremos técnicas de *splitting* y de relajación, lo que resultará en la resolución de dos sistemas, en lugar de uno, en cada paso de tiempo. Ambas estrategias nos permiten desacoplar los fenómenos acústicos y de transporte presentes en nuestras ecuaciones, así como diseñar de forma natural esquemas implícito-explícitos. Esto resulta especialmente útil en la aproximación de flujos subsónicos o de número de Froude pequeño, en que la restricción CFL habitual debida a las ondas acústicas conduce a pasos de tiempo muy pequeños. El uso de esquemas implícitos para el sistema de presión hace que la restricción CFL se reduzca al paso de transporte y se evite en el paso acústico, más restrictivo.

En general, el procedimiento que seguiremos comenzará por considerar la discretización espacial mediante el método de volúmenes finitos del sistema estudiado, particionando el dominio del problema en una serie de celdas computacionales que conformarán nuestra malla. Tras esto obtendremos un esquema semi-discreto en tiempo, que constituye un sistema de ecuaciones diferenciales ordinarias. Finalmente, se aplicará un resolvidor en tiempo a ese sistema de EDOs. Como se verá en el Capítulo 1, en función del resolvidor que se considere se obtendrá un tipo de esquema u otro: explícito, implícito o semi-implícito.

Esta tesis se apoya en las siguientes publicaciones:

- C. Caballero-Cárdenas, M.J. Castro, T. Morales de Luna, and M.L. Muñoz-Ruiz. Implicit and implicit-explicit Lagrange-projection finite volume schemes exactly well-balanced for 1d shallow water system. *Applied Mathematics and Computation*, 443:127784, 2023.
- C. Caballero-Cárdenas, M.J. Castro, T. Morales de Luna, and M.L. Muñoz-Ruiz.

*Lagrange-projection exactly well-balanced finite volume schemes for the Ripa model.* SEMA/SMAI Springer series book, aceptado.

- C. Caballero-Cárdenas, M.J. Castro, C. Chalons, T. Morales de Luna, and M.L. Muñoz-Ruiz. A semi-implicit fully well-balanced relaxation scheme for shallow water system. Enviado para su revisión a una revista de alto impacto.

C. Caballero-Cárdenas ha disfrutado del contrato predoctoral FPI2019/087773 financiado por MCIN/AEI/10.13039/501100011033 y por FSE invierte en tu futuro.

A continuación se describe la organización en capítulos de esta tesis, haciéndose un breve resumen de sus contenidos.

## Capítulo 1: Preliminares

Este capítulo inicial tiene como objetivo establecer las bases teóricas en que se asienta la memoria.

En la primera sección se realiza un breve repaso de los métodos numéricos de tipo volúmenes finitos para sistemas unidimensionales de leyes de equilibrio.

Se definen las leyes de equilibrio (I. 1) y las leyes de conservación (I. 2), y se introduce el concepto de sistema hiperbólico en el marco de los sistemas de leyes de conservación, así como el de campo característico linealmente degenerado y el de campo característico genuinamente no lineal. Se define también el problema de Cauchy para una ley de conservación, esto es, el problema de valor inicial dado por

$$\begin{cases} U_t + f(U)_x = 0, \\ U(x, 0) = U_0(x). \end{cases} \quad (\text{I. 3})$$

De una función suficientemente regular  $U : \mathbb{R} \times [0, \infty) \rightarrow \Omega$  que satisfaga (I. 3) diremos que es una solución clásica del problema de Cauchy. Sin embargo, nos enfrentamos a una dificultad, y es que se sabe que aunque la condición inicial  $U_0$  sea suave, el problema de Cauchy anterior puede no tener solución en el sentido clásico. Es por ello que necesitamos introducir el concepto de solución débil, para lo que se hace uso de la formulación variacional del problema. Junto con esto se aportarán una serie de resultados teóricos relacionados con las soluciones débiles de un problema (I. 3), como el establecimiento de la condición de Rankine-Hugoniot para soluciones que presenten discontinuidades de tipo salto, que nos aporta información acerca de la velocidad de propagación de la discontinuidad. Otro de los problemas a los que nos enfrentamos al resolver (I. 3) es que en general la solución débil no es necesariamente única, de modo que para seleccionar de entre todas las posibles soluciones aquella con sentido físico necesitaremos hacer uso de una condición de entropía.

Los resultados anteriores se extenderán también al caso de sistemas de leyes de equilibrio, de los que se presentan algunos ejemplos, como el de shallow water y el modelo de Ripa que nos ocupan en esta memoria.

Finalmente, se describe en qué consiste un método numérico de volúmenes finitos. Para ello el espacio se discretiza usando un conjunto de celdas  $[x_{i-1/2}, x_{i+1/2})$  con un paso de malla  $\Delta x$ , siendo  $x_{i+1/2} = i\Delta x$  y  $x_i = (x_{i-1/2} + x_{i+1/2})/2$  las interceldas y los centros de las celdas, respectivamente, para  $i \in \mathbb{Z}$ . De la misma manera, el paso temporal se denota por  $t_n = n\Delta t$ , con  $n \in \mathbb{N}$ . Comenzamos considerando leyes de conservación y finalmente comentamos cómo se procedería en el caso de leyes de equilibrio. Para leyes de conservación, se comienza integrando (I. 2) en cada celda  $I_i$ , y denotando por  $U_i(t)$  la aproximación del promedio de la solución en la celda  $I_i$  y en el tiempo  $t$ ,

$$U_i(t) = \int_{x_{i-1/2}}^{x_{i+1/2}} U(x, t) dx,$$

podemos escribir un método numérico semi-discreto de la siguiente manera:

$$\frac{dU_i(t)}{dt} = -\frac{1}{\Delta x} (F_{i+1/2}^t - F_{i-1/2}^t), \quad (\text{I. 4})$$

siendo  $F_{i+1/2}^t$  una aproximación del promedio del flujo en  $x_{i+1/2}$ .

A continuación se introducen tres herramientas que serán clave en este trabajo: los operadores de reconstrucción, los resolvedores de relajación y los integradores en tiempo.

En el apartado destinado a los operadores de reconstrucción se introduce su definición y se utilizan para obtener el flujo  $F_{i+1/2}^t$  en (I. 4) de la siguiente manera:

$$F_{i+1/2}^t = \mathbb{F}(U_{i+1/2-}^t, U_{i+1/2+}^t), \quad (\text{I. 5})$$

siendo  $\mathbb{F}$  el flujo numérico, evaluado en  $U_{i+1/2-}^t$  y  $U_{i+1/2+}^t$ , los estados reconstruidos en las interceldas. Se describen además los operadores de reconstrucción que se utilizarán, que son el operador constante para primer orden y el operador MUSCL en el caso del segundo orden, sin olvidar el uso de los limitadores correspondientes para evitar oscilaciones en el caso de presencia de discontinuidades. El uso de operadores de reconstrucción puede ampliarse al caso de leyes de equilibrio, escribiendo el método numérico semi-discreto como

$$\frac{dU_i(t)}{dt} = -\frac{1}{\Delta x} (F_{i+1/2}^t - F_{i-1/2}^t) + \frac{1}{\Delta x} S_i^t, \quad (\text{I. 6})$$

donde el flujo numérico viene dado por (I. 5) y  $S_i^t$  corresponde a una aproximación de la integral del término fuente en la celda  $i$ , es decir:

$$S_i^t \approx \int_{x_{i-1/2}}^{x_{i+1/2}} S(P_i^t(x)) H_x dx, \quad (\text{I. 7})$$

siendo  $P_i^t$  un operador de reconstrucción.

En el apartado de resolvedores de relajación se describe en qué consisten estos métodos, utilizados en trabajos como [63, 39, 3, 18, 32]. La idea de los esquemas de relajación

consiste en considerar, para un sistema (I. 2) dado, otro sistema en una dimensión mayor que el original, de modo que sea posible relacionar ambos mediante un operador lineal tal que la solución del sistema relajado sea una aproximación de la solución del original, siempre que se cumplan unas ciertas condiciones. Algunos esquemas de relajación bien conocidos son el esquema HLL o el flujo de Rusanov. En este caso, nos centramos en los sistemas de relajación de tipo Suliciu, considerados en [86, 87, 39, 18, 32, 3], por mencionar algunos ejemplos. Para introducir el sistema de relajación de Suliciu, se consideran las ecuaciones isentrópicas de dinámica de gases, dadas por

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + p(\rho)) = 0, \end{cases} \quad (\text{I. 8})$$

donde  $\rho$  es la densidad,  $u$  la velocidad y  $p(\rho)$  la presión. Tras realizar unos sencillos cálculos que se especifican en el Capítulo 1 se llega al sistema de relajación

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + \pi) = 0, \\ \partial_t(\rho \pi) + \partial_x(\rho \pi u) + a^2 \partial_x u = 0, \end{cases} \quad (\text{I. 9})$$

siendo  $\pi = p(\rho)$  y  $a$  una constante que debe satisfacer la condición subcaracterística, esto es:

$$|\lambda_j(U)| \leq a,$$

con  $\lambda_j(U)$  los autovalores del sistema original. La ventaja de resolver el sistema (I. 9) en lugar de (I. 8) radica en la mayor facilidad de resolver el problema de Riemann para (I. 9), al tener este último sistema todos los campos linealmente degenerados.

En el apartado correspondiente a los integradores temporales éstos se aplican para aproximar la solución del sistema de EDOs (I. 6). Introduciremos integradores de tipo explícito, implícito y semi-implícito. En cada caso, consideraremos integradores de primer y segundo orden. También presentaremos el conocido como *Strang splitting*, que será utilizado en el Capítulo 4 para obtener esquemas de segundo orden.

En la segunda sección del capítulo introductorio nos centramos en el concepto de métodos bien equilibrados, haciendo hincapié en la importancia de esta propiedad y describiendo con detalle la estrategia que se seguirá para obtenerlos, basada en la idea que se presenta en [25].

Esta estrategia consiste en la utilización de un operador de reconstrucción que sea bien equilibrado, esto es, que cuando se aplica a los promedios de las celdas de una solución estacionaria, las aproximaciones obtenidas deben coincidir con dicha solución estacionaria. En general, los operadores de reconstrucción no tienen la propiedad anterior, pero construiremos uno que sí la verifique a partir del proceso siguiente:

1. Se busca, si es posible, una solución estacionaria de forma que su promedio en la celda coincida con el valor de nuestra aproximación en dicha celda.

2. Se calculan las fluctuaciones en las celdas del *stencil* de un operador de reconstrucción estándar de nuestra elección, esto es, las diferencias entre los promedios de la solución y los de la solución estacionaria obtenida en el paso anterior. Se aplica a las fluctuaciones el operador de reconstrucción estándar elegido.
3. Se suman la solución estacionaria obtenida en el primer paso y la reconstrucción del segundo para obtener el operador de reconstrucción bien equilibrado.

El operador así obtenido es bien equilibrado y del mismo orden que el operador de reconstrucción estándar utilizado en el segundo paso, supuesto que las soluciones estacionarias son suficientemente regulares y que el operador de reconstrucción estándar es exacto en el caso de funciones nulas.

Asimismo, hay que adaptar la escritura de la integral del término fuente de manera que el esquema final que se obtiene sea bien equilibrado.

La última sección se dedica a introducir las coordenadas Lagrangianas, que serán utilizadas en los Capítulos 2 y 3. Dado que estas coordenadas siguen las trayectorias de las partículas del flujo, se considera una partícula  $\xi$  y se definen las curvas características

$$\begin{cases} \frac{\partial x}{\partial t}(\xi, t) = u(x(\xi, t), t), \\ x(\xi, 0) = \xi. \end{cases}$$

Además, dada una función en coordenadas Eulerianas  $(x, t) \mapsto \mathbf{U}(x, t)$ , se define su equivalente en coordenadas Lagrangianas como

$$\bar{\mathbf{U}}(\xi, t) = \mathbf{U}(x(\xi, t), t).$$

Por último, se define también el Jacobiano de la aplicación Lagrangiana

$$L(\xi, t) = \frac{\partial x}{\partial \xi}(\xi, t),$$

y se derivan una serie de propiedades que relacionan las derivadas parciales en las distintas coordenadas. En esta sección se considera el caso de las ecuaciones de Euler, que se reescriben en términos de las coordenadas Lagrangianas, y se particulariza finalmente en el caso de las ecuaciones de aguas someras.

## Capítulo 2: Esquemas de volúmenes finitos Lagrangiano-Projectados implícitos e implícitos-explicitos exactamente bien equilibrados para las ecuaciones de aguas someras unidimensionales

En este capítulo se presenta el trabajo publicado en [22], en el que consideramos la técnica Lagrangiano-Projectado en el marco de esquemas de volúmenes finitos aplicados al sistema de aguas someras. Como ya se ha mencionado, esta técnica consta de dos pasos: en

primer lugar, se resuelve el sistema en coordenadas Lagrangianas y, tras ello, el resultado se proyecta a coordenadas Eulerianas. El primer paso se conoce como paso Lagrangiano y el segundo, como paso de proyección o de transporte. El objetivo de volver a coordenadas Eulerianas en cada paso es que el uso de coordenadas Lagrangianas puras y el seguimiento de mallas en movimiento puede ser engorroso, y se pueden dar situaciones complejas para la configuración de las celdas en movimiento, especialmente pensando en la extensión 2D. Además, de este modo podemos desacoplar los fenómenos acústicos y de transporte, y así diseñar esquemas implícito-explicitos y tomar pasos de tiempo grandes de manera natural. Esto resulta interesante sobre todo en el caso de número de Froude pequeño, ya que al utilizar esquemas implícitos o implícito-explicitos, la restricción CFL se reduce únicamente a las ondas de transporte lentas en lugar de a las acústicas, que son más restrictivas.

Se consideran dos versiones del esquema para el paso Lagrangiano: una implícita no lineal y otra implícita-explicita, basada en cómo se trata el término fuente geométrico. El paso de transporte siempre se realiza de manera explícita. Se presentan versiones exactamente bien equilibradas de primer y segundo orden de los esquemas, en las que se preservan las soluciones de agua en reposo. Para ello, se consideran operadores de reconstrucción exactamente bien equilibrados construidos utilizando la estrategia descrita en [25].

Tras llevar a cabo distintos tests numéricos, se concluye que el esquema en el que el paso Lagrangiano se resuelve de manera implícita-explicita es más eficiente que el que es implícito no lineal, ya que en el último caso hay que resolver un sistema no lineal y son necesarias varias iteraciones de un algoritmo de punto fijo. Además, como es de esperar, los esquemas de primer orden son más difusivos que los de segundo orden. Por último, se observa una mejora en la eficiencia de los esquemas en los que el paso Lagrangiano se efectúa de manera implícita (tanto no lineal como implícita-explicita) con respecto a aquellos en los que se realiza de manera explícita, en el caso de bajo número de Froude, ya que en el caso implícito se pueden tomar pasos de tiempo mucho mayores.

### Capítulo 3: Esquemas de volúmenes finitos Lagrangiano-Proyectados implícitos bien equilibrados para el modelo de Ripa

En el tercer capítulo proponemos una estrategia para resolver numéricamente el modelo de Ripa aplicando, de nuevo, la técnica Lagrangiano-Proyectado. Esto es, en primer lugar se resuelve el sistema de Ripa en coordenadas Lagrangianas y a continuación se proyecta la solución en coordenadas Eulerianas. Al igual que en el Capítulo 2, el sistema Lagrangiano se resuelve implícitamente mientras que el paso de proyección se hace de manera explícita.

Se diseñan esquemas de volúmenes finitos de primer orden que son bien equilibrados para este modelo, preservando los conocidos como estados estacionarios hidrostáticos, que son los que corresponden a  $u = 0$ , y que satisfacen

$$\partial_x p = -2p \frac{\partial_x z}{h}.$$

Casos particulares de estos estados estacionarios son los que corresponden al agua en reposo,

$$u = 0, \quad h + z = \text{constant}, \quad \theta = \text{constant},$$

así como los estados estacionarios isobáricos, en los que la presión es constante y el fondo es plano:

$$u = 0, \quad \frac{g}{2}h^2\theta = \text{constant}, \quad z = \text{constant}.$$

Puesto que en este caso las soluciones estacionarias quedan definidas fijando un perfil para  $h$  o  $\theta$ , resulta imposible diseñar esquemas numéricos exactamente bien equilibrados para todas las soluciones hidrostáticas. Por eso hemos optado por usar una técnica de aproximación, de forma que una vez elegido un perfil discreto para  $h$ , procedemos a la aproximación de la presión utilizando para ello un método de colocación (ver [56]) para aproximar las soluciones de la EDO

$$\partial_\xi p = -2p \frac{\partial_\xi z}{h}.$$

Una vez recuperada la presión, y conocido el perfil discreto de  $h$  es posible recuperar el perfil discreto estacionario de  $\theta$ . Los esquemas que se obtienen son, por tanto, bien equilibrados pero no exactamente bien equilibrados.

Se incluye una sección de tests numéricos en la que se compara el esquema propuesto con otros dos esquemas de tipo Lagrangiano-Proyectado: uno no bien equilibrado y otro exactamente bien equilibrado para los estados estacionarios de agua en reposo y los estados estacionarios isobáricos, pero no para cualquier estado estacionario hidrostático.

#### Capítulo 4: Esquemas de volúmenes finitos exactamente bien equilibrados para las ecuaciones de aguas someras unidimensionales que preservan todas las soluciones estacionarias

En este capítulo, el objetivo es el diseño de esquemas de primer y segundo orden implícitos que preserven todas las soluciones estacionarias de las ecuaciones de aguas someras. Para ello, se aplica una técnica de *splitting* mediante la cual se obtienen dos sistemas a resolver: un sistema de presión y otro de transporte. A continuación se considera un sistema relajado del sistema de presión. Este se resolverá de manera implícita, mientras que el otro, el de transporte, se resolverá explícitamente. Es posible resolver primero el sistema de presión seguido del de transporte o viceversa.

En el caso de segundo orden, se utiliza el conocido como *Strang splitting* junto con reconstrucciones de segundo orden en espacio y primer orden en tiempo. Este *splitting* consiste en efectuar un paso de uno de los sistemas con paso de tiempo  $\Delta t/2$ , seguido de un paso del segundo sistema con paso de tiempo  $\Delta t$  y finalizando con un paso del primer sistema con paso de tiempo  $\Delta t/2$ . Por tanto, en el caso en el que se comience por el sistema de presión, serán necesarios dos pasos implícitos mientras que en el que se

comience por el sistema de transporte únicamente el paso segundo será implícito, por lo que cabe esperar que este último sea más eficiente.

De esta manera, como en los capítulos anteriores, se obtienen esquemas con los que es posible realizar pasos de tiempo mayores que en el caso de un esquema explícito, y que por tanto, resultan más eficientes que los explícitos en el caso de número de Froude bajo.

Al estudiar los resultados de los esquemas propuestos con varios tests y aumentar el valor del CFL, se han observado distintos comportamientos en función de qué sistema se resolvía primero. En el caso del primer orden, se observa un mejor comportamiento en el caso del esquema en el que se resuelve primero el sistema de presión, mientras que en el segundo orden, la estabilidad es mejor en el esquema que comienza con el sistema de transporte.

Además, se han efectuado tests en los que se consideran perturbaciones de una solución subcrítica, así como de una solución transcítica suave, observándose en ambos los casos buenos resultados.

## **Capítulo 5: Conclusiones y trabajo futuro**

En esta sección final del documento se presentan las conclusiones derivadas de los resultados obtenidos a lo largo de esta tesis.

En cuanto a los trabajos futuros, se plantean varios aspectos en los que se trabajará, con el objetivo de aplicar lo desarrollado en esta tesis a otros problemas, como la extensión al caso bidimensional, el diseño de esquemas de orden mayor que dos o la aplicación de las estrategias estudiadas a otros sistemas.



# Introduction

Computational fluid mechanics is a valuable tool in the simulation of all kinds of natural phenomena. Indeed, although mathematics allows us to describe the evolution of fluids by means of systems of partial differential equations, the exact resolution of the problems posed is generally not possible, which makes it necessary to use advanced numerical methods to obtain suitable approximate solutions.

Numerical methods have important applications in fields as diverse as oceanography, biology, meteorology, climatology and aeronautics, to name but a few. A solid understanding of both the physical characteristics of fluids and the mathematical properties of the systems that describe them is essential for the development of effective methods.

One of the research lines of the EDANYA group (Ecuaciones Diferenciales, Análisis Numérico Y Aplicaciones) of the University of Málaga, within whose activity this PhD thesis has been developed, is the numerical resolution of systems of hyperbolic balance laws. These non-linear systems of evolution partial differential equations, which include a flux and a source term, are used in the simulation of numerous fluid dynamics problems. This is the case for shallow water models, multiphase fluid models, gas dynamics, magnetohydrodynamics, etc.

This thesis deals with some problems related to the numerical resolution of hyperbolic systems of balance laws. In particular, it deals with the shallow water equations system and the Ripa system, which corresponds to a variation of the system of shallow water equations in which temperature variations are considered in a special way.

The shallow water equations are derived from the Navier-Stokes equations, which are used in fluid mechanics to describe the motion of a viscous fluid. The Navier-Stokes equations describe the physical laws of conservation of mass, quantity of motion and energy, taking into account an equation of state relating pressure, energy and density.

The derivation of the one-dimensional system of shallow water equations studied here is carried out from the Navier-Stokes equations by means of a vertical integration procedure in which a number of assumptions are made:

- The water is homogeneous and incompressible.
- The pressure is hydrostatic, which implies that it increases linearly with depth.
- The only internal force acting on the fluid is pressure, so viscous effects are neglected.

- The bottom over which the fluid flows can be represented by a function that depends only on one of the horizontal variables,  $x$ , while the free surface depends on that horizontal variable  $x$  and on time  $t$ .
- The velocity of the fluid only depends on  $x$  and  $t$ , neglecting the vertical variations of the horizontal components of the velocity.

The one-dimensional model thus obtained, which was first derived in 1871 by the French engineer and mathematician Adhémar Jean Claude Barré de Saint-Venant (see [41]), which is why it is also known as the Saint-Venant system, represents the flow of a thin layer of fluid in which the horizontal dimension is considerably larger than the vertical one.

We therefore consider the shallow water system given by the following equations:

$$\begin{cases} \partial_t h + \partial_x(hu) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + g\frac{h^2}{2}\right) = -gh\partial_x z, \end{cases}$$

where  $h(x, t)$  represents the thickness of the water layer,  $u(x, t)$  is the horizontal velocity averaged in the vertical direction,  $z(x)$  is a known smooth function, denoting a certain topography measured from a reference level, and  $g > 0$  is the gravitational constant. The free surface is usually denoted by  $\eta$  and is given by  $\eta = h + z$ . The first equation corresponds to the conservation of mass and the second to the quantity of motion.

As for Ripa's model, its derivation is based on the consideration of multilayer ocean models, vertically integrating the density, the horizontal pressure gradient and the velocity fields in each layer (see [79, 80]). This model incorporates the effect of horizontal temperature gradients. The equations are the following:

$$\begin{cases} \partial_t h + \partial_x(hu) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + \frac{g}{2}h^2\theta\right) = -gh\theta\partial_x z, \\ \partial_t(h\theta) + \partial_x(h\theta u) = 0, \end{cases}$$

where  $\theta(x, t)$  is the potential temperature field.

Both the shallow water system and the Ripa system can be placed within the framework of balance law systems, which respond to an expression of the following form:

$$U_t + f(U)_x = S(U)H_x, \quad (\text{I. 10})$$

where  $U$  is a vector function that takes values in an open  $\Omega \subset \mathbb{R}^N$ ,  $f : \Omega \rightarrow \mathbb{R}^N$  is a function known as the flux function and  $S(U)H_x$  is the source term. Note that in certain situations  $H$  could be the identity function. In the case where  $S(U) = 0$  or  $H$  is a constant function, we have what is known as a conservation law:

$$U_t + f(U)_x = 0. \quad (\text{I. 11})$$

As already mentioned, the above systems cannot be solved exactly, so in order to perform simulations of flows governed by these equations it is necessary to use appropriate numerical methods. The main objective of this thesis is the design of implicit numerical methods for these systems. The advantage of implicit methods over explicit ones in this case has to do with computational efficiency in situations where the Froude number is low. The Froude number is a dimensionless number that relates the effect of inertial forces and gravity forces acting on the fluid. In the case of shallow water the Froude number is defined as  $F_r = \frac{|u|}{\sqrt{gh}}$ . One can find numerous works concerning the design of semi-implicit or IMEX schemes, such as, [67, 10, 70, 71, 44, 15, 16, 17].

Indeed, explicit numerical methods present greater restrictions in the choice of time step to achieve numerical stability, while the use of implicit methods allows us to consider larger time steps and thus fewer time iterations than when explicit schemes are considered.

We will also be particularly concerned that the developed schemes are well-balanced. That is, schemes that preserve in some sense the equilibrium solutions, also called steady states, which are those that do not depend on time, and which therefore satisfy

$$f(U)_x = S(U)H_x.$$

The importance of the well-balanced property lies in the fact that, in many scenarios, the flows to be simulated arise as a result of perturbation of an equilibrium solution. This is the case, for example, in the particular case of tsunami simulation. If the discretisation errors of the numerical method itself are of the same order as the initial perturbation, it would be impossible to distinguish between the waves generated by the perturbation and those caused by the discretisation errors. While the method's own errors could be reduced by refining the spatial grid under consideration, this can be computationally expensive. For this reason, it is crucial that the numerical methods considered are capable of accurately preserving the equilibrium solutions, which will allow an adequate simulation to be carried out in situations such as those described above.

In the case in which a numerical method preserves a certain family of steady states we will say that it is *well-balanced*, while if it preserves every possible steady state we will call it *fully well-balanced*.

In the particular field of shallow water models, it was in the work of Bermúdez and Vázquez-Cendón [4] that the development of well-balanced methods was first introduced, as well as the C-property, which consists of preserving stationary solutions of zero velocity. In addition to this, there are many other works in which this property has been studied and schemes of this type have been designed. Some of them are, for example, [61, 77, 51, 50, 7, 1, 38, 2, 9, 45, 49], in which the steady states of water at rest are preserved. There are also papers that present schemes preserving every steady state, such as [57, 20, 93, 92, 5, 6, 24, 74, 81, 8].

The main objective of this thesis will be the design of different well-balanced semi-implicit finite volume schemes for the shallow water equations and for the Ripa system.



# Abstract

This thesis addresses the problem of the design of first and second order finite volume semi-implicit well-balanced schemes for the shallow water equations as well as for the Ripa system.

To ensure the well-balanced character of the numerical methods developed, the approach outlined in [25] has been adopted. This strategy, employed in various models, enables the creation of schemes that preserve all stationary solutions of systems of balance laws when an explicit or implicit expression of these solutions is available. Some examples are the shallow water model (see [24]), the blood flow model in blood vessels (see [73]), the Ripa model (see [82]) or the Euler model with gravity (see [48, 65]).

The design of well-balanced implicit numerical methods is addressed in this work using two different strategies. On the one hand, in Chapters 2 and 3 we will apply the Lagrange-Projection strategy, already used in other works such as [23, 72, 33], which consists of, in each time iteration, first solving the system in Lagrangian coordinates and then projecting the solution thus obtained in Eulerian coordinates. On the other hand, in Chapter 4 we will apply splitting and relaxation techniques, which will result in solving two systems, instead of one, at each time step. Both strategies allow us to decouple the acoustic and transport phenomena present in our equations, as well as to design implicit-explicit schemes in a natural way. This is especially useful in the approximation of subsonic or small Froude number flows, where the usual CFL constraint due to acoustic waves leads to very small time steps. The use of implicit schemes for the pressure system means that the CFL constraint is reduced to the transport step and avoided in the more restrictive acoustic step.

In general, the procedure we will follow will begin by considering the spatial discretisation by means of the finite volume method of the system under study and partitioning the problem domain into a series of computational cells that will conform our grid. After this, we will obtain a semi-discrete scheme in time, which constitutes a system of ordinary differential equations. Finally, a time integrator will be applied to the system of ODEs. As will be seen in Chapter 1, depending on the integrator considered, a different type of scheme will be obtained: explicit, implicit or semi-implicit.

This thesis is supported by the following publications:

- C. Caballero-Cárdenas, M. J. Castro, T. Morales de Luna, and M. L. Muñoz-Ruiz. Implicit and implicit-explicit Lagrange-projection finite volume schemes exactly well- balanced for 1d shallow water system. *Applied Mathematics and Computation*, 443:127784, 2023.
- C. Caballero-Cárdenas, M. J. Castro, T. Morales de Luna, and M. L. Muñoz-Ruiz. *Lagrange-projection exactly well-balanced finite volume schemes for the Ripa model*. SEMA/SMAI Springer series book, accepted for publication.
- C. Caballero-Cárdenas, M. J. Castro, C. Chalons, T. Morales de Luna, and M. L. Muñoz-Ruiz. A semi-implicit fully well-balanced relaxation scheme for shallow water system. Sent for review to a high impact journal.

C. Caballero-Cárdenas is supported by the grant FPI2019/087773 funded by MCIN/AEI/10.13039/501100011033 and “ESF Investing in your future”.

The organisation of this thesis into chapters is described below, with a brief summary of its contents.

## Chapter 1: Mathematical settings

The aim of this initial chapter is to establish the theoretical foundations on which the thesis is based.

In the first section, a brief review of finite volume numerical methods for one-dimensional systems of balance laws is given. The balance laws (I. 10) and the conservation laws (I. 11) are defined, and the concept of hyperbolic system in the framework of conservation law systems is introduced, as well as that of linearly degenerate characteristic field and that of genuinely nonlinear characteristic field. The Cauchy problem for a conservation law is also defined, i.e. the initial value problem given by

$$\begin{cases} U_t + f(U)_x = 0, \\ U(x, 0) = U^0(x). \end{cases} \quad (\text{I. 12})$$

Given a sufficiently smooth function  $U : \mathbb{R} \times [0, \infty) \rightarrow \Omega$  satisfying (I. 12), we will say that it is a classical solution of the Cauchy problem. However, we face a difficulty, since it is known that even if the initial condition  $U^0$  is smooth, the above Cauchy problem may not have a solution in the classical sense. This is why we need to introduce the concept of weak solution, for which we make use of the variational formulation of the problem. Along with this, a series of theoretical results related to the weak solutions of a problem (I. 12) will be provided, such as the establishment of the Rankine-Hugoniot condition for solutions with jump discontinuities, which gives us information about the speed of propagation of the discontinuity. Another problem we face when solving (I. 12) is that in

general the weak solution is not necessarily unique, so in order to select from all possible solutions the one with physical sense we will need to make use of an entropy condition.

The previous results will also be extended to the case of balance law systems, of which some examples are presented, such as the shallow water and the Ripa model that we are dealing with in this report.

Finally, a finite volume numerical method is described. For this, the space is discretised using a set of cells  $[x_{i-1/2}, x_{i+1/2})$  with a mesh step  $\Delta x$ , where  $x_{i+1/2} = i\Delta x$  and  $x_i = (x_{i-1/2} + x_{i+1/2})/2$  are the intercells and cell centres, respectively, for  $i \in \mathbb{Z}$ . In the same way, the time step is denoted by  $t_n = n\Delta t$ , with  $n \in \mathbb{N}$ . We begin by considering conservation laws and finally discuss how to proceed in the case of balance laws. For conservation laws, we start by integrating (I. 11) in each cell  $I_i$ , and denote by  $U_i(t)$  the approximation of the average of the solution in cell  $I_i$  and time  $t$ ,

$$U_i(t) = \int_{x_{i-1/2}}^{x_{i+1/2}} U(x, t) dx.$$

We can then write a semi-discrete numerical method as follows:

$$\frac{dU_i(t)}{dt} = -\frac{1}{\Delta x} (F_{i+1/2}^t - F_{i-1/2}^t), \quad (\text{I. 13})$$

where  $F_{i+1/2}^t$  is an approximation of the average flux at  $x_{i+1/2}$ .

We then introduce three tools that will be key in this work: reconstruction operators, relaxation solvers and time integrators.

In the subsection dedicated to reconstruction operators, their definition is introduced and used to obtain the flux  $F_{i+1/2}^t$  in (I. 13) as follows:

$$F_{i+1/2}^t = \mathbb{F}(U_{i+1/2-}^t, U_{i+1/2+}^t), \quad (\text{I. 14})$$

where  $\mathbb{F}$  is the numerical flux, evaluated at  $U_{i+1/2-}^t$  and  $U_{i+1/2+}^t$ , the reconstructed states in the intercells. The reconstruction operators to be used are described, including the constant operator for first order and the MUSCL operator for second order, with the incorporation of appropriate limiters to prevent the appearance of spurious oscillations in the presence of discontinuities. The use of reconstruction operators can be extended to balance laws, expressing the semi-discrete numerical method as

$$\frac{dU_i(t)}{dt} = -\frac{1}{\Delta x} (F_{i+1/2}^t - F_{i-1/2}^t) + \frac{1}{\Delta x} S_i^t, \quad (\text{I. 15})$$

where the numerical flux is given by (I. 14), and  $S_i^t$  corresponds to an approximation of the integral of the source term in cell  $i$ , that is:

$$S_i^t \approx \int_{x_{i-1/2}}^{x_{i+1/2}} S(P_i^t(x)) H_x dx, \quad (\text{I. 16})$$

with  $P_i^t$  being a reconstruction operator.

In the subsection on relaxation solvers, we describe what these methods, used in works such as [63, 39, 3, 18, 32], consist of. The idea of relaxation schemes consists in considering, for a given system (I. 11), another system in a higher dimension than the original one, so that it is possible to relate the two by means of a linear operator such that the solution of the relaxed system is an approximation of the solution of the original one, provided that certain conditions are met. Some well-known relaxation operators are the HLL scheme or the Rusanov flow. In this case, we focus on Suliciu-type relaxation systems, considered in [86, 87, 39, 18, 32, 3], to mention some examples. To introduce the Suliciu relaxation system, we consider the isentropic equations of gas dynamics, given by

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + p(\rho)) = 0, \end{cases} \quad (\text{I. 17})$$

where  $\rho$  is the density,  $u$  the velocity and  $p(\rho)$  the pressure. After some simple calculations specified in Chapter 1, we arrive at the relaxation system

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + \pi) = 0, \\ \partial_t(\rho \pi) + \partial_x(\rho \pi u) + a^2 \partial_x u = 0, \end{cases} \quad (\text{I. 18})$$

where  $\pi = p(\rho)$  and  $a$  is a constant which must satisfy the subcharacteristic condition, i.e:

$$|\lambda_j(U)| \leq a,$$

with  $\lambda_j(U)$  representing the eigenvalues of the original system, solving the system (I. 18) instead of (I. 17) offers the advantage of easier resolution of the Riemann problem for (I. 18). This arises from the fact that all the fields become linearly degenerate in the relaxed system.

In the subsection dedicated to time integrators, they are applied to approximate the solution of the system of ODEs (I. 15). Explicit, implicit, and semi-implicit integrators are introduced in each case, considering both first and second-order integrators. The well-known "Strang splitting" is also presented, which will be used in Chapter 4 to derive second-order schemes.

The second section of the introductory chapter focuses on the concept of well-balanced methods, emphasizing the importance of this property and describing in detail the strategy to achieve it, based on the idea presented in [25].

This strategy involves using a reconstruction operator that is well-balanced, meaning that when applied to the cell averages of a stationary solution, the approximations obtained should coincide with that stationary solution. Reconstruction operators do not possess this property in general, but we construct one that does by following the process outlined below:

1. Seek, if possible, a stationary solution such that its average in the cell matches the value of our approximation in that cell.
2. Calculate the fluctuations in the cells of the stencil of a standard reconstruction operator of our choice, i.e., the differences between the averages of the solution and those of the stationary solution obtained in the first step. Apply the chosen standard reconstruction operator to these fluctuations.
3. Sum the stationary solution obtained in the first step and the reconstruction from the second step to obtain the well-balanced reconstruction operator.

The obtained operator is well-balanced and of the same order as the standard reconstruction operator used in the second step, assuming that stationary solutions are sufficiently smooth and that the standard reconstruction operator is exact in the case of zero functions.

Additionally, it is necessary to adapt the expression of the source term integral so that the final scheme obtained is well-balanced.

The last section is dedicated to introducing Lagrangian coordinates, which will be used in Chapters 2 and 3. Since these coordinates follow the trajectories of the flow particles, consider a particle  $\xi$  and define the characteristic curves as

$$\begin{cases} \frac{\partial x}{\partial t}(\xi, t) = u(x(\xi, t), t), \\ x(\xi, 0) = \xi. \end{cases}$$

Moreover, given a function in Eulerian coordinates  $(x, t) \mapsto \mathbf{U}(x, t)$ , its equivalent in Lagrangian coordinates is defined as

$$\bar{\mathbf{U}}(\xi, t) = \mathbf{U}(x(\xi, t), t).$$

Finally, the Jacobian of the Lagrangian mapping is defined as

$$L(\xi, t) = \frac{\partial x}{\partial \xi}(\xi, t),$$

and a series of properties relating the partial derivatives in different coordinates are derived. This section considers the case of the Euler equations, which are reformulated in terms of Lagrangian coordinates and are eventually particularised to the case of shallow water equations.

## Chapter 2: Implicit and implicit-explicit Lagrange-projection exactly well-balanced finite volume schemes for 1D shallow water system

This chapter presents the work published in [22], where the Lagrange-Projection technique is considered within the framework of finite volume schemes applied to the shallow water

system. As mentioned before, this technique consists of two steps: first, the system in Lagrangian coordinates is solved, and then the result is projected to Eulerian coordinates. These steps are known as the Lagrangian step and the projection or transport step, respectively. The reason for returning to Eulerian coordinates in each step is that using pure Lagrangian coordinates and tracking moving meshes can be cumbersome, especially in two-dimensional configurations. Additionally, this methodology allows the decoupling of acoustic and transport phenomena, allowing the design of implicit-explicit schemes and enabling naturally larger time steps. This is particularly beneficial in the case of a small Froude number, as using implicit or implicit-explicit schemes reduces the CFL restriction to slow transport waves rather than acoustic ones, which are more restrictive.

In the Lagrangian step, two versions of the scheme are considered: a nonlinear implicit version and an implicit-explicit version, based on how the geometric source term is treated. The transport step is always performed explicitly. Exactly well-balanced first and second-order versions of the schemes are presented, where water at rest solutions are preserved. To achieve this, exactly well-balanced reconstruction operators are used, constructed using the strategy described in [25].

After conducting various numerical tests, it is concluded that the scheme where the Lagrangian step is solved implicitly-explicitly is more efficient than the nonlinear implicit one, as the latter requires solving a nonlinear system and several iterations of a fixed-point algorithm. Additionally, as expected, first-order schemes are more diffusive than second-order ones. Finally, an improvement in the efficiency of schemes in which the Lagrangian step is performed implicitly (both nonlinear and implicit-explicit) is observed compared to those performed explicitly, especially in the case of a low Froude number, where implicit schemes allow much larger time steps.

### Chapter 3: Implicit Lagrange-projection well-balanced finite volume scheme for the Ripa model

In the third chapter, we propose a strategy to numerically solve the Ripa model by once again applying the Lagrange-Projection technique. This consists on solving the Ripa system in Lagrangian coordinates first and then projecting the solution to Eulerian coordinates. Similarly to Chapter 2, the Lagrangian system is solved implicitly, while the projection step is done explicitly.

First order well-balanced finite volume schemes are designed for this model, preserving the so-called hydrostatic steady states corresponding to  $u = 0$ , and satisfying

$$\partial_x p = -2p \frac{\partial_x z}{h}.$$

Particular cases of these steady states include those corresponding to water at rest,

$$u = 0, \quad h + z = \text{constant}, \quad \theta = \text{constant},$$

as well as isobaric steady states, where pressure is constant and the bottom is flat:

$$u = 0, \quad \frac{g}{2}h^2\theta = \text{constant}, \quad z = \text{constant}.$$

Since in this case the stationary solutions are determined by setting a profile for  $h$  or  $\theta$ , it is impossible to design exactly well-balanced numerical schemes for all hydrostatic solutions. That is why we have chosen to use an approximation technique, so that once a discrete profile for  $h$  has been chosen, we proceed to approximate the pressure using a collocation method (see [56]) to approximate the solutions of the ODE

$$\partial_\xi p = -2p \frac{\partial_\xi z}{h}.$$

Once the pressure is recovered at the quadrature points, and the discrete profile of  $h$  is known, it is possible to recover the stationary discrete profile of  $\theta$ . The schemes obtained are, therefore, well balanced but not exactly well balanced.

A section of numerical tests is included, comparing the proposed scheme with two other Lagrangian-Projected schemes: one that is not well-balanced and another that is exactly well-balanced for the water at rest steady states and the isobaric steady states but not for any hydrostatic steady state.

## Chapter 4: Semi-implicit fully exactly well-balanced finite volume schemes for the 1D shallow water system

In this chapter, the goal is the design of implicit first and second order schemes that preserve all steady state solutions of the shallow water equations. To achieve this, a splitting technique is applied, resulting in two systems to solve: a pressure system and a transport system. A relaxed system of the pressure system is then considered. This will be solved implicitly, while the transport system will be solved explicitly. It is possible to solve the pressure system first, followed by the transport system, or vice versa. Therefore, for each order, two versions of the scheme are considered.

For second order, the well-known "Strang splitting" is used along with second order reconstructions in space and first order in time. This splitting involves taking one step of one system with a time step  $\Delta t/2$ , followed by a step of the second system with a time step  $\Delta t$ , and ending with a step of the first system with a time step  $\Delta t/2$ . Therefore, in the case where the pressure system is solved first, two implicit steps will be necessary, while in the case where the transport system is solved first, only the second step will be implicit, making the latter potentially more efficient.

Thus, as in the previous chapters, schemes that allow larger time steps than explicit schemes are obtained, making them more efficient, especially for low Froude numbers.

When studying the results of the proposed schemes with various tests and increasing the CFL value, different behaviors have been observed depending on which system is solved

first. For first order, better performance is observed in the scheme where the pressure system is solved first, while for second order, stability is better in the scheme that starts with the transport system.

Additionally, tests have been performed considering perturbations of a subcritical solution, as well as a smooth transcritical solution, with good results observed in both cases.

## Chapter 5: Conclusions and future work

In this final section of the document, conclusions drawn from the results obtained throughout this thesis are presented.

Regarding future work, several aspects are outlined for further exploration with the aim of applying the developments in this thesis to other problems. These include extending the methods to the two-dimensional case, designing schemes of order higher than two, and applying the studied strategies to other systems.

# Chapter 1

## Mathematical settings

### 1.1 Finite volume numerical methods for 1D systems of balance laws

In this section we will present a brief review on systems of balance laws, which provide a framework for understanding numerous important phenomena within fluid dynamics. Once we describe these systems and general finite volume methods for them, we will also present the particular numerical strategies that will be used afterwards in the following chapters.

#### 1.1.1 Systems of balance laws

A system of balance laws in one space dimension takes the form

$$U_t + f(U)_x = S(U)H_x, \quad (1.1.1)$$

where  $U(x, t)$  is a vector function that takes values in an open set  $\Omega \subset \mathbb{R}^N$ , called the set of states, and  $f : \Omega \rightarrow \mathbb{R}^N$  is the flux function. The term on the right hand side,  $S(U)H_x$ , is the source term, where  $S : \Omega \rightarrow \mathbb{R}$  and  $H : \mathbb{R} \rightarrow \mathbb{R}$  are known functions. Note that function  $H$  could be the identity function.

In the particular case where  $S(U) = 0$  or  $H$  is a constant function, system (1.1.1) reduces to

$$U_t + f(U)_x = 0, \quad (1.1.2)$$

and it is called a system of conservation laws.

Note that a system of balance laws can be seen as a system of conservation laws with a source term. Therefore, many of the concepts we introduce below refer initially to conservation laws, and are later extended to equilibrium laws.

**Definition 1.1.1.** *The system (1.1.2) is said to be hyperbolic if for every  $U \in \Omega$ ,  $J_f(U)$ , the jacobian matrix of  $f$ , has  $N$  real eigenvalues*

$$\lambda_1(U) \leq \dots \leq \lambda_N(U),$$

with associated eigenvectors

$$R_1(U), \dots, R_N(U).$$

If all the eigenvalues are distinct, that is, if

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

the system (1.1.2) is said to be strictly hyperbolic.

Let us now define when a characteristic field is linearly degenerate or genuinely nonlinear:

**Definition 1.1.2.** *The characteristic field  $R_i(U)$  is said to be linearly degenerate if*

$$\nabla \lambda_i(U) \cdot R_i(U) = 0, \quad \forall U \in \Omega, \quad (1.1.3)$$

where  $\nabla \lambda_i(U)$  denotes the gradient of  $\lambda_i(U)$ :

$$\nabla \lambda_i(U) = \left( \frac{\partial \lambda_i}{\partial u_1}, \dots, \frac{\partial \lambda_i}{\partial u_N} \right).$$

The characteristic field  $R_i(U)$  is said to be genuinely nonlinear if

$$\nabla \lambda_i(U) \cdot R_i(U) \neq 0, \quad \forall U \in \Omega. \quad (1.1.4)$$

In the following, we will suppose that the characteristic fields are of the type linearly degenerate or genuinely nonlinear.

Let us now focus on Cauchy problems for conservation laws, which are represented as follows:

$$\begin{cases} U_t + f(U)_x = 0, \\ U(x, 0) = U^0(x), \end{cases} \quad (1.1.5)$$

where we consider the particular case  $\Omega = \mathbb{R}$ .

In this context, a function  $U : \mathbb{R} \times [0, \infty) \rightarrow \Omega$  that is  $C^1$  and satisfies (1.1.5) is referred to as a classical solution of that Cauchy problem.

As (1.1.5) may not have a classical solution, the concept of weak solution is introduced by means of a variational formulation:

**Definition 1.1.3.** Let us consider the Cauchy problem (1.1.5) with initial condition  $U^0 \in \mathcal{L}_{loc}^\infty(\mathbb{R})$ , where  $\mathcal{L}_{loc}^\infty$  denotes the space of locally bounded measurable functions. A function  $U \in \mathcal{L}_{loc}^\infty(\mathbb{R} \times [0, \infty))$  such that  $U \in \Omega$  almost everywhere is said to be a weak solution of (1.1.5) if for any  $C^1$  test function  $\Phi$  with compact support in  $\mathbb{R} \times [0, \infty)$ , one has

$$\int_0^\infty \int_{\mathbb{R}} \left( U(x, t) \Phi_t(x, t) + f(U(x, t)) \Phi_x(x, t) \right) dx dt + \int_{\mathbb{R}} U^0(x) \Phi(x, 0) dx = 0. \quad (1.1.6)$$

A weak solution satisfies the Cauchy problem (1.1.5) in the sense of the distributions theory (see [54]). In the case in which the function  $U$  is piecewise  $C^1$ , the following result holds:

**Theorem 1.1.1.** Let  $\mathcal{C}$  be a  $C^1$  curve in  $\mathbb{R}^2$  defined by  $x = \xi(t)$ , that cuts the open set  $\Omega \subset \mathbb{R}^2$  in two open sets  $\Omega_-$  and  $\Omega_+$ , defined respectively by  $x < \xi(t)$  and  $x \geq \xi(t)$ . A piecewise  $C^1$  function  $U : \mathbb{R} \times [0, \infty) \rightarrow \Omega$  is a weak solution of the Cauchy problem (1.1.5) if and only if:

- $U$  is a classical solution of (1.1.5) in  $\Omega_-$  and in  $\Omega_+$ ;
- the condition

$$\sigma(U_+ - U_-) = f(U_+) - f(U_-) \quad (1.1.7)$$

is satisfied on  $\mathcal{C} \cap \Omega$ , where  $\sigma := \dot{\xi}$  is the speed of the propagation of the discontinuity depending on  $U_-$ ,  $U_+$  which are, respectively, the left and right limit states of the solution at the discontinuity.

The condition (1.1.7) is known as the Rankine-Hugoniot condition, and it is usually written as

$$\sigma[U] = [f(U)], \quad (1.1.8)$$

where

$$[U] = U_+ - U_-, \quad [f(U)] = f(U_+) - f(U_-).$$

Since, in general, there isn't a unique weak solution for (1.1.5), we require an entropy condition to select physically meaningful solutions from the possible candidates. Specifically, it is common to consider an entropy pair  $(\eta, q)$ , where  $\eta : \Omega \rightarrow \mathbb{R}$  is a convex function referred to as the *entropy*, and  $q : \Omega \rightarrow \mathbb{R}$  is a smooth function known as the *entropy flux*. This pair should satisfy the following relationship:

$$\nabla \eta(U)^T J_f(U) = \nabla q(U)^T, \quad \forall U \in \Omega, \quad (1.1.9)$$

being

$$\nabla \eta(U) = \left( \frac{\partial \eta}{\partial u_1}, \dots, \frac{\partial \eta}{\partial u_N} \right)^T.$$

It can be proved (see [66]) that smooth solutions of (1.1.2) satisfy the equation

$$\eta(U)_t + q(U)_x = 0, \quad x \in \mathbb{R}, t \in [0, \infty). \quad (1.1.10)$$

A weak solution of (1.1.2) is considered an *entropy solution* if the inequality

$$\eta(U)_t + q(U)_x \leq 0 \quad (1.1.11)$$

holds in the sense of distributions.

Additionally, for a piecewise  $C^1$  weak solution  $U$  of (1.1.2), condition (1.1.11) is satisfied if and only if

$$\sigma[\eta(U)] \geq [q(U)] \quad (1.1.12)$$

holds across all the discontinuities.

Going back now to systems of balance laws (1.1.1) in which  $H$  is continuous, one can extend the definition of weak solution for

$$\begin{cases} U_t + f(U)_x = S(U)H_x, \\ U(x, 0) = U^0(x), \end{cases} \quad (1.1.13)$$

as follows:

**Definition 1.1.4.** *Let us consider the Cauchy problem (1.1.13) with initial condition  $U^0 \in \mathcal{L}_{loc}^\infty(\mathbb{R})$  and  $H$  be a continuous function. A function  $U \in \mathcal{L}_{loc}^\infty(\mathbb{R} \times [0, \infty))$  such that  $U \in \Omega$  almost everywhere is said to be a weak solution of (1.1.13) if for any  $C^1$  function  $\Phi$  with compact support in  $\mathbb{R} \times [0, \infty)$ , one has*

$$\begin{aligned} & \int_0^\infty \int_{\mathbb{R}} \left( U(x, t) \Phi_t(x, t) + f(U(x, t)) \Phi_x(x, t) \right) dx dt \\ & - \int_0^\infty \int_{\mathbb{R}} S(U(x, t)) H_x(x) \Phi(x, t) dx dt + \int_{\mathbb{R}} U^0(x) \Phi(x, 0) dx = 0. \end{aligned} \quad (1.1.14)$$

## 1.1.2 Some systems of balance laws models

In this section we will introduce and briefly discuss the models that will be considered in this thesis: the shallow water model, the Ripa model, the isentropic gas dynamics model and the compressible Euler system. Even though we will only derive numerical methods for the shallow water and the Ripa model, the other two systems will also be presented here, since the isentropic gas dynamics one will be used to introduce how relaxation solvers can be applied and the Euler equations will be used to introduce the use of Lagrangian coordinates.

### The shallow water model

The shallow water model is a simplification of the Navier-Stokes one, obtained after performing a vertical averaging of the Navier-Stokes equations under a set of assumptions:

- Water is homogeneous and incompressible.
- The pressure is hydrostatic, increasing with depth.
- The only internal force acting on the fluid is pressure, with viscous effects being neglected.
- Both the bottom over which the fluid evolves and the free surface can be represented by functions depending on one of the horizontal variables,  $x$ , in the case of the bottom, and on  $x$  and time  $t$ , in the case of the free surface.
- Fluid velocity depends only on  $x$  and  $t$ , with vertical variations of the horizontal components of velocity being neglected.

The shallow water equations (SWE) describe the flow of liquids with a single thin layer, i.e., when the horizontal dimension is considerably larger than the vertical dimension. In the one-dimensional case, these equations are also known as the Saint-Venant equations, as this French engineer and mathematician first derived them in 1871 (see [41]). The one dimensional shallow water system is given by the following equations:

$$\begin{cases} \partial_t h + \partial_x(hu) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + g\frac{h^2}{2}\right) = -gh\partial_x z, \end{cases} \quad (1.1.15)$$

where  $h(x, t)$  is the thickness of the water layer,  $u(x, t)$  is the horizontally averaged velocity in the vertical direction,  $z(x)$  denotes a certain smooth topography measured from a reference level, and  $g > 0$  is the gravitational constant. The free surface is usually denoted by  $\eta$  and is given by  $\eta = h + z$ . The first equation corresponds to mass conservation, and the second one to the momentum equation.

As far as the eigenstructure of the shallow water system is concerned, the system is strictly hyperbolic over the phase space  $\Omega = \{(h, hu)^T \in \mathbb{R}^2 \mid h > 0\}$  and it is composed of two genuinely nonlinear fields associated with the eigenvalues  $u - c$  and  $u + c$ , where  $c = \sqrt{gh}$  is the sound speed. The regions where  $u^2 < c^2$  (resp.  $u^2 > c^2$ ) are called subcritical (resp. supercritical). Defining the Froude number as

$$Fr = \frac{|u|}{\sqrt{gh}}, \quad (1.1.16)$$

we can clearly relate the subcritical and supercritical regions with it as: if  $Fr < 1$  (resp.  $Fr > 1$ ) the region is subcritical (resp. supercritical).

Moreover, the equilibrium or steady states of the SWE are given by

$$\begin{cases} \frac{d}{dx}(hu)^e = 0, \\ \frac{d}{dx}\left(hu^2 + g\frac{h^2}{2}\right)^e = -gh^e z'. \end{cases} \quad (1.1.17)$$

In particular, smooth equilibria satisfy

$$\begin{cases} (hu)^e = C_1, \\ \frac{(u^e)^2}{2} + g(h^e + z) = C_2, \end{cases} \quad (1.1.18)$$

where  $C_1$  and  $C_2$  are two real constants. Of course, when  $C_1 = 0$ , we obtain the so-called water at rest steady states, for which

$$q^e = (hu)^e = 0, \quad \eta^e = h^e + z = \text{cst}. \quad (1.1.19)$$

Remark that for any two fixed constants  $C_1$  and  $C_2$ , system (1.1.18) is equivalent to setting  $q^e(x) = C_1$  and  $h^e(x)$  solution of the cubic equation

$$(h^e)^3 + \left(z - \frac{C_2}{g}\right)(h^e)^2 + \frac{C_1^2}{2g} = 0. \quad (1.1.20)$$

Note that equation (1.1.20) does not always has a physical solution. It always has a negative root, but it does not always have positive ones.

Moreover, for fixed values  $C_1$  and  $C_2$ , one may define the function

$$\begin{aligned} f_{C_1} : (0, \infty] &\rightarrow \mathbb{R} \\ h &\mapsto \frac{C_1^2}{2h^2} + gh, \end{aligned}$$

and write (1.1.20) equivalently as

$$f_{C_1}(h) = C_2 - gz.$$

Following the study done in [20], for any fixed value  $C_1$ ,  $f_{C_1}$  has a global minimum at

$$h_{crit}(C_1) = \frac{|C_1|^{2/3}}{g^{1/3}}. \quad (1.1.21)$$

Let us denote by  $m_s(C_1) = f_{C_1}(h_{crit}(C_1))$ . Then one can find three different possibilities when solving (1.1.20):

1. If  $C_2 - gz < m_s(C_1)$ , then there exists no real positive solution.
2. If  $C_2 - gz = m_s(C_1)$ , then there exist a unique real positive solution given by  $h = h_{crit}(C_1)$ .
3. If  $C_2 - gz > m_s(C_1)$ , then there are exactly two positive solutions: the subcritical one, that satisfies  $h > h_{crit}(C_1)$  and the supercritical one, that verifies  $h < h_{crit}(C_1)$ .

### The Ripa model

The Ripa model was derived from the compressible Euler model in [79, 80] to incorporate horizontal temperature gradients. The equations of the Ripa system are the following:

$$\begin{cases} \partial_t h + \partial_x(hu) = 0, \\ \partial_t(hu) + \partial_x\left(hu^2 + \frac{g}{2}h^2\theta\right) = -gh\theta\partial_x z, \\ \partial_t(h\theta) + \partial_x(h\theta u) = 0, \end{cases} \quad (1.1.22)$$

where  $\theta$  is the potential temperature field.

This system is also strictly hyperbolic, with eigenvalues given by  $u - c$ ,  $u$ , and  $u + c$ , where now  $c = \sqrt{gh\theta}$ .

Among all the possible steady states of the Ripa system, which satisfy

$$\begin{cases} \partial_x(hu)^e = 0, \\ \partial_x\left(hu^2 + \frac{g}{2}h^2\theta\right)^e = -g(h\theta)^e\partial_x z, \\ \partial_x(h\theta u)^e = 0, \end{cases} \quad (1.1.23)$$

we can highlight the hydrostatic steady states, that is, those steady states corresponding to

$$u = 0,$$

which satisfy

$$\partial_x p = -2p \frac{\partial_x z}{h}, \quad (1.1.24)$$

with  $p = \frac{g}{2}h^2\theta$ .

Particular cases of these steady states are the ones corresponding to "water at rest",

$$u = 0, \quad h + z = cst., \quad \theta = cst., \quad (1.1.25)$$

as well as the isobaric steady states corresponding to constant pressure and flat topography,

$$u = 0, \quad \frac{g}{2}h^2\theta = cst., \quad z = cst. \quad (1.1.26)$$

### Euler equations

The Euler equations can be written in the one-dimensional case and assuming slab symmetry as

$$\begin{cases} \frac{\partial \rho}{\partial t} + \frac{\partial}{\partial x}(\rho u) = 0, \\ \frac{\partial}{\partial t}(\rho u) + \frac{\partial}{\partial x}(\rho u^2 + p) = 0, \\ \frac{\partial}{\partial t}(\rho e) + \frac{\partial}{\partial x}((\rho e + p)u) = 0, \end{cases} \quad (1.1.27)$$

where  $\rho$  is the density,  $u$  is the velocity,  $p$  is the pressure,  $e = \varepsilon + \frac{|u|^2}{2}$  is the specific total energy and  $\varepsilon$  is the specific internal energy. They represent conservation of mass, conservation of momentum and conservation of total energy, respectively.

An equation for the pressure, known as the equation of state, has to be added to (1.1.27) to close the system. In general,

$$p = p(\rho, \varepsilon).$$

The eigenvalues of the previous system are given by  $u - c$ ,  $u$  and  $u + c$ , where  $c = \sqrt{\frac{\partial p}{\partial \rho} + \frac{p}{\rho} \frac{\partial p}{\partial \varepsilon}}$  is the sound speed.

The case of an isentropic gas, where  $p = p(\rho)$  is a particular one of special interest. In that case it is enough to solve the equations of conservation of mass and momentum (1.1.28).

### Isentropic gas dynamics model

Isentropic gas dynamics refers to the study of the motion and behavior of gases with constant entropy. The isentropic gas dynamics equations are given by

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + p(\rho)) = 0, \end{cases} \quad (1.1.28)$$

where  $\rho$  is the density,  $u$  the velocity and  $p(\rho)$  the pressure. Its eigenvalues are given by  $u - c$  and  $u + c$ , where now  $c = \sqrt{p'(\rho)}$ .

### 1.1.3 Finite volume numerical methods

In this section, again, we will first talk about finite volume numerical methods for systems of conservation laws (1.1.2) and after that, we will give the details about the case of balance laws (1.1.1).

The aim is to obtain approximations of the solution of the equations by discrete values  $U_i(t)$ ,  $i \in \mathbb{Z}$ . To do that, space is discretized in a set of cells or finite volumes  $I_i = [x_{i-1/2}, x_{i+1/2})$  using a space step  $\Delta x$ , where  $x_{i+1/2} = i\Delta x$  and  $x_i = (x_{i-1/2} + x_{i+1/2})/2$  are, respectively, the cell interfaces and cell centers for  $i \in \mathbb{Z}$ . Similarly, the time step is denoted as  $\Delta t$  and the instants are written as  $t_n = n\Delta t$ , with  $n \in \mathbb{N}$ . The approximation of the averages at time  $t_n$  will be denoted as  $U_i^n$ . There exists a restriction on the time step to prevent the blow up of the numerical schemes. This restriction, called the CFL condition after Courant, Friedrichs and Levy (see [40]), is the following:

$$\Delta t \leq \frac{\Delta x \cdot \text{CFL}}{a}, \quad (1.1.29)$$

where  $a$  is an approximation of the speed of propagation. In the explicit cases,  $\text{CFL} \leq 1$ . However, in the implicit ones, CFL can take values greater than 1, allowing the consideration of bigger time steps.

In a finite volume method, the discrete values  $U_i(t)$  are approximations of the averages of the exact solution over the cells, that is:

$$U_i(t) \approx \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U(x, t) dx. \quad (1.1.30)$$

Now, integrating (1.1.2) over every cell  $I_i$  gives

$$\frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U_t dx = -\frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} f(U)_x dx. \quad (1.1.31)$$

We can then write

$$\frac{dU_i(t)}{dt} = -\frac{1}{\Delta x} (F_{i+1/2}^t - F_{i-1/2}^t), \quad (1.1.32)$$

being  $F_{i+1/2}^t$  an approximation of the average of the flux at  $x_{i+1/2}$ , called *numerical flux*, which constitutes a semi-discrete finite volume conservative scheme for solving (1.1.2).

Consistency is as a fundamental prerequisite for a scheme to guarantee an appropriate approximation of the equation. In the context of conservative schemes, consistency is defined as follows:

**Definition 1.1.5.** *The scheme (1.1.32) for system (1.1.2) is said to be consistent if the numerical flux*

$$F_{i+1/2}^t = \mathbb{F}(U_{i-q}(t), \dots, U_{i+p}(t)),$$

where  $\mathbb{F}$  is a Lipschitz continuous function satisfies

$$\mathbb{F}(U, \dots, U) = f(U), \quad \forall U \in \Omega.$$

### 1.1.3.1 Reconstruction operators

We will now give a definition of what is known as a reconstruction operator, that will play a fundamental role in the numerical methods that we consider:

**Definition 1.1.6.** *A reconstruction operator of order  $p$  is an operator that, given a family of cell values  $\{U_i\}$ , provides at every cell  $I_i$  a smooth function that depends on the values at some neighbour cells whose indexes belong to the so-called stencil  $S_i$ :*

$$P_i(x) = P_i(x; \{U_j\}_{j \in S_i}), \quad (1.1.33)$$

so that, if the cell values are the averages of a smooth function  $U$ :

$$U_i = \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U(x) dx, \quad \forall i \in \mathbb{Z}, \quad (1.1.34)$$

then

$$U_{i+\frac{1}{2}-} = U(x_{i+\frac{1}{2}}) + \mathcal{O}(\Delta x^p), \quad (1.1.35)$$

$$U_{i-\frac{1}{2}+} = U(x_{i-\frac{1}{2}}) + \mathcal{O}(\Delta x^p), \quad (1.1.36)$$

being

$$U_{i+\frac{1}{2}-} = P_i(x_{i+\frac{1}{2}}), \quad (1.1.37)$$

$$U_{i-\frac{1}{2}+} = P_i(x_{i-\frac{1}{2}}), \quad (1.1.38)$$

where the states  $U_{i+\frac{1}{2}-}$  and  $U_{i-\frac{1}{2}+}$  are the so called reconstructed states at the intercells.

We will consider that the value  $F_{i+1/2}^t$  in expression (1.1.32) is computed as

$$F_{i+1/2}^t = \mathbb{F}(U_{i+1/2-}^t, U_{i+1/2+}^t), \quad (1.1.39)$$

so the numerical flux  $\mathbb{F}$  is evaluated at the reconstructed states at the intercells.

In this thesis, first and second order reconstruction operators will be used. For the first order case, a constant reconstruction operator is used. For the second order case a MUSCL (Monotone Upwind Scheme for Conservation Laws) reconstruction operator has been considered. This operator was introduced in [90] and is based on a piecewise linear reconstruction of the form

$$P_i(x) = U_i + \Delta_i U(x - x_i), \quad (1.1.40)$$

where  $\Delta_i U$  is an approximation of the first-order spatial derivative at  $x_i$ , computed by means of a limiter that avoids the appearance of spurious oscillations in case of discontinuities. Two slope limiters have been used in this work: minmod and avg (see [83]). Both limiters are computed componentwise. They are defined as follows:

- The minmod limiter is given by

$$\Delta_i U = \text{minmod} \left( \frac{U_{i+1} - U_i}{\Delta x}, \frac{U_i - U_{i-1}}{\Delta x} \right), \quad (1.1.41)$$

where

$$\text{minmod}(a, b) = \begin{cases} \min(a, b) & \text{if } a > 0, b > 0, \\ \max(a, b) & \text{if } a < 0, b < 0, \\ 0 & \text{otherwise.} \end{cases} \quad (1.1.42)$$

- The avg or harmod limiter is given by

$$\Delta_i U = \text{avg} \left( \frac{U_{i+1} - U_i}{\Delta x}, \frac{U_i - U_{i-1}}{\Delta x} \right), \quad (1.1.43)$$

where

$$\text{avg}(a, b) = \begin{cases} \frac{|a|b + |b|a}{|a| + |b|} & \text{if } |a| + |b| > 0, \\ 0 & \text{otherwise.} \end{cases} \quad (1.1.44)$$

Let us note that in view of (1.1.41) and (1.1.43), in the computation of (1.1.40), the stencils are composed by two neighbour cells and, given a cell  $i$ , the reconstructed states at the intercells are given by

$$\begin{aligned} U_{i-1/2+} &= U_i - \frac{\Delta x}{2} \Delta_i U, \\ U_{i+1/2-} &= U_i + \frac{\Delta x}{2} \Delta_i U. \end{aligned} \tag{1.1.45}$$

Now, supposing a  $H$  continuous function, we can proceed similarly for the resolution of a system of balance laws (1.1.1), writing the semi-discrete numerical method as

$$\frac{dU_i(t)}{dt} = -\frac{1}{\Delta x} (F_{i+1/2}^t - F_{i-1/2}^t) + \frac{1}{\Delta x} S_i^t, \tag{1.1.46}$$

where the numerical fluxes are defined as (1.1.39) and  $S_i^t$  stands for the following approximation of the integral of the source term at cell  $i$ :

$$S_i^t = \int_{x_{i-1/2}}^{x_{i+1/2}} S(P_i^t(x)) H_x dx, \tag{1.1.47}$$

where  $P_i^t$  is a reconstruction operator corresponding to time  $t$ .

### 1.1.3.2 Relaxation solvers

We will now introduce the idea of relaxation solvers (see [19]), that will be applied in the different methods presented in this thesis.

Let us start by presenting the so called approximate Riemann solvers in the sense of Harten, Lax, Van Leer [62]. A Riemann problem for (1.1.2) consists on solving (1.1.5) with initial condition

$$U^0(x) = \begin{cases} U_l & \text{if } x < 0, \\ U_r & \text{if } x > 0, \end{cases} \tag{1.1.48}$$

where  $U_l$  and  $U_r$  are two constant states. It is not difficult to see that the solution can be written as a function of  $x/t$ .

**Definition 1.1.7.** *An approximate Riemann solver for (1.1.2) is a vector function  $R(x/t, U_l, U_r)$  that is an approximation of the solution to the Riemann problem, satisfying*

- *Consistency relation:*

$$R(x/t, U, U) = U.$$

- *Conservativity identity:*

$$F_l(U_l, U_r) = F_r(U_l, U_r), \tag{1.1.49}$$



where the left and right numerical fluxes are defined by

$$F_l(U_l, U_r) = f(U_l) - \int_{-\infty}^0 (R(v, U_l, U_r) - U_l) dv, \quad (1.1.50)$$

$$F_r(U_l, U_r) = f(U_r) + \int_0^{\infty} (R(v, U_l, U_r) - U_r) dv. \quad (1.1.51)$$

Let us consider a sequence of cell averages  $\{U_i^0\}$  of a function  $U^0(x)$  that we suppose constant in each cell  $I_i$ . The idea of the approximate Riemann solvers is to consider that close to each interface a translated Riemann problem has to be solved. Therefore, the approximate solution for  $U(x, t)$  is given by:

$$U(x, t) = R\left(\frac{x - x_{i+1/2}}{t}, U_i^0, U_{i+1}^0\right) \text{ if } x_i < x < x_{i+1}. \quad (1.1.52)$$

The approximation of the solution at time  $t$ ,  $U_i(t)$ , is the average over the cell  $I_i$  of this approximate solution at time  $t$ , and using the definitions of  $F_l$  and  $F_r$  (1.1.50)-(1.1.51) and the conservativity assumption (1.1.49), the scheme writes as

$$\frac{dU_i(t)}{dt} = -\frac{\Delta t}{\Delta x} (F(U_i^0, U_{i+1}^0) - F(U_{i-1}^0, U_i^0)).$$

The use of the consistency assumption implies that the scheme is conservative.

The exact resolution of the Riemann problem is excessively complex and costly, particularly for systems of significant dimensions. Consequently, we opt for the use of approximate solvers instead. Simple solvers and relaxation solvers will now be introduced.

### Simple solvers

**Definition 1.1.8.** *A simple solver is an approximate Riemann solver  $R(x/t, U_l, U_r)$  that consists on a set of finitely many simple discontinuities, that is, there exists a finite number  $m \geq 1$  of speeds*

$$\sigma_0 = -\infty < \sigma_1 < \dots < \sigma_m < \sigma_{m+1} = \infty$$

and intermediate states

$$U_0 = U_l, U_1, \dots, U_{m-1}, U_m = U_r$$

that depend on  $U_l$  and  $U_r$  such that

$$R(x/t, U_l, U_r) = U_k \text{ if } \sigma_k < x/t < \sigma_{k+1}. \quad (1.1.53)$$

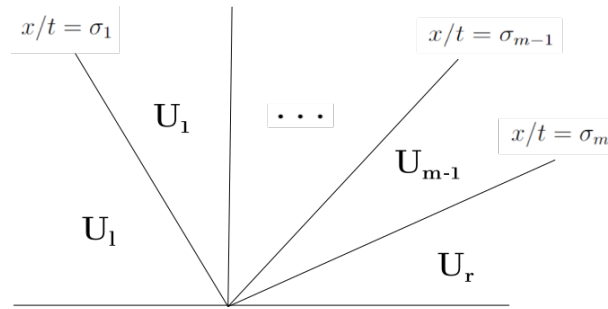


Figure 1.1: Sketch of a simple Riemann solver

Using the conservativity identity (1.1.49) it is possible to check that a simple solver satisfies:

$$\sum_{k=1}^m \sigma_k (U_k - U_{k-1}) = f(U_r) - f(U_l).$$

Since in all the works presented in this thesis the relaxation method will be applied, will now introduce the procedure to be applied.

### Relaxation solvers

Relaxation methods have been used in plenty of previous works such us [63, 39, 3, 18, 32]. Here we keep following [19].

**Definition 1.1.9.** A relaxation system for (1.1.2) is another system of conservation laws in higher dimension  $q > s$ ,

$$\partial_t g + \partial_x (\mathcal{A}(g)) = 0, \tag{1.1.54}$$

where  $g(x, t) \in \mathbb{R}^q$  and  $\mathcal{A}(g) \in \mathbb{R}^q$ . This system is assumed to be hyperbolic.

The link between (1.1.2) and (1.1.54) is made by the assumption that we have a linear operator

$$L : \mathbb{R}^q \rightarrow \mathbb{R}^s$$

and for any  $U$ , an equilibrium  $M(U)$  such that

$$LM(U) = U, \tag{1.1.55}$$

$$L\mathcal{A}(M(U)) = f(U). \tag{1.1.56}$$

Once system (1.1.54) has been solved, we define

$$U \equiv Lg. \tag{1.1.57}$$

The idea of the relaxation schemes is that  $U = Lg$  should be an approximation of the solution of (1.1.2) when  $g$  is a solution of (1.1.54) and is close to the equilibrium.



There are many different well known relaxation solvers such as Rusanov flux or HLL flux. Here we will focus on the Suliciu relaxation system, considered in works such as [86, 87, 39, 18, 32, 3]. In order to introduce the Suliciu relaxation system we will consider the isentropic gas dynamics equations given by (1.1.28). The first equation in (1.1.28) can be written for smooth solutions as

$$\partial_t \rho + u \partial_x \rho + \rho \partial_x u = 0, \quad (1.1.58)$$

and after multiplying by  $p'(\rho)$  we obtain

$$\partial_t p(\rho) + u \partial_x p(\rho) + \rho p'(\rho) \partial_x u = 0. \quad (1.1.59)$$

Finally, multiplying (1.1.58) by  $p(\rho)$  and (1.1.59) by  $\rho$  and summing, we obtain

$$\partial_t(\rho p(\rho)) + \partial_x(\rho p(\rho)u) + \rho^2 p'(\rho) \partial_x u = 0. \quad (1.1.60)$$

If we now define a new variable  $\pi = p(\rho)$  and we replace  $\rho^2 p'(\rho)$  by a constant  $a^2$ , we get the relaxation system

$$\begin{cases} \partial_t \rho + \partial_x(\rho u) = 0, \\ \partial_t(\rho u) + \partial_x(\rho u^2 + p(\rho)) = 0, \\ \partial_t(\rho \pi) + \partial_x(\rho \pi u) + a^2 \partial_x u = 0, \end{cases} \quad (1.1.61)$$

which has  $q = 3$  unknowns with  $g = (\rho, \rho u, \rho \pi)$  for  $s = 2$  unknowns  $(\rho, \rho u)$  for the original system. Then, in this case

$$\mathcal{A}(\rho, \rho u, \rho \pi) = (\rho u, \rho u^2 + \pi, \rho \pi u + a^2 u), \quad (1.1.62)$$

the linear operator is

$$L(g_1, g_2, g_3) = (g_1, g_2) \quad (1.1.63)$$

and the equilibrium is given by

$$M(\rho, \rho u) = (\rho, \rho u, \rho p(\rho)). \quad (1.1.64)$$

The constant  $a$  has to verify the known as subcharacteristic condition (see [19]), that consists on imposing that the eigenvalues of the original system lie between the eigenvalues of the relaxation one, so

$$|\lambda_j(U)| \leq a. \quad (1.1.65)$$

The benefit of solving system (1.1.61) instead of (1.1.28) is that the exact resolution of the Riemann problem for (1.1.61) is much simpler since it can be rewritten equivalently as a transport system:

$$\begin{cases} \partial_t(\pi + au) + (u + a/\rho) \partial_x(\pi + au) = 0, \\ \partial_t(\pi - au) + (u - a/\rho) \partial_x(\pi - au) = 0, \\ \partial_t(1/\rho + \pi/a^2) + u \partial_x(1/\rho + \pi/a^2) = 0. \end{cases} \quad (1.1.66)$$

### 1.1.3.3 Time integrators

Since the semi-discrete numerical method (1.1.46) is an ODE system, we will apply a numerical solver in order to approximate its solution. In this work, different types of solvers will be considered: explicit, implicit and semi-implicit. For each case we will focus on first and second order solvers. We will also introduce here the idea of the Strang splitting which will be used in the schemes proposed in Chapter 4.

#### Explicit time integrators

Let us suppose that we can write the semi-discrete scheme (1.1.46) as

$$\frac{dU_i(t)}{dt} = A(U_i(t)), \quad (1.1.67)$$

where  $A(U_i(t))$  represents the right-hand side of the scheme. Then, the first and second order explicit schemes can be written as:

- First order:

$$U_i^{n+1} = U_i^n + \Delta t A(U_i^n), \quad (1.1.68)$$

where  $U_i^n$  and  $U_i^{n+1}$  are the discrete averages approximations at time  $t^n$  and  $t^{n+1}$ , respectively.

- Second order: Applying the mid-point quadrature rule we obtain

$$U_i^{n+1} = U_i^n + \Delta t A(U_i^{n+1/2}). \quad (1.1.69)$$

Here  $U_i^{n+1/2}$  represents the discrete averages approximations at time

$$t^{n+1/2} = \frac{t^n + t^{n+1}}{2},$$

which are computed by applying an Euler step as

$$U_i^{n+1/2} = U_i^n + \frac{\Delta t}{2} A(U_i^n).$$

#### Implicit time integrators

In the implicit case, we will give a first order solver and two different second order ones applied to (1.1.46):

- First order:

$$U_i^{n+1} = U_i^n + \Delta t A(U_i^{n+1}). \quad (1.1.70)$$

- Second order (trapezoidal rule): applying the trapezoidal rule we obtain

$$U_i^{n+1} = U_i^n + \frac{\Delta t}{2} (A(U_i^n) + A(U_i^{n+1})). \quad (1.1.71)$$

- Second order (DIRK): we can also obtain a second order integrator by using a DIRK (Diagonally Implicit Runge-Kutta) scheme (see [76]), which in the second order case has the following Butcher tableau:

$$\begin{array}{c|cc} \gamma & \gamma & 0 \\ 1 & 1-\gamma & \gamma \\ \hline & 1-\gamma & \gamma, \end{array} \quad (1.1.72)$$

where  $\gamma = 1 - \frac{\sqrt{2}}{2}$ . The numerical method then writes:

$$\begin{aligned} U_i^1 &= U_i^n + \Delta t \gamma A(U_i^1), \\ U_i^{n+1} &= U_i^n + \frac{(1-\gamma)}{\gamma} U_i^1 + \Delta t \gamma A(U_i^{n+1}). \end{aligned} \quad (1.1.73)$$

### Semi-implicit time integrators

In the case in which the balance law (1.1.1) can be written as

$$U_t + f^1(U)_x + f^2(U)_x = S^1(U)H_x + S^2(U), \quad (1.1.74)$$

being  $f^1$  and  $S^1$  non stiff and  $f^2$  and  $S^2$  stiff, we could write the semi-discrete numerical method in the following way:

$$\frac{dU_i(t)}{dt} = A^1(U_i(t)) + A^2(U_i(t)), \quad (1.1.75)$$

where  $A^1$  contains the non stiff terms and  $A^2$  the stiff ones. The semi-implicit time integrators consist on treating the function  $A^1$  explicitly and the function  $A^2$  implicitly.

- First order: the non stiff terms are treated explicitly as in (1.1.68), and the stiff ones are treated implicitly as in (1.1.70).
- Second order: for higher order one can use the IMEX methods (see [14], [12]). Here, we consider the SSP2(2,2,2) IMEX scheme defined by the Butcher tableau (see [76]):

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}, \quad \begin{array}{c|cc} \gamma & \gamma & 0 \\ 1-\gamma & 1-2\gamma & \gamma \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}, \quad (1.1.76)$$

where  $\gamma = 1 - \frac{\sqrt{2}}{2}$ .

### Strang splitting

After having presented different explicit, implicit and semi-implicit time integrators that will be used in the following, we will introduce another strategy, called Strang splitting, that will be used in Chapter 4 and that can be applied explicitly or semi-implicitly, as we will see.

Let us suppose that our equation can be written as

$$\partial_t U = S_A(U) + S_B(U). \quad (1.1.77)$$

Then, instead of solving (1.1.77) directly we could consider a splitting that would consist on solving each of the two systems

$$\partial_t U = S_A(U), \quad (1.1.78)$$

and

$$\partial_t U = S_B(U), \quad (1.1.79)$$

sequentially. We could either solve the system defined by  $S_A$  first, followed by the one defined by  $S_B$  or vice versa, so we will have two different versions of the schemes depending on the order we consider. Moreover, we could decide to solve each of the systems explicitly or implicitly as we see fit.

Let us denote by  $S_A^\tau, S_B^\tau$  the approximate solution operators in the interval  $[t, t + \tau]$  of the corresponding exact solution operators to systems (1.1.78) and (1.1.79), respectively. Then, the first version of a second order scheme can be written as

$$U(x, t + \Delta t) = S_A^{\frac{\Delta t}{2}} \circ S_B^{\Delta t} \circ S_A^{\frac{\Delta t}{2}}(U(x, t)), \quad (1.1.80)$$

while the second version corresponds to

$$U(x, t + \Delta t) = S_B^{\frac{\Delta t}{2}} \circ S_A^{\Delta t} \circ S_B^{\frac{\Delta t}{2}}(U(x, t)). \quad (1.1.81)$$

In each of the steps we need to consider second order approximations in space, while the time stepping is just first order within the step, the second order in time being obtained thanks to Strang method.

## 1.2 Well-balanced methods

In this section we will focus on the design of numerical methods that are well-balanced. Let us consider a system of balance laws (1.1.1) that admits a non-trivial stationary solution satisfying

$$f(U)_x = S(U)H_x. \quad (1.2.1)$$

The solutions of (1.1.1) that satisfy the previous equation are called steady states or stationary solutions. Numerical schemes that preserve a family of steady states (resp. all steady states) are called well-balanced (resp. fully well-balanced). It is worth distinguishing between exactly preserving the steady states or a discrete approximation of them. The former are called exactly well-balanced while the latter are simply well-balanced schemes (see [56]).

The need to use well-balanced methods comes from the fact that when standard methods are used to solve (1.1.1) with an initial condition that reflects a perturbation to a stable solution, the discretization errors perturb the steady state throughout the entire computational domain from the very first time step. If these errors are approximately as significant as the initial perturbation, it becomes impossible to differentiate between the waves we intend to simulate and those that emerge as a consequence of the discretization errors. Furthermore, even though we can reduce discretization errors by refining the mesh or opting for higher-order methods, the resulting increase in computational expenses may become prohibitive.

In the different works presented in this thesis we will apply the strategy proposed in [25] for the design of well-balanced schemes. A brief review of it will now be presented.

**Definition 1.2.1.** *Given a stationary solution  $U^e$  of (1.1.1) a reconstruction operator is said to be exactly well-balanced for  $U^e$  if, for every cell index  $i$ , the following equality holds:*

$$P_i(x) = U^e(x), \quad \forall x \in I_i,$$

where  $P_i$  is the approximation of  $U^e$  obtained by applying the reconstruction operator to the sequence of cell averages  $\{U_i^e\}$  of  $U^e$ :

$$U_i^e = \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U^e(x) dx.$$

It can be proved that the following result for exactly well-balanced reconstruction operators holds:

**Theorem 1.2.1.** *If the reconstruction operator  $P_i$  is exactly well-balanced for a continuous stationary solution  $U^e$  of (1.1.1), then the numerical method (1.1.46) with*

$$S_i^t = \int_{x_{i-1/2}}^{x_{i+1/2}} S(P_i^t(x)) H_x(x) dx, \quad (1.2.2)$$

*is also exactly well-balanced for  $U^e$ , i.e., the sequence of cell averages of  $U^e$  is an equilibrium of the ODE given by the numerical method (1.1.46) with (1.2.2), where  $P_i^t$  is defined as (1.1.33).*

*Proof.* Let us consider a continuous stationary solution  $U^e$  of (1.1.1) at time  $t$  and a exactly well-balanced reconstruction operator  $P_i^t$  applied to the cell averages  $\{U_i^e\}$ . We will prove that  $\{U_i^e\}$  is an equilibrium of (1.1.46).

Since  $P_i$  is exactly well-balanced, for every  $i$  we have:

$$P_i^t(x) = U^e(x) \quad \forall x \in I_i,$$

and

$$\begin{aligned} U_{i+1/2-}^t &= P_i^t(x_{i+1/2}) = U^e(x_{i+1/2}), \\ U_{i+1/2+}^t &= P_{i+1}^t(x_{i+1/2}) = U^e(x_{i+1/2}). \end{aligned}$$

Therefore

$$\begin{aligned} &F_{i+1/2}^t - F_{i-1/2}^t - S_i^t \\ &= \mathbb{F}(U_{i+1/2-}^t, U_{i+1/2+}^t) - \mathbb{F}(U_{i-1/2-}^t, U_{i-1/2+}^t) - \int_{x_{i-1/2}}^{x_{i+1/2}} S(P_i^t(x)) H_x(x) dx \\ &= \mathbb{F}(U^e(x_{i+1/2}), U^e(x_{i+1/2})) - \mathbb{F}(U^e(x_{i-1/2}), U^e(x_{i-1/2})) - \int_{x_{i-1/2}}^{x_{i+1/2}} S(U^e(x)) H_x(x) dx \\ &= f(U^e(x_{i+1/2})) - f(U^e(x_{i-1/2})) - \int_{x_{i-1/2}}^{x_{i+1/2}} S(U^e(x)) H_x(x) dx = 0, \end{aligned}$$

where in the last equality we have applied the consistency of the numerical flux and the equation satisfied by the stationary solution (1.2.1).

Then, the numerical method (1.1.46) is exactly well-balanced for  $U^e$ .  $\square$

In view of Theorem 1.2.1, we will now focus on the design of reconstruction operators that are exactly well-balanced, following the idea developed in [25]. To do so, we will start from standard reconstruction operators that we will denote as

$$Q_i(x) = Q_i(x; \{U_j\}_{j \in \mathcal{S}_i}).$$

We will now present an algorithm that allows the construction of exactly well-balanced operators from standard ones:

**Algorithm 1.2.1.** *Given a family of cell values  $\{U_i\}$ , at every cell  $I_i$ :*

1. *Find, if possible, a stationary solution  $U_i^e(x)$  in the cells belonging to the stencil of  $I_i$ ,  $\cup_{j \in \mathcal{S}_i} I_j$ , such that:*

$$\frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U_i^e(x) dx = U_i. \quad (1.2.3)$$

*Otherwise, take  $U_i^e \equiv 0$ .*

2. *Apply a standard reconstruction operator of order  $p$  to the cell values  $\{V_j\}_{j \in \mathcal{S}_i}$ , called fluctuations, given by*

$$V_j = U_j - \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} U_i^e(x) dx, \quad j \in \mathcal{S}_i,$$

*to obtain:*

$$Q_i(x) = Q_i(x; \{V_j\}_{j \in \mathcal{S}_i}).$$

3. Finally, define

$$P_i(x) = U_i^e(x) + Q_i(x). \quad (1.2.4)$$

The reconstruction operator defined in Algorithm 1.2.1 satisfies the following result:

**Theorem 1.2.2.** *The reconstruction operator  $P_i$  in (1.2.4) is exactly well-balanced for any stationary solution of (1.1.1) provided that  $Q_i$  is exact for the zero function. Additionally,  $P_i$  is conservative provided that  $Q_i$  is conservative and it is high-order accurate provided that the stationary solutions are smooth.*

Since the cell averages are usually approximated by applying quadrature formulas, Algorithm 1.2.1 can be updated by using the appropriate quadrature formulas when needed.

As in this thesis we focus on first and second order schemes, we apply the midpoint rule to approximate the integrals. Therefore:

$$\frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U(x) dx = U(x_i) + O(\Delta x^2), \quad (1.2.5)$$

which means that we identify cell averages with centered point values up to second order.

The way to proceed would be to first compute the sequence of initial cell averages  $\{U_i^0\}$  defined by the initial condition  $U^0$ , that is:

$$U_i^0 = U^0(x_i),$$

and then apply Algorithm 1.2.1 but now using the midpoint rule. Therefore, in the first and second order case, the algorithm applied to obtain an exactly well-balanced reconstruction operator would be:

**Algorithm 1.2.2.** *Given a family of cell values  $\{U_i\}$ , at every cell  $I_i$ :*

1. Find, if possible, a stationary solution  $U_i^e(x)$  defined in the cells belonging to the stencil of  $I_i$ ,  $\cup_{j \in \mathcal{S}_i} I_j$ , such that:

$$U_i^e(x_i) = U_i, \quad (1.2.6)$$

*i.e., look for  $U_i^e(x)$  in the cells of the stencil of  $I_i$  which solves the Cauchy problem*

$$\begin{cases} f(U)_x = S(U)H_x, \\ U(x_i) = U_i. \end{cases} \quad (1.2.7)$$

*Otherwise, take  $U_i^e \equiv 0$ .*

2. Apply a standard reconstruction operator of first/second order to the fluctuations given by

$$V_j = U_j - U_i^e(x_j), \quad j \in \mathcal{S}_i,$$

to obtain:

$$Q_i(x) = Q_i(x; \{V_j\}_{j \in \mathcal{S}_i}).$$

(Observe that  $V_i = 0$ ).

3. Finally, define

$$P_i(x) = U_i^e(x) + Q_i(x). \quad (1.2.8)$$

In our case, for the first order schemes the reconstruction operator used will be the piecewise constant one. Therefore, we have

$$Q_i(x) = V_i = 0,$$

and the exactly well-balanced reconstruction operator in this case would just be

$$P_i(x) = U_i^e(x).$$

For the second order case, we consider the MUSCL reconstruction operator (1.1.40):

$$Q_i(x) = V_i + \Delta_i V(x - x_i) = \Delta_i V(x - x_i),$$

where  $\Delta_i V$  is an approximation of the spatial derivatives of the fluctuations computed by means of a limiter (in our case, minmod or avg). Then, the exactly well-balanced reconstruction operator can be written as

$$P_i(x) = U_i^e(x) + \Delta_i V(x - x_i).$$

Finally, we still have the issue of computing the integral of the source term

$$S_i^t = \int_{x_{i-1/2}}^{x_{i+1/2}} S(P_i^t(x)) H_x(x) dx. \quad (1.2.9)$$

Applying quadrature formulas to approximate it could destroy the well-balance character of the method and the way to avoid this is by rewriting the integral as

$$\begin{aligned} \int_{x_{i-1/2}}^{x_{i+1/2}} S(P_i^t(x)) H_x(x) dx &= f(U_i^{t,e}(x_{i+1/2})) - f(U_i^{t,e}(x_{i-1/2})) \\ &+ \int_{x_{i-1/2}}^{x_{i+1/2}} ((S(P_i^t(x)) - S(U_i^{t,e}(x))) H_x(x) dx, \end{aligned} \quad (1.2.10)$$

where  $U_i^{t,e}$  denotes again the stationary solution found in the first step of the exactly well-balanced reconstruction procedure in Algorithm 1.2.1, at the cell  $I_i$  and at time  $t$ , and  $P_i^t$  is the reconstruction operator defined as (1.1.33).

In the first and second order cases, applying the midpoint rule to (1.2.10) we obtain

$$\begin{aligned} \int_{x_{i-1/2}}^{x_{i+1/2}} S(P_i^t(x)) H_x(x) dx &= f(U_i^{t,e}(x_{i+1/2})) - f(U_i^{t,e}(x_{i-1/2})) \\ &+ \Delta x ((S(P_i^t(x_i)) - S(U_i^{t,e}(x_i))) H_x(x) dx \\ &= f(U_i^{t,e}(x_{i+1/2})) - f(U_i^{t,e}(x_{i-1/2})), \end{aligned}$$

since  $P_i^t(x_i) = U_i^{t,e}(x_i)$ .

### 1.3 Systems in Lagrangian coordinates

The use of Lagrangian coordinates in fluid dynamics partial differential equations may be of special interest since they are a good tool for tracking individual fluid particles as they move through the flow and this can allow the design of schemes with interesting properties.

In this thesis we will use them in Chapters 2 and 3 for obtaining implicit and semi-implicit schemes for the shallow water equations and the Ripa system.

We will introduce the Lagrangian coordinates through a simple example following [55] and considering the equations of gas dynamics in Lagrangian coordinates.

In Eulerian coordinates, we can consider the Euler equations for a compressible inviscid fluid (1.1.27).

We will now rewrite system (1.1.27) in Lagrangian coordinates, which follow the trajectories of the particles within the flow. In particular, consider any "fluid particle",  $\xi$ , and define the characteristic curves

$$\begin{cases} \frac{\partial x}{\partial t}(\xi, t) = u(x(\xi, t), t), \\ x(\xi, 0) = \xi. \end{cases} \quad (1.3.1)$$

Now, consider any function defined in Eulerian coordinates  $(x, t) \mapsto U(x, t)$ . We then define by

$$\bar{U}(\xi, t) = U(x(\xi, t), t)$$

its equivalent in the Lagrangian variables  $(\xi, t)$ .

Moreover, we define

$$L(\xi, t) = \frac{\partial x}{\partial \xi}(\xi, t),$$

the Jacobian of the Lagrangian map, which satisfies

$$\partial_t L(\xi, t) = \partial_\xi \bar{u}(\xi, t), \quad (1.3.2)$$

$$L(\xi, 0) = 1, \quad (1.3.3)$$

where  $\bar{u}(\xi, t) = u(x(\xi, t), t)$ .

Additionally, for all  $U$  we have the following relations between the partial derivatives:

$$\partial_\xi \bar{U} = L \bar{\partial}_x U, \quad (1.3.4)$$

$$\partial_t \bar{U} = \bar{\partial}_t U + u \bar{\partial}_x U. \quad (1.3.5)$$

We can then rewrite system (1.1.27) as:

$$\begin{cases} \partial_t(L\bar{\rho}) = 0, \\ \partial_t(L\bar{\rho}u) + \partial_\xi \bar{p} = 0, \\ \partial_t(L\bar{\rho}e) + \partial_\xi(\bar{p}u) = 0. \end{cases} \quad (1.3.6)$$

From now on we will suppress the overline notation for simplicity. From the first equation of (1.3.6) we deduce that

$$L\rho = \rho_0, \quad (1.3.7)$$

where  $\rho_0(\xi) = \rho(\xi, 0)$ .

Let us now introduce the variable

$$\tau = \frac{1}{\rho}, \quad (1.3.8)$$

which represents the specific volume. Then,

$$L = \rho_0 \tau \quad (1.3.9)$$

and (1.3.2) is equivalent to

$$\rho_0 \partial_t \tau - \partial_\xi u = 0. \quad (1.3.10)$$

Using again (1.3.9), the equations of conservation of momentum and energy can be written as:

$$\begin{aligned} \rho_0 \partial_t u + \partial_\xi p &= 0, \\ \rho_0 \partial_t e + \partial_\xi(pu) &= 0. \end{aligned}$$

Finally, introducing a mass variable  $m$  such that

$$dm = \rho_0 d\xi, \quad (1.3.11)$$

we can write the equations of gas dynamics with slab symmetry in Lagrangian coordinates as

$$\begin{cases} \partial_t \tau - \partial_m u = 0, \\ \partial_t u + \partial_m p = 0, \\ \partial_t e + \partial_m(pu) = 0. \end{cases} \quad (1.3.12)$$

As previously said, in the case of an isentropic flow it is enough to solve the system of conservation of mass and momentum:

$$\begin{cases} \partial_t \tau - \partial_m u = 0, \\ \partial_t u + \partial_m p = 0, \end{cases} \quad (1.3.13)$$

which is usually called the  $p$ -system. This system has two real distinct eigenvalues

$$\lambda_1 = -\sqrt{(-p')} < \lambda_2 = \sqrt{(-p')}$$

and is strictly hyperbolic given that  $p'(\tau) < 0$ .

In practice, instead of considering  $\partial_m$ , we will write everything in terms of  $\partial_\xi$ , so system (1.3.13) will be written as

$$\begin{cases} \partial_t \tau - \tau_0 \partial_\xi u = 0, \\ \partial_t u + \tau_0 \partial_\xi p = 0, \end{cases} \quad (1.3.14)$$

where  $\tau_0 = \frac{1}{\rho_0}$ .

Shallow water equations are actually equivalent to the particular case of Euler equations in which the flow is isentropic, the density is identified with the fluid height  $\rho \equiv h$  and the pressure takes the form  $p = \frac{1}{2}gh^2$ . In this thesis we will mainly focus on these equations in the non-flat topography case, which in Eulerian coordinates are given by (1.1.15).

Proceeding in a similar way to what has been done for the Euler equations, system (1.1.15) can be written in Lagrangian coordinates, in a similar way to (1.3.6), as follows:

$$\begin{cases} \partial_t(L\bar{h}) = 0, \\ \partial_t(L\bar{h}u) + \partial_\xi \bar{p} + g\bar{h}\partial_\xi \bar{z} = 0. \end{cases} \quad (1.3.15)$$

## Chapter 2

# Implicit and implicit-explicit Lagrange-projection exactly well-balanced finite volume schemes for 1D shallow water system

In this chapter we will focus on the design of finite volume schemes for the one-dimensional shallow water system. As discussed in Section 1.2, it is important that these schemes satisfy the well-balanced property. In fact, this has been the subject of study in many different works such as [4, 61, 77, 51, 1, 38, 7, 9, 2]. Here, we will focus on preserving the so-called water at rest steady states given by (1.1.19). That is, we want our schemes to preserve the steady states such that  $u = 0$ .

Nevertheless, when dealing with low Froude number situations, that is when  $Fr \ll 1$ , the time step due to the CFL condition makes explicit schemes inefficient (see (1.1.29)). Indeed, in such situations a large final time will be needed and performing small time steps will require a lot of iterations. To overcome this difficulty, implicit or implicit-explicit schemes allow the use of a larger time step and therefore less time iterations are needed. In [27, 28, 29, 30] implicit-explicit methods are proposed for three-dimensional shallow water flows. The extension into the DG framework to obtain high-order schemes is described in [46]. The technique is also extended for bedload sediment transport in [52]. Another possibility is to use implicit-explicit schemes (IMEX) which have successfully been applied to hyperbolic systems (see for instance [14, 13]). In [9] an implicit-explicit scheme is proposed for the shallow water flows in the low Froude number limit.

In this chapter, we propose the use of the so-called Lagrange-Projection decomposition in order to construct exactly well-balanced implicit finite volume schemes for system (1.1.15). A Lagrange-Projection type scheme is a two-step algorithm in which the system is first solved in Lagrangian coordinates (see Section 1.3) and the results are then projected into Eulerian coordinates. The first step is usually called the Lagrangian step and the

## 26 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

second one is referred to as the Projection or transport step.

In some works such as [23], these two systems to be solved were interpreted as an splitting of the original system. However, in later works such as [72] the authors consider the first system as the writing in Lagrangian coordinates of the original system, which facilitates obtaining properties such as high order or well-balancing. In addition, the use of the Lagrangian formulation in the first step is interesting since it simplifies the resulting scheme, which allows us to use simple Riemann solvers with good properties. Moreover, the use of the relaxation technique together with the use of Riemann invariants in the implicit system allows a simpler resolution of this step, avoiding the appearance of non-linearities in the pressure term.

The main purpose of going back to Eulerian coordinates is that using pure Lagrangian coordinates and keeping track of moving meshes may be cumbersome and many complex situations may arise for the configuration of the moving cells, especially thinking on the extension to 2D. Moreover, this strategy allows us to decouple the acoustic and transport phenomena related to our equations and to design implicit-explicit and large time step schemes in a natural way. Indeed it is very useful for approximating subsonic or low Froude number flows, where the usual CFL time step driven by the acoustic waves can be very restrictive. Using the implicit or implicit-explicit schemes means that the CFL restriction reduces only to the slow transport step rather the more restrictive acoustic one. This way, we will obtain implicit exactly well-balanced schemes for the shallow water equations that outperform the explicit ones.

The strategy proposed here follows the lines of [23], where an explicit fully well-balanced finite volume method is proposed in the Lagrange-Projection framework. In that case, only the explicit case was studied and its extension to an implicit scheme was hinted as future work. That strategy was then extended in [72] by proposing an explicit high-order Lagrange-Projection scheme, but again for the explicit case. Those two papers set the basis for this work, where we intend to describe implicit schemes based on such approach. Let us recall that in [34, 35, 36] one can find implicit-explicit schemes in the Lagrangian framework, where the source term is always treated explicitly. Those schemes were only first order accurate. One of the objectives here will be to extend the technique to second order implicit and implicit-explicit exactly well-balanced schemes. Moreover, we set the basis for their extension to higher order and fully exactly well-balanced schemes.

This chapter is organized as follows: in Section 2.1 the main ideas concerning the Lagrange-Projection techniques are introduced. Then, in Section 2.2 two new second order numerical schemes are proposed for the Lagrangian step based on an implicit or implicit-explicit approach. The schemes are exactly well-balanced for water at rest steady states. Next, in Section 2.3 the projection step is described. Finally, some numerical simulations are shown in Section 2.4 in order to study the accuracy and efficiency of the new schemes. For the sake of completeness, we include the description of the explicit scheme in 2.2.3.

## 2.1 The Lagrange-Projection strategy

Let us consider the shallow water equations in Lagrangian coordinates presented in (1.3.15) as well as the equivalent system (1.3.14). From the numerical point of view, it will be useful to consider a relaxation approach of the Lagrangian system (1.3.14) (see [36, 23]):

$$\begin{cases} \partial_t \tau - \tau_0 \partial_\xi u = 0, \\ \partial_t u + \tau_0 \partial_\xi \pi + g \tau_0 h \partial_\xi z = 0, \\ \partial_t \pi + a^2 \tau_0 \partial_\xi u = 0, \end{cases} \quad (2.1.1)$$

where  $a$  is a constant satisfying the subcharacteristic condition  $a > h\sqrt{gh}$  (see (1.1.65)) and  $\pi = p(\tau)$ .

Now, defining two new variables  $\vec{w} = \pi + au$  and  $\overleftarrow{w} = \pi - au$ , system (2.1.1) can be written as

$$\begin{cases} \partial_t \tau - \tau_0 \partial_\xi u = 0, \\ \partial_t \vec{w} + a \tau_0 \partial_\xi \vec{w} = -ga \tau_0 h \partial_\xi z, \\ \partial_t \overleftarrow{w} - a \tau_0 \partial_\xi \overleftarrow{w} = ga \tau_0 h \partial_\xi z. \end{cases} \quad (2.1.2)$$

The introduction of these two new variables allows one to obtain a system in which the second and third equations are just simple transport equations with a geometric source term.

Moreover, we can easily recover  $\pi$  and  $u$  from  $\vec{w}$  and  $\overleftarrow{w}$ :

$$\pi = \frac{\vec{w} + \overleftarrow{w}}{2}, \quad u = \frac{\vec{w} - \overleftarrow{w}}{2a}. \quad (2.1.3)$$

### 2.1.1 The Lagrange-Projection numerical algorithm

In order to state the Lagrange-Projection numerical scheme, we shall proceed as it is usually done for finite volume schemes (see Section 1.1.3), although we distinguish here between Lagrangian and Eulerian coordinates. Space (in Lagrangian framework) will be discretized by means of a fixed space step  $\Delta\xi$ . The cells  $I_i = [\xi_{i-1/2}, \xi_{i+1/2}]$  for  $i \in \mathbb{Z}$  are then considered and we define the set of times  $t^n = n\Delta t$  for  $n \in \mathbb{N}$ . In the Eulerian framework, the same space discretization will be used.

The Eulerian and Lagrangian space discretization will be then related using the following notation (see Figure 2.1):

$$x_i^*(t) = x(\xi_i, t), \quad x_{i+1/2}^*(t) = x(\xi_{i+1/2}, t).$$

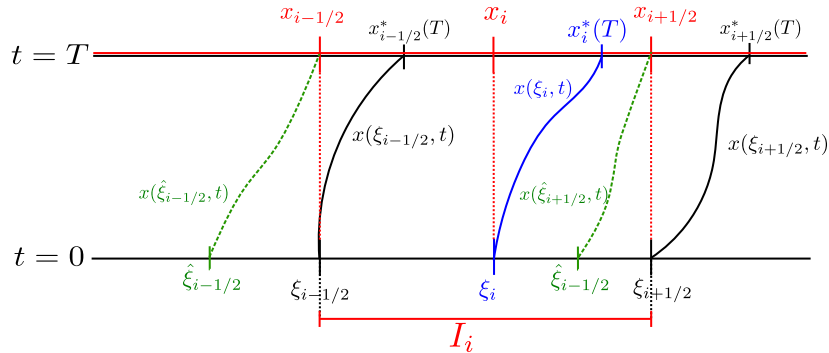


Figure 2.1: Sketch of the relation between Eulerian and Lagrangian coordinates

Remark that at time  $t = 0$  we then have  $x_i^*(0) = \xi_i = x_i$  and  $x_{i+1/2}^*(0) = \xi_{i+1/2} = x_{i+1/2}$ .

Let  $U = (h, hu)^T$ . Assume that the initial condition  $x \mapsto U^0(x)$  is given. Then, define the discrete initial data  $U_i^0$ , where

$$U_i^0 \approx \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U^0(x) dx, \quad i \in \mathbb{Z}.$$

The objective is to compute the values

$$U_i^n \approx \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U(x, t^n) dx,$$

where  $x \mapsto U(x, t^n)$  corresponds to the solution of (1.1.15) at time  $t^n$ ,  $n \in \mathbb{N}$ .

The Lagrange-Projection consists on the following two-step process:

1. Lagrangian step: update  $U_i^n$  to  $\bar{U}_i^{n+1}$  by numerically solving (1.3.15);
2. Projection step: update  $\bar{U}_i^{n+1}$  to  $U_i^{n+1}$  by projecting back to Eulerian framework.

## 2.2 The Lagrangian step

At this step we need to solve the shallow water system in Lagrangian coordinates (1.3.15). Although the aim here is to propose an implicit approximation of this system, we shall introduce, for the sake of completeness, an explicit approximation as well in 2.2.3. We will present first and second order schemes. For the implicit case, we will propose two versions of the scheme: nonlinear implicit and implicit-explicit.

Let us start by working on the second equation of system (1.3.15). The last term of this equation may be rewritten as

$$\overline{gh\partial_\xi z} = gL\overline{h\partial_x z} = gh_0\overline{\partial_x z}, \quad (2.2.1)$$

where we have used the first equation in (1.3.15) and denoted  $h_0 = h|_{t=0}$ . Now, since  $z$  does not depend on time,

$$\overline{\partial_x z}(\xi_i, t) = \partial_x z(x(\xi_i, t)) = z'(x_i(t)),$$

where we have used the notation  $z'$  for  $\partial_x z$ . Therefore, the second equation in (1.3.15) is equivalent to

$$\partial_t(L\overline{hu}) + \partial_\xi \overline{\pi} + gh_0 z'(x_i(t)) = 0.$$

Now, integrating (1.3.15) in the interval  $I_i = [\xi_{i-1/2}, \xi_{i+1/2}]$  we can write a semi-discrete scheme for the system in Lagrangian coordinates (1.3.15):

$$\begin{cases} (L\overline{h})'_i(t) = 0, \\ (L\overline{hu})'_i(t) = -\frac{1}{\Delta\xi} \left( \pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t) \right) - \frac{1}{\Delta\xi} \int_{\xi_{i-1/2}}^{\xi_{i+1/2}} gP_{i,h_0}(\xi) z'(x(\xi, t)) d\xi, \end{cases} \quad (2.2.2)$$

where we set  $\overline{h}_i(0) = h_i^n$ , and  $(\overline{hu})_i(0) = (hu)_i^n$  as initial conditions. In this system,  $\pi_{i\pm\frac{1}{2}}^*(t) \approx \overline{\pi}(\xi_{i\pm\frac{1}{2}}, t)$  and  $P_{i,h_0}(\xi)$  is a reconstruction operator obtained from the sequence of cell values  $\{h_i^n\}$ :

$$P_{i,h_0}(\xi) = P_{i,h_0}(\xi; \{h_j^n\}_{j \in \mathcal{S}_i}),$$

being  $\mathcal{S}_i$  the set of indexes belonging to the stencil corresponding to the  $i$ -th cell.

Now we will consider a semi-discretization of the relaxed system (2.1.2), that will play an important role in the definition of the discretization of the Lagrangian formulation of the SWE, since the values  $\overrightarrow{w}$  and  $\overleftarrow{w}$  will give us the approximations that we need for  $u^*$  and  $\pi^*$ .

Using (2.2.1), system (2.1.2) can also be discretized in space using a first or second order scheme as follows:

$$\begin{cases} \tau'_i(t) = \frac{1}{h_{0,i}\Delta\xi} \int_{\xi_{i-1/2}}^{\xi_{i+1/2}} P_{i,u}(\xi, t) d\xi, \\ \overrightarrow{w}'_i(t) = -\frac{a}{h_{0,i}\Delta\xi} (\overrightarrow{w}_{i+1/2}(t) - \overrightarrow{w}_{i-1/2}(t)) - \frac{ga}{\Delta\xi} \int_{\xi_{i-1/2}}^{\xi_{i+1/2}} z'(x(\xi, t)) d\xi, \\ \overleftarrow{w}'_i(t) = \frac{a}{h_{0,i}\Delta\xi} (\overleftarrow{w}_{i+1/2}(t) - \overleftarrow{w}_{i-1/2}(t)) + \frac{ga}{\Delta\xi} \int_{\xi_{i-1/2}}^{\xi_{i+1/2}} z'(x(\xi, t)) d\xi, \end{cases} \quad (2.2.3)$$

### 30 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

with initial conditions  $\tau_i(0) = \frac{1}{h_{0,i}}$ ,  $\vec{w}_i(0) = \frac{1}{2}gh_{0,i}^2 + au_{0,i}$  and  $\overleftarrow{w}_i(0) = \frac{1}{2}gh_{0,i}^2 - au_{0,i}$ .  $P_{i,u}(\xi, t)$  denotes the reconstruction operators of  $u$  defined from the cell values  $\{u_i(t)\}$  on a given stencil. Finally,  $\vec{w}_{i+1/2}(t)$  and  $\overleftarrow{w}_{i+1/2}(t)$  are the numerical fluxes of the second and third equations of system (2.1.2), for which we will consider the upwind numerical fluxes given by:

$$\begin{aligned}\vec{w}_{i+1/2}(t) &= P_{i,\vec{w}}(\xi_{i+1/2}, t) = \vec{w}_{i+1/2-}^t, \\ \overleftarrow{w}_{i+1/2}(t) &= P_{i+1,\overleftarrow{w}}(\xi_{i+1/2}, t) = \overleftarrow{w}_{i+1/2+}^t,\end{aligned}$$

where  $P_{i,\vec{w}}$  and  $P_{i+1,\overleftarrow{w}}$  are their respective reconstruction operators.

Note that there is no need to solve the first equation in (2.2.3), since the three equations are decoupled and only the values of  $\vec{w}$  and  $\overleftarrow{w}$  will be needed to perform the Lagrangian and the projection step.

Finally, making use of the relations (2.1.3), we can define the following numerical fluxes for  $\pi$  and  $u$  at the interfaces:

$$\begin{aligned}\pi_{i+1/2}^*(t) &= \frac{P_{i,\vec{w}}(\xi_{i+1/2}, t) + P_{i+1,\overleftarrow{w}}(\xi_{i+1/2}, t)}{2}, \\ u_{i+1/2}^*(t) &= \frac{P_{i,\vec{w}}(\xi_{i+1/2}, t) - P_{i+1,\overleftarrow{w}}(\xi_{i+1/2}, t)}{2a}.\end{aligned}\tag{2.2.4}$$

As an example, let us now consider a first order in space and time implicit finite volume scheme for the SWE in Lagrangian coordinates that uses the previous formalism. In that case, the system is written as

$$\begin{aligned}\vec{w}_i^{n+1} &= \vec{w}_i^n - \frac{a\Delta t}{h_i^n \Delta \xi} \left( \vec{w}_{i+1/2-}^{n+1} - \vec{w}_{i-1/2-}^{n+1} \right) - \frac{a\Delta t}{\Delta \xi} g z' (x_i^{*,n+1}), \\ \overleftarrow{w}_i^{n+1} &= \overleftarrow{w}_i^n + \frac{a\Delta t}{h_i^n \Delta \xi} \left( \overleftarrow{w}_{i+1/2+}^{n+1} - \overleftarrow{w}_{i-1/2+}^{n+1} \right) + \frac{a\Delta t}{\Delta \xi} g z' (x_i^{*,n+1}),\end{aligned}\tag{2.2.5}$$

where we set  $\vec{w}_{i+1/2-}^{n+1} = \vec{w}_i^{n+1}$  and  $\overleftarrow{w}_{i+1/2+}^{n+1} = \overleftarrow{w}_{i+1}^{n+1}$ . That is, the reconstruction operators reduce to the standard cell values at time  $t^{n+1}$ .

It should be noted that the presence of the source term requires the evaluation of  $z'$  at a point that we have denoted as  $x_i^{*,n+1}$ , which corresponds to a first order approximation of  $x(\xi_i, t^{n+1})$ . Here we propose the following approximation

$$x_i^{*,n+1} = \xi_i + \Delta t u_i^{*,n+1},\tag{2.2.6}$$

where

$$u_i^{*,n+1} = \frac{u_{i-1/2}^{*,n+1} + u_{i+1/2}^{*,n+1}}{2}\tag{2.2.7}$$

and  $u_{i\pm 1/2}^{*,n+1}$  can be defined using the relations (2.2.4) as follows:

$$u_{i+1/2}^{*,n+1} = \frac{\overrightarrow{w}_i^{n+1} - \overleftarrow{w}_{i+1}^{n+1}}{2a}. \quad (2.2.8)$$

Note that (2.2.5)-(2.2.8) define a coupled non-linear system to be solved. Here we use a simple fixed-point algorithm where we first solve (2.2.5), fixing the value  $x_i^{*,n+1}$  and then we update it by computing the velocity at the interfaces using (2.2.8). Observe that (2.2.5) reduces to two uncoupled linear systems that are simple to solve when  $x_i^{*,n+1}$  is fixed. We have no theoretical proof of the convergence of this fixed-point algorithm. However, as it will be shown in the numerical results, no convergence problems have been found in practice.

Once  $\overrightarrow{w}_i^{n+1}$  and  $\overleftarrow{w}_i^{n+1}$  are computed,  $(L\overline{hu})_i^{n+1}$  would be obtained as follows:

$$(L\overline{hu})_i^{n+1} = (hu)_i^n - \frac{\Delta t}{\Delta \xi} \left( \pi_{i+1/2}^{*,n+1} - \pi_{i-1/2}^{*,n+1} \right) - g\Delta t h_i^n z' \left( x_i^{*,n+1} \right), \quad (2.2.9)$$

where  $\pi_{i+1/2}^{*,n+1}$  is computed using the relations (2.2.4), which reduces to

$$\pi_{i+1/2}^{*,n+1} = \frac{\overrightarrow{w}_i^{n+1} + \overleftarrow{w}_{i+1}^{n+1}}{2} \quad (2.2.10)$$

for a first order numerical scheme.

Finally, we define

$$L_i^{n+1} = 1 + \frac{\Delta t}{\Delta \xi} \left( u_{i+1/2}^{*,n+1} - u_{i-1/2}^{*,n+1} \right). \quad (2.2.11)$$

Notice that we could avoid the solution of the non-linear system (2.2.5)-(2.2.8) by considering the following explicit-implicit first order numerical scheme

$$\begin{aligned} \overrightarrow{w}_i^{n+1} &= \overrightarrow{w}_i^n - \frac{a\Delta t}{h_i^n \Delta \xi} \left( \overrightarrow{w}_{i+1/2-}^{n+1} - \overrightarrow{w}_{i-1/2-}^{n+1} \right) - \frac{a\Delta t}{\Delta \xi} g z' \left( x_i^n \right), \\ \overleftarrow{w}_i^{n+1} &= \overleftarrow{w}_i^n + \frac{a\Delta t}{h_i^n \Delta \xi} \left( \overleftarrow{w}_{i+1/2+}^{n+1} - \overleftarrow{w}_{i-1/2+}^{n+1} \right) + \frac{a\Delta t}{\Delta \xi} g z' \left( x_i^n \right). \end{aligned} \quad (2.2.12)$$

This way, we only have to solve one system for  $\overrightarrow{w}^{n+1}$  and another one for  $\overleftarrow{w}^{n+1}$ , avoiding the use of fixed point iterations.

Once  $\overleftarrow{w}_i^{n+1}$  and  $\overrightarrow{w}_i^{n+1}$  are computed, then  $(L\overline{hu})_i^{n+1}$  is obtained as follows:

$$(L\overline{hu})_i^{n+1} = (hu)_i^n - \frac{\Delta t}{\Delta \xi} \left( \pi_{i+1/2}^{*,n+1} - \pi_{i-1/2}^{*,n+1} \right) - g\Delta t h_i^n z' \left( x_i^n \right). \quad (2.2.13)$$

Finally,  $L_i^{n+1}$  is updated using (2.2.11).

Fully explicit versions of the previous numerical scheme are straightforward and their extension to second order is also possible (see Section 2.2.3).

## 32 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

Let us remark that, in order to ensure stability, the parameter  $a$  must be chosen sufficiently large according to the so-called subcharacteristic condition  $a > h\sqrt{gh}$  (see (1.1.65)). Moreover, the explicit Lagrangian step is stable provided the following CFL condition is satisfied

$$a\Delta t \leq \frac{1}{2}h\Delta\xi. \quad (2.2.14)$$

The use of an implicit approach for the Lagrangian step makes that condition (2.2.14) is no longer required.

Unfortunately, neither the numerical scheme defined by (2.2.5)-(2.2.9) nor (2.2.12)-(2.2.13) are well-balanced for the water at rest solution. In order to define exactly well-balanced numerical schemes for the water at rest solution we follow the procedure described in [25], detailed in Algorithm 1.2.2. More explicitly, at every cell  $I_i = [\xi_{i-1/2}, \xi_{i+1/2}]$  we look for a steady state that matches with the given cell average. As we are interested in first and second order schemes and we focus on water at rest steady states, this condition reduces to defining an in-cell stationary water height  $h_i^{e,n}(\xi)$  such that  $h_i^{e,n}(\xi_i) = h_i^n$ . Taking into account the special expression of the water at rest solutions for the SWE,  $h_i^{e,n}(\xi)$  is given by

$$h_i^{e,n}(\xi) = h_i^n + z(\xi_i) - z(\xi). \quad (2.2.15)$$

We also define  $\pi_i^{e,n}$  as

$$\pi_i^{e,n}(\xi) = \frac{1}{2}g(h_i^{e,n}(\xi))^2. \quad (2.2.16)$$

Notice that  $h_i^{e,n}(\xi)$  exactly satisfies

$$\partial_\xi \overline{\pi_i^{e,n}} + gL\overline{h_i^{e,n} z'} = 0. \quad (2.2.17)$$

Now, the second equation of (2.2.2) could be rewritten equivalently as follows:

$$\begin{aligned} (L\overline{hu})'_i(t) &= -\frac{1}{\Delta\xi} (\pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t)) + \frac{1}{\Delta\xi} (\pi_i^{e,n}(x_{i+1/2}(t)) - \pi_i^{e,n}(x_{i-1/2}(t))) \quad (2.2.18) \\ &\quad - \frac{1}{\Delta\xi} \int_{\xi_{i-1/2}}^{\xi_{i+1/2}} g(P_{i,h_0}(\xi) - L(\xi, t)h_i^{e,n}(x(\xi, t))) z'(x(\xi, t))d\xi, \end{aligned}$$

where  $x_{i+1/2}(t) = x(\xi_{i+1/2}, t)$ .

Proceeding analogously with the second and third equations of (2.2.3), one has:

$$\left\{ \begin{array}{l} \vec{w}'_i(t) = -\frac{a}{h_{0,i}\Delta\xi} (\vec{w}_{i+1/2}(t) - \vec{w}_{i-1/2}(t)) \\ \quad + \frac{a}{h_{0,i}\Delta\xi} (\pi_i^{e,n}(x_{i+1/2}(t)) - \pi_i^{e,n}(x_{i-1/2}(t))) \\ \quad - \frac{ga}{\Delta\xi} \int_{\xi_{i-1/2}}^{\xi_{i+1/2}} \left( 1 - \frac{L(\xi, t)h_i^{e,n}(x(\xi, t))}{P_{i,h_0}(\xi)} \right) z'(x(\xi, t))d\xi, \\ \overleftarrow{w}'_i(t) = \frac{a}{h_{0,i}\Delta\xi} (\overleftarrow{w}_{i+1/2}(t) - \overleftarrow{w}_{i-1/2}(t)) \\ \quad - \frac{a}{h_{0,i}\Delta\xi} (\pi_i^{e,n}(x_{i+1/2}(t)) - \pi_i^{e,n}(x_{i-1/2}(t))) \\ \quad + \frac{ga}{\Delta\xi} \int_{\xi_{i-1/2}}^{\xi_{i+1/2}} \left( 1 - \frac{L(\xi, t)h_i^{e,n}(x(\xi, t))}{P_{i,h_0}(\xi)} \right) z'(x(\xi, t))d\xi. \end{array} \right. \quad (2.2.19)$$

Note that the semi-discrete scheme described in (2.2.18) and (2.2.19) is equivalent to (2.2.2) and (2.2.3). Now, we need to describe the procedure to define the reconstructions operators of the unknowns that are exactly well-balanced in the sense defined in [56].

### 2.2.1 Exactly well-balanced reconstruction operators

In the next paragraphs we describe the reconstruction procedure that we define in order to achieve a numerical scheme that is exactly well-balanced for the water at rest solutions.

Let us suppose that  $\{h_i^n\}$  and  $\{(hu)_i^n\}$  are known. We shall consider reconstruction operator for variables  $h_0$ ,  $\vec{w}$  and  $\overleftarrow{w}$  denoted as  $P_{i,h_0}$ ,  $P_{i,\overleftarrow{w}}^n$  and  $P_{i,\vec{w}}^n$  and defined as follows:

$$P_{i,h_0}(\xi) = h_i^{e,n}(\xi) + Q_i^n(\xi; \{h_j^{n,f}\}_{j \in \mathcal{S}_i}), \quad (2.2.20)$$

where  $h_i^{e,n}(\xi)$  is given by (2.2.15),  $h_j^{n,f} = h_j^n - h_i^{e,n}(\xi_j)$ ,  $j \in \mathcal{S}_i$ , that is  $h_j^{n,f}$  are the fluctuations with respect to the steady state at the stencil, and  $Q_i^n(\xi)$  is a standard first or second order reconstruction operator. Note that as we are only interested in first and second order numerical methods, the fluctuations are computed using punctual evaluations at the cell centers. If order higher than two were required, then cell averages would be computed with a quadrature formula.

Now,  $P_{i,\vec{w}}^n$  is defined equally as

$$P_{i,\vec{w}}^n(\xi) = \pi_i^{e,n}(\xi) + Q_i^n(\xi; \{\vec{w}_j^{n,f}\}_{j \in \mathcal{S}_i}), \quad (2.2.21)$$

where  $\vec{w}_j^{n,f}$  is given by  $\vec{w}_j^{n,f} = \vec{w}_j^n - \pi_i^{e,n}(\xi_j)$ ,  $j \in \mathcal{S}_i$ , that is, they are again the fluctuations with respect to the local steady state at the stencil.  $P_{i,\overleftarrow{w}}^n$  is defined analogously.

### 34 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

Finally,  $P_{i,u}^n$  is a standard first or second order reconstruction operator computed with the cell values  $\{u_i^n\}$ . Note that as we are only interested in water at rest steady state, no special care has to be taken for the reconstruction operator of  $u$  or  $hu$ . This would not be the case for moving equilibria. Moreover, if order higher than two were required, then a more sophisticated procedure should be used as the cell averages of  $u$  cannot be obtained by simply setting  $u_i^n = \frac{(hu)_i^n}{h_i^n}$ .

Observe that  $P_{i,h_0}$ ,  $P_{i,\overleftarrow{w}}^n$ ,  $P_{i,\overrightarrow{w}}^n$  and  $P_{i,u}^n$  are exactly well-balanced operators in the sense that if  $\{h_i^n\}$  and  $\{(hu)_i^n\}$  are the cell averages of a given water at rest steady state that are computed with the mid-point rule, that is  $h_i^n = h^e(\xi_i)$ ,  $(hu)_i^n = 0$ , then

$$P_{i,h_0}(\xi) = h^e(\xi)|_{I_i}, \quad P_{i,\overleftarrow{w}}^n(\xi) = P_{i,\overrightarrow{w}}^n(\xi) = \pi^e(\xi)|_{I_i}, \quad \text{and} \quad P_{i,u}^n(\xi) = 0.$$

Given our interest in the definition of implicit solvers, we also need to provide a reconstruction procedure for any value  $t \in [t^n, t^{n+1}]$  that is able to preserve water at rest steady states. In particular we need to define  $P_{i,\overleftarrow{w}}(\xi, t)$  and  $P_{i,\overrightarrow{w}}(\xi, t)$ . Let us define  $P_{i,\overrightarrow{w}}(\xi, t)$  and the definition of  $P_{i,\overleftarrow{w}}(\xi, t)$  will be analogous.  $P_{i,\overrightarrow{w}}(\xi, t)$  will be defined using the exactly well-balanced reconstruction operator at time  $t = t^n$ ,  $P_{i,\overrightarrow{w}}^n$ , described previously, and a standard reconstruction operator for the time fluctuations, that is:

$$P_{i,\overrightarrow{w}}(\xi, t) = P_{i,\overrightarrow{w}}^n(\xi) + \tilde{Q}_i(\xi, t), \quad (2.2.22)$$

where  $\tilde{Q}_i(\xi, t)$  is a standard reconstruction operator defined in terms of the time fluctuations, that is

$$\tilde{Q}_i(\xi, t) = \tilde{Q}_i(\xi; \{\overrightarrow{w}_j^{t,f}\}_{j \in \mathcal{S}_i}), \quad \text{where} \quad \overrightarrow{w}_j^{t,f} = \overrightarrow{w}_j(t) - \overrightarrow{w}_j^n, \quad j \in \mathcal{S}_i.$$

Let us show a first order reconstruction operator for  $\overrightarrow{w}$  at cell  $I_i$  for any value  $t \in [t^n, t^{n+1}]$

$$P_{i,\overrightarrow{w}}^{o1}(\xi, t) = P_{i,\overrightarrow{w}}^{n,o1}(\xi) + \tilde{Q}_i^{o1}(\xi, t),$$

where  $\tilde{Q}_i^{o1}(\xi, t)$  is the standard first order reconstruction operator at cell  $I_i$ , that is,  $\tilde{Q}_i^{o1}(\xi, t) = \overrightarrow{w}_i^{t,f} = \overrightarrow{w}_i(t) - \overrightarrow{w}_i^n$ . Now, taking into account the definition of  $P_{i,\overrightarrow{w}}^n(\xi)$ , we have that

$$P_{i,\overrightarrow{w}}^{o1}(\xi, t) = \pi_i^{e,n}(\xi) - \pi_i^{e,n}(\xi_i) + \overrightarrow{w}_i(t). \quad (2.2.23)$$

The second order reconstruction operator is defined in the same way, that is

$$P_{i,\overrightarrow{w}}^{o2}(\xi, t) = P_{i,\overrightarrow{w}}^{n,o2}(\xi) + \tilde{Q}_i^{o2}(\xi, t), \quad (2.2.24)$$

where  $\tilde{Q}_i^{o2}(\xi, t)$  is a standard second order reconstruction operator defined in terms of the time fluctuations  $\{\overrightarrow{w}_j^{t,f}\}_{j \in \mathcal{S}_i}$ . Now, taking into account the definition of  $P_{i,\overrightarrow{w}}^{n,o2}$  we have that

$$P_{i,\overrightarrow{w}}^{o2}(\xi, t) = \pi_i^{e,n}(\xi) + Q_i^{n,o2}(\xi; \{\overrightarrow{w}_j^{n,f}\}_{j \in \mathcal{S}_i}) + \tilde{Q}_i^{o2}(\xi, t; \{\overrightarrow{w}_j^{t,f}\}_{j \in \mathcal{S}_i}),$$

where  $\vec{w}_j^{n,f} = \vec{w}_j^n - \pi_i^{e,n}(\xi_j)$ . In the previous expression,  $Q_i^{n,o2}$  and  $\tilde{Q}_i^{o2}$  are standard second order reconstruction operators. In this case, to avoid the non-linear dependency of the limiters at any time  $t$ , we consider here that both  $Q_i^{n,o2}$  and  $\tilde{Q}_i^{o2}$  use the same limiters computed at time  $t = t^n$ .

In particular,  $Q_i^{n,o2}(\xi)$  is defined as follows:

$$Q_i^{n,o2}(\xi) = \vec{w}_i^n - \pi_i^{e,n}(\xi_i) + \Delta^n \vec{w}_i^{n,f}(\xi - \xi_i), \quad (2.2.25)$$

where

$$\Delta^n \vec{w}_i^{n,f} = \frac{1}{\Delta\xi} \left( \phi_{i+}^n \left( \vec{w}_i^{n,f} - \vec{w}_{i-1}^{n,f} \right) + \phi_{i-}^n \left( \vec{w}_{i+1}^{n,f} - \vec{w}_i^{n,f} \right) \right), \quad (2.2.26)$$

with  $\phi_{i-}^n$  and  $\phi_{i+}^n$  slope limiters. Here we use the following:

$$\phi_{i-}^n = \begin{cases} \frac{|d_{i-}|}{|d_{i-}| + |d_{i+}|} & \text{if } |d_{i-}| + |d_{i+}| > 0, \\ 0 & \text{otherwise,} \end{cases}$$

$$\phi_{i+}^n = \begin{cases} \frac{|d_{i+}|}{|d_{i-}| + |d_{i+}|} & \text{if } |d_{i-}| + |d_{i+}| > 0, \\ 0 & \text{otherwise,} \end{cases}$$

where

$$d_{i-} = \vec{w}_i^{n,f} - \vec{w}_{i-1}^{n,f}, \quad d_{i+} = \vec{w}_{i+1}^{n,f} - \vec{w}_i^{n,f}.$$

$\tilde{Q}_i^{o2}(\xi, t)$  is defined as follows

$$\tilde{Q}_i^{o2}(\xi, t) = \vec{w}_i(t) - \vec{w}_i^n + \Delta^t \vec{w}_i^{t,f}(\xi - \xi_i),$$

where

$$\Delta^t \vec{w}_i^{t,f} = \frac{1}{\Delta\xi} \left( \tilde{\phi}_{i-}^n \left( \vec{w}_{i+1}^{t,f} - \vec{w}_i^{t,f} \right) + \tilde{\phi}_{i+}^n \left( \vec{w}_i^{t,f} - \vec{w}_{i-1}^{t,f} \right) \right),$$

with  $\tilde{\phi}_{i\pm}^n = \phi_{i\pm}^n$ .

Taking into account the definitions of  $Q_i^{n,o2}$  and  $\tilde{Q}_i^{t,o2}$ , we obtain that

$$P_{i,\vec{w}}^{o2}(\xi, t) = \pi_i^{e,n}(\xi) - \pi_i^{e,n}(\xi_i) + \vec{w}_i(t) + \Delta^n \vec{w}_i^{n,f}(\xi - \xi_i) + \Delta^t \vec{w}_i^{t,f}(\xi - \xi_i). \quad (2.2.27)$$

## 2.2.2 Implicit and implicit-explicit exactly well-balanced Lagrangian schemes

We are now going to describe the implicit and implicit-explicit first and second order exactly well-balanced Lagrangian schemes. In the implicit and implicit-explicit case we use the mid-point rule for approximating the integrals in (2.2.18) and (2.2.19) and we obtain the following semi-discrete formulation:

$$(L\bar{hu})'_i(t) = -\mathcal{L}_i(t) - \mathcal{G}_i(t) \quad (2.2.28)$$

where

$$\mathcal{L}_i(t) = \frac{1}{\Delta\xi} (\pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t)) \quad (2.2.29)$$

$$\begin{aligned} \mathcal{G}_i(t) &= -\frac{1}{\Delta\xi} (\pi_i^{e,n}(x_{i+1/2}(t)) - \pi_i^{e,n}(x_{i-1/2}(t))) \\ &\quad + g(P_{i,h_0}(\xi_i) - L_i(t)h_i^{e,n}(x(\xi_i, t))) z'(x(\xi_i, t)) \end{aligned} \quad (2.2.30)$$

and

$$\begin{cases} \vec{w}'_i(t) = -(\mathcal{L}_{\vec{w}})_i(t) - \frac{a}{h_{0,i}}\mathcal{G}_i(t), \\ \overleftarrow{w}'_i(t) = -(\mathcal{L}_{\overleftarrow{w}})_i(t) + \frac{a}{h_{0,i}}\mathcal{G}_i(t), \end{cases} \quad (2.2.31)$$

where

$$(\mathcal{L}_{\vec{w}})_i(t) = \frac{a}{h_{0,i}\Delta\xi} (\vec{w}_{i+1/2}(t) - \vec{w}_{i-1/2}(t)) \quad (2.2.32)$$

$$(\mathcal{L}_{\overleftarrow{w}})_i(t) = -\frac{a}{h_{0,i}\Delta\xi} (\overleftarrow{w}_{i+1/2}(t) - \overleftarrow{w}_{i-1/2}(t)), \quad (2.2.33)$$

where we recall that  $x_{i\pm 1/2}(t) = x(\xi_{i\pm 1/2}, t)$ .

Now, we are ready to start the definition of first and second order exactly well-balanced implicit and implicit-explicit schemes in Lagrangian coordinates. Here we suppose that  $\{h_i^n\}$  and  $\{(hu)_i^n\}$  are known and we integrate the system in the interval  $[t^n, t^{n+1}]$  to compute the new states at  $t^{n+1}$ .

We are going to consider first and second order schemes. In each case, implicit and implicit-explicit schemes will be considered, in which the operators  $\mathcal{L}$  and  $\mathcal{G}$  will be evaluated in each case at the appropriate time as shown in Section (1.1.3.3).

### 2.2.2.1 First order schemes

**Implicit scheme** Firstly, let us consider the definition of the first order implicit scheme. In that case, the scheme reads as follows:

$$(L\bar{hu})_i^{n+1} = (hu)_i^n - \Delta t (\mathcal{L}_i^{n+1} + \mathcal{G}_i^{n+1}), \quad (2.2.34)$$

where

$$\mathcal{L}_i^{n+1} = \frac{1}{\Delta\xi} (\pi_{i+1/2}^{*,n+1} - \pi_{i-1/2}^{*,n+1}), \quad (2.2.35)$$

and

$$\mathcal{G}_i^{n+1} = -\frac{1}{\Delta\xi} (\pi_i^{e,n}(x_{i+1/2}^{*,n+1}) - \pi_i^{e,n}(x_{i-1/2}^{*,n+1})) + g(h_i^n - L_i^{n+1}h_i^{e,n}(x_i^{*,n+1})) z'(x_i^{*,n+1}), \quad (2.2.36)$$

with  $L_i^{n+1}$  defined as in (2.2.11),  $x_i^{*,n+1}$  defined as in (2.2.6) and  $x_{i\pm 1/2}^{*,n+1}$  defined as

$$x_{i\pm 1/2}^{*,n+1} = \xi_{i\pm 1/2} + \Delta t u_{i\pm 1/2}^{*,n+1}. \quad (2.2.37)$$

We recall that  $\pi_{i\pm 1/2}^{*,n+1}$  and  $u_{i\pm 1/2}^{*,n+1}$  are defined by (2.2.4), that is,

$$\pi_{i+1/2}^{*,n+1} = \frac{P_{i,\vec{w}}^{o1}(\xi_{i+1/2}, t^{n+1}) + P_{i+1,\overleftarrow{w}}^{o1}(\xi_{i+1/2}, t^{n+1})}{2} = \frac{\vec{w}_{i+1/2-}^{n+1} + \overleftarrow{w}_{i+1/2+}^{n+1}}{2} \quad (2.2.38)$$

and

$$u_{i+1/2}^{*,n+1} = \frac{P_{i,\vec{w}}^{o1}(\xi_{i+1/2}, t^{n+1}) - P_{i+1,\overleftarrow{w}}^{o1}(\xi_{i+1/2}, t^{n+1})}{2a} = \frac{\vec{w}_{i+1/2-}^{n+1} - \overleftarrow{w}_{i+1/2+}^{n+1}}{2a} \quad (2.2.39)$$

where

$$\begin{aligned} \vec{w}_{i+1/2-}^{n+1} &= P_{i,\vec{w}}^{o1}(\xi_{i+1/2}, t^{n+1}) = \pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_i) + \vec{w}_i^{n+1}, \\ \overleftarrow{w}_{i+1/2+}^{n+1} &= P_{i+1,\overleftarrow{w}}^{o1}(\xi_{i+1/2}, t^{n+1}) = \pi_{i+1}^{e,n}(\xi_{i+1/2}) - \pi_{i+1}^{e,n}(\xi_i) + \overleftarrow{w}_{i+1}^{n+1}. \end{aligned} \quad (2.2.40)$$

Note that  $\pi_{i\pm 1/2}^{*,n+1}$  and  $u_{i\pm 1/2}^{*,n+1}$  are defined in terms of  $\vec{w}^{n+1}$  and  $\overleftarrow{w}^{n+1}$ , that are the solutions of the non-linear system

$$\begin{cases} \vec{w}_i^{n+1} &= \vec{w}_i^n - \Delta t \left( (\mathcal{L}_{\vec{w}})_{i+1/2}^{n+1} + \frac{a}{h_i^n} \mathcal{G}_i^{n+1} \right), \\ \overleftarrow{w}_i^{n+1} &= \overleftarrow{w}_i^n - \Delta t \left( (\mathcal{L}_{\overleftarrow{w}})_{i-1/2}^{n+1} - \frac{a}{h_i^n} \mathcal{G}_i^{n+1} \right), \end{cases} \quad (2.2.41)$$

where

$$\begin{aligned} (\mathcal{L}_{\vec{w}})_{i+1/2}^{n+1} &= \frac{a}{h_i^n \Delta \xi} \left( \vec{w}_{i+1/2-}^{n+1} - \vec{w}_{i-1/2-}^{n+1} \right), \\ (\mathcal{L}_{\overleftarrow{w}})_{i-1/2}^{n+1} &= -\frac{a}{h_i^n \Delta \xi} \left( \overleftarrow{w}_{i+1/2+}^{n+1} - \overleftarrow{w}_{i-1/2+}^{n+1} \right). \end{aligned} \quad (2.2.42)$$

In order to apply the previous numerical scheme, we first solve the non-linear system (2.2.41), and then, (2.2.34) is updated using  $\{\vec{w}_i^{n+1}\}$  and  $\{\overleftarrow{w}_i^{n+1}\}$ .

The non-linear system (2.2.41) is solved by means of a fixed-point iteration where  $\{\vec{w}_i^n\}$  and  $\{\overleftarrow{w}_i^n\}$  are given as initial guess. Here (2.2.41) is solved as follows: we fix the values of  $x_{i\pm 1/2}^{*,n+1}$  and  $x_i^{*,n+1}$  in  $(\mathcal{G}_w)_i^{n+1}$  and then we solve the two linear systems in (2.2.41). Next  $u_{i\pm 1/2}^{*,n+1}$  are updated as well as  $x_{i\pm 1/2}^{*,n+1}$  and  $x_i^{*,n+1}$  with the new computed values of  $\{\overleftarrow{w}_i^{n+1}\}$  and  $\{\vec{w}_i^{n+1}\}$ .

**Theorem 2.2.1.** *Given any water at rest stationary solution of system (1.1.15) with a continuous bottom topography  $z$ , the scheme defined by (2.2.34) and (2.2.41) preserves it and the scheme is exactly well-balanced.*

### 38 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

*Proof.* Let  $h_i^n = h^{e,n}(\xi_i)$ ,  $\pi_i^n = \pi^{e,n}(\xi_i)$ ,  $u_i^n = 0$ . In the first iteration of the fixed-point algorithm,  $x_{i\pm 1/2}^{*,n+1} = \xi_{i\pm 1/2}$ ,  $x_i^{*,n+1} = \xi_i$ . Let  $\overrightarrow{w}_i^{n+1,l}$  be the result of the  $l$ -th iteration of the fixed-point algorithm. Then,

$$\begin{aligned} \overrightarrow{w}_i^{n+1,1} &= \overrightarrow{w}_i^n - \frac{\Delta ta}{h_i^n \Delta \xi} (\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_i) + \overrightarrow{w}_i^{n+1,1} - \pi_{i-1}^{e,n}(\xi_{i-1/2}) + \pi_{i-1}^{e,n}(\xi_{i-1}) - \overrightarrow{w}_{i-1}^{n+1,1}) \\ &\quad + \frac{\Delta ta}{h_i^n \Delta \xi} (\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_{i-1/2})) - \frac{ga}{h_i^n} (h_i^n - L_i^{n+1} h_i^{e,n}(\xi_i)) z'(\xi_i) \\ &= \overrightarrow{w}_i^n - \frac{\Delta ta}{h_i^n \Delta \xi} (\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_i) + \overrightarrow{w}_i^{n+1,1} - \pi_{i-1}^{e,n}(\xi_{i-1/2}) + \pi_{i-1}^{e,n}(\xi_{i-1}) - \overrightarrow{w}_{i-1}^{n+1,1}) \\ &\quad + \frac{\Delta ta}{h_i^n \Delta \xi} (\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_{i-1/2})) \\ &= \overrightarrow{w}_i^n - \frac{\Delta ta}{h_i^n \Delta \xi} (-\pi_i^{e,n}(\xi_i) + \overrightarrow{w}_i^{n+1,1} + \pi_{i-1}(\xi_{i-1}) - \overrightarrow{w}_{i-1}^{n+1,1}), \end{aligned}$$

where we have used that  $\pi_{i-1}^{e,n}(\xi_{i-1/2}) = \pi_i^{e,n}(\xi_{i-1/2})$  and that  $L_i^{n+1} = 1$  since in the first iteration  $u_{i\pm 1/2}^{*,n+1} = 0$ . Note that at this point we are using that the bottom topography is continuous what implies that  $h^e$  is continuous as well as  $\pi^e$ , since we are considering a water at rest situation.

It is clear that  $\overrightarrow{w}_i^{n+1,1} = \pi_i^{e,n}(\xi_i)$ ,  $\overrightarrow{w}_{i-1}^{n+1,1} = \pi_{i-1}^{e,n}(\xi_{i-1})$  is a solution of the previous equation and therefore,  $\overrightarrow{w}_i^{n+1} = \overrightarrow{w}_i^n = \pi_i^{e,n}(\xi_i)$ . The result for  $\overleftarrow{w}$  is analogous.

Once we know  $\overrightarrow{w}^{n+1}$  and  $\overleftarrow{w}^{n+1}$ , we need to compute  $\pi_{i+1/2}^{*,n+1}$  and  $\pi_{i+1/2}^{*,n+1}$  making use of (2.2.38), (2.2.39) and (2.2.23) and we obtain

$$\begin{aligned} \pi_{i+1/2}^{*,n+1} &= \frac{\pi_i^{e,n}(\xi_{i+1/2}) + \pi_{i+1}^{e,n}(\xi_{i+1/2})}{2} = \frac{\pi_i^{e,n}(\xi_{i+1/2}) + \pi_i^{e,n}(\xi_{i+1/2})}{2} = \pi_i^{e,n}(\xi_{i+1/2}), \\ u_{i+1/2}^{*,n+1} &= \frac{\pi_i^{e,n}(\xi_{i+1/2}) - \pi_{i+1}^{e,n}(\xi_{i+1/2})}{2a} = \frac{\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_{i+1/2})}{2a} = 0. \end{aligned}$$

Then, using (2.2.11) and (2.2.37), we obtain  $L_i^{n+1} = 1$  and  $x_{i\pm 1/2}^{*,n+1} = \xi_{i\pm 1/2}$ . Therefore,

$$\begin{aligned} \mathcal{L}_i^{n+1} &= \frac{1}{\Delta \xi} (\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_{i-1/2})), \\ \mathcal{G}_i^{n+1} &= -\frac{1}{\Delta \xi} (\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_{i-1/2})) + g (h_i^n - h_i^{e,n}(\xi_i)) z'(\xi_i) \\ &= -\frac{1}{\Delta \xi} (\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_{i-1/2})). \end{aligned}$$

Finally,  $(L\bar{h})_i^{n+1} = h_i^n$ ,  $(L\bar{h}u)_i^{n+1} = (hu)_i^n$  and, using the fact that  $L_i^{n+1} = 1$ , we get that the scheme is exactly well-balanced.  $\square$

**Implicit-explicit scheme** We can also consider the following first order implicit-explicit scheme:

$$(\overline{Lhu})_i^{n+1} = (hu)_i^n - \Delta t (\mathcal{L}_i^{n+1} + \mathcal{G}_i^n), \quad (2.2.43)$$

where  $\mathcal{L}_i^{n+1}$  is defined as in (2.2.35) and

$$\mathcal{G}_i^n = -\frac{1}{\Delta\xi} (\pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_{i-1/2})) \quad (2.2.44)$$

since the second term of  $\mathcal{G}$  evaluated at time  $t^n$  is zero.

In this case,  $\vec{w}^{n+1}$  and  $\overleftarrow{w}^{n+1}$  are the solutions of the following linear system

$$\begin{cases} \vec{w}_i^{n+1} = \vec{w}_i^n - \Delta t \left( (\mathcal{L}_{\vec{w}})_i^{n+1} + \frac{a}{h_i^n} \mathcal{G}_i^n \right), \\ \overleftarrow{w}_i^{n+1} = \overleftarrow{w}_i^n - \Delta t \left( (\mathcal{L}_{\overleftarrow{w}})_i^{n+1} - \frac{a}{h_i^n} \mathcal{G}_i^n \right) \end{cases} \quad (2.2.45)$$

with  $(\mathcal{L}_{\vec{w}})_i^{n+1}$ ,  $(\mathcal{L}_{\overleftarrow{w}})_i^{n+1}$  defined as in (2.2.42). Remark that now the solution of (2.2.45) is straightforward and the fixed point iterations are avoided, making this implicit-explicit version more efficient.

**Theorem 2.2.2.** *Given any water at rest stationary solution of system (1.1.15) with a continuous bottom topography  $z$ , the scheme defined by (2.2.43) and (2.2.45) preserves it and the scheme is exactly well-balanced.*

*Proof.* The proof is analogous to the previous one. □

### 2.2.2.2 Second order schemes

**Implicit scheme** Let us now define the second order exactly well-balanced implicit scheme. In this case, applying the trapezoidal rule

$$(\overline{Lhu})_i^{n+1} = (hu)_i^n - \frac{\Delta t}{2} (\mathcal{L}_i^n + \mathcal{L}_i^{n+1} + \mathcal{G}_i^n + \mathcal{G}_i^{n+1}), \quad (2.2.46)$$

where

$$\begin{aligned} \mathcal{L}_i^n &= \frac{1}{\Delta\xi} (\pi_{i+1/2}^{*,n} - \pi_{i-1/2}^{*,n}), \\ \mathcal{L}_i^{n+1} &= \frac{1}{\Delta\xi} (\pi_{i+1/2}^{*,n+1} - \pi_{i-1/2}^{*,n+1}), \end{aligned}$$

## 40 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

and  $\mathcal{G}_i^n$  is computed as in (2.2.44) and  $\mathcal{G}_i^{n+1}$  as in (2.2.36) where now we use the following second order approximations

$$\begin{aligned} x_{i\pm 1/2}^{*,n+1} &= x_{i\pm 1/2} + \frac{\Delta t}{2} \left( u_{i\pm 1/2}^{*,n+1} + u_{i\pm 1/2}^{*,n} \right), \\ x_i^{*,n+1} &= x_i + \frac{\Delta t}{4} \left( u_{i+1/2}^{*,n+1} + u_{i-1/2}^{*,n+1} + u_{i+1/2}^{*,n} + u_{i-1/2}^{*,n} \right), \\ L_i^{n+1} &= 1 + \frac{\Delta t}{2\Delta\xi} \left( u_{i+1/2}^{*,n+1} - u_{i-1/2}^{*,n+1} + u_{i+1/2}^{*,n} - u_{i-1/2}^{*,n} \right). \end{aligned}$$

The values  $\pi_{i+1/2}^{*,n}$ ,  $\pi_{i+1/2}^{*,n+1}$ ,  $u_{i+1/2}^{*,n}$ ,  $u_{i+1/2}^{*,n+1}$  can be obtained as follows

$$\pi_{i+1/2}^{*,n} = \frac{P_{i,\vec{w}}^{o2}(\xi_{i+1/2}, t^n) + P_{i+1,\overleftarrow{w}}^{o2}(\xi_{i+1/2}, t^n)}{2} = \frac{\vec{w}_{i+1/2-}^n + \overleftarrow{w}_{i+1/2+}^n}{2} \quad (2.2.47)$$

$$\pi_{i+1/2}^{*,n+1} = \frac{P_{i,\vec{w}}^{o2}(\xi_{i+1/2}, t^{n+1}) + P_{i+1,\overleftarrow{w}}^{o2}(\xi_{i+1/2}, t^{n+1})}{2} = \frac{\vec{w}_{i+1/2-}^{n+1} + \overleftarrow{w}_{i+1/2+}^{n+1}}{2} \quad (2.2.48)$$

$$u_{i+1/2}^{*,n} = \frac{P_{i,\vec{w}}^{o2}(\xi_{i+1/2}, t^n) - P_{i+1,\overleftarrow{w}}^{o2}(\xi_{i+1/2}, t^n)}{2a} = \frac{\vec{w}_{i+1/2-}^n - \overleftarrow{w}_{i+1/2+}^n}{2a} \quad (2.2.49)$$

$$u_{i+1/2}^{*,n+1} = \frac{P_{i,\vec{w}}^{o2}(\xi_{i+1/2}, t^{n+1}) - P_{i+1,\overleftarrow{w}}^{o2}(\xi_{i+1/2}, t^{n+1})}{2a} = \frac{\vec{w}_{i+1/2-}^{n+1} - \overleftarrow{w}_{i+1/2+}^{n+1}}{2a} \quad (2.2.50)$$

where in order to compute  $\vec{w}_{i+1/2-}^n$ ,  $\vec{w}_{i+1/2-}^{n+1}$ ,  $\overleftarrow{w}_{i+1/2+}^n$  and  $\overleftarrow{w}_{i+1/2+}^{n+1}$  we use the reconstruction proposed in (2.2.27).

Similarly,  $\vec{w}^{n+1}$  and  $\overleftarrow{w}^{n+1}$  are the solutions of the following non-linear system

$$\begin{cases} \vec{w}_i^{n+1} = \vec{w}_i^n - \frac{\Delta t}{2} \left( (\mathcal{L}_{\vec{w}}^n)_i + (\mathcal{L}_{\vec{w}}^{n+1})_i + \frac{a}{h_i^n} (\mathcal{G}_i^n + \mathcal{G}_i^{n+1}) \right), \\ \overleftarrow{w}_i^{n+1} = \overleftarrow{w}_i^n - \frac{\Delta t}{2} \left( (\mathcal{L}_{\overleftarrow{w}}^n)_i + (\mathcal{L}_{\overleftarrow{w}}^{n+1})_i - \frac{a}{h_i^n} (\mathcal{G}_i^n + \mathcal{G}_i^{n+1}) \right) \end{cases} \quad (2.2.51)$$

where in this case

$$\begin{aligned} (\mathcal{L}_{\vec{w}}^n)_i &= \frac{a}{2h_i^n \Delta\xi} (\vec{w}_{i+1/2-}^n - \vec{w}_{i-1/2-}^n), \\ (\mathcal{L}_{\overleftarrow{w}}^n)_i &= -\frac{a}{2h_i^n \Delta\xi} (\overleftarrow{w}_{i+1/2+}^n - \overleftarrow{w}_{i-1/2+}^n), \\ (\mathcal{L}_{\vec{w}}^{n+1})_i &= \frac{a}{2h_i^n \Delta\xi} (\vec{w}_{i+1/2-}^{n+1} - \vec{w}_{i-1/2-}^{n+1}), \\ (\mathcal{L}_{\overleftarrow{w}}^{n+1})_i &= -\frac{a}{2h_i^n \Delta\xi} (\overleftarrow{w}_{i+1/2+}^{n+1} - \overleftarrow{w}_{i-1/2+}^{n+1}) \end{aligned} \quad (2.2.52)$$

and we use the second order reconstruction introduced in (2.2.27).

The equations in (2.2.51) define two systems with four diagonals each: the ones corresponding to  $\vec{w}_{i-2}$ ,  $\vec{w}_{i-1}$ ,  $\vec{w}_i$ ,  $\vec{w}_{i+1}$  in the case of the system for  $\vec{w}^{n+1}$  and the ones corresponding to  $\overleftarrow{w}_{i-1}$ ,  $\overleftarrow{w}_i$ ,  $\overleftarrow{w}_{i+1}$ ,  $\overleftarrow{w}_{i+2}$  in the case of the system for  $\overleftarrow{w}^{n+1}$ .

**Theorem 2.2.3.** *Given any water at rest stationary solution of system (1.1.15) with a continuous bottom topography  $z$ , the scheme defined by (2.2.46) and (2.2.51) preserves it and the scheme is exactly well-balanced.*

*Proof.* The proof is analogous to the previous one.  $\square$

**Implicit-explicit scheme** Finally, the implicit-explicit second order scheme is obtained by treating the function  $\mathcal{L}$  implicitly and the function  $\mathcal{G}$  explicitly, considering the SSP2(2,2,2) IMEX scheme (see equation (1.1.76)).

Similarly to the implicit second order case, the systems that we have to solved for this scheme have four diagonals each.

Note that the use of the IMEX scheme avoids the costly nonlinear system inversion.

A similar theorem as Theorem 2.2.3 is still valid in this case.

## 2.2.3 Explicit exactly well-balanced Lagrangian schemes

As previously said, even though our aim is the design of schemes that perform the Lagrangian step implicitly, we will describe, for the sake of completeness, first and second order explicit Lagrangian schemes.

### 2.2.3.1 First order explicit Lagrangian scheme

The first order explicit scheme can be written similarly to how it has been done for the implicit schemes, with the difference that in this case functions  $\mathcal{L}$  and  $\mathcal{G}$  are evaluated at time  $t^n$ .

$$(\overline{Lhu})_i^{n+1} = (hu)_i^n - \Delta t (\mathcal{L}_i^n + \mathcal{G}_i^n),$$

where

$$\mathcal{L}_i^n = \frac{1}{\Delta \xi} \left( \pi_{i+1/2}^{*,n} - \pi_{i-1/2}^{*,n} \right),$$

and

$$\mathcal{G}_i^n = -\frac{1}{\Delta \xi} (\pi_i^{e,n}(x_{i+1/2}) - \pi_i^{e,n}(x_{i-1/2})).$$

Moreover,

$$L_i^n = 1 + \frac{\Delta t}{\Delta \xi} \left( u_{i+1/2}^{*,n} - u_{i-1/2}^{*,n} \right).$$

The needed approximations of  $\pi_{i+1/2}^{*,n}$  and  $u_{i+1/2}^{*,n}$  are obtained as

$$\pi_{i+1/2}^{*,n} = \frac{P_{i,\overline{w}}^{o1}(\xi_{i+1/2}, t^n) + P_{i+1,\overline{w}}^{o1}(\xi_{i+1/2}, t^n)}{2} = \frac{\overline{w}_{i+1/2-}^n + \overleftarrow{w}_{i+1/2+}^n}{2}$$

and

$$u_{i+1/2}^{*,n} = \frac{P_{i,\overline{w}}^{o1}(\xi_{i+1/2}, t^n) - P_{i+1,\overline{w}}^{o1}(\xi_{i+1/2}, t^n)}{2a} = \frac{\overline{w}_{i+1/2-}^n - \overleftarrow{w}_{i+1/2+}^n}{2a}$$

where

$$\begin{aligned}\vec{w}_{i+1/2-}^n &= P_{i,\vec{w}}^{o1}(\xi_{i+1/2}, t^n) = \pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_i) + \vec{w}_i^n, \\ \overleftarrow{w}_{i+1/2+}^n &= P_{i+1,\overleftarrow{w}}^{o1}(\xi_{i+1/2}, t^n) = \pi_{i+1}^{e,n}(\xi_{i+1/2}) - \pi_{i+1}^{e,n}(\xi_i) + \overleftarrow{w}_{i+1}^n.\end{aligned}$$

### 2.2.3.2 Second order explicit Lagrangian scheme

For the second order explicit case we propose using a mid-point quadrature rule in time in order to obtain

$$(L\bar{h}u)_i^{n+1} = (hu)_i^n - \Delta t \left( \mathcal{L}_i^{n+1/2} + \mathcal{G}_i^{n+1/2} \right),$$

where

$$\begin{aligned}\mathcal{L}_i^{n+1/2} &= \frac{1}{\Delta\xi} \left( \pi_{i+1/2}^{*,n+1/2} - \pi_{i-1/2}^{*,n+1/2} \right), \\ \mathcal{G}_i^{n+1/2} &= -\frac{1}{\Delta\xi} \left( \pi_i^e(x_{i+1/2}^{*,n+1/2}) + \pi_i^e(x_{i-1/2}^{*,n+1/2}) \right) \\ &\quad + g \left( h_i^n - L_i^{n+1/2} h^e(x_i^{*,n+1/2}) \right) z'(x_i^{*,n+1/2})\end{aligned}$$

and

$$L_i^{n+1/2} = 1 + \frac{\Delta t}{2\Delta\xi} (u_{i+1/2}^{*,n} - u_{i-1/2}^{*,n}).$$

The values  $\pi_{i+1/2}^{*,n}$  and  $u_{i+1/2}^{*,n}$  are approximated following the idea used for the implicit schemes:

$$\begin{aligned}\pi_{i+1/2}^{*,n} &= \frac{P_{i,\vec{w}}^{o2}(\xi_{i+1/2}, t^n) + P_{i+1,\overleftarrow{w}}^{o2}(\xi_{i+1/2}, t^n)}{2} = \frac{\vec{w}_{i+1/2-}^n + \overleftarrow{w}_{i+1/2+}^n}{2}, \\ u_{i+1/2}^{*,n} &= \frac{P_{i,\vec{w}}^{o2}(\xi_{i+1/2}, t^n) - P_{i+1,\overleftarrow{w}}^{o2}(\xi_{i+1/2}, t^n)}{2a} = \frac{\vec{w}_{i+1/2-}^n - \overleftarrow{w}_{i+1/2+}^n}{2a}\end{aligned}$$

where for the reconstructions  $\vec{w}_{i+1/2-}^n, \overleftarrow{w}_{i+1/2+}^n$  we use (2.2.27).

Once we have computed the previous values, we can obtain the approximations of the position at time  $t^{n+1/2}$  of characteristic curves as

$$\begin{aligned}x_i^{*,n+1/2} &= \xi_i + \frac{\Delta t}{4} \left( u_{i-1/2}^{*,n} + u_{i+1/2}^{*,n} \right), \\ x_{i\pm 1/2}^{*,n+1/2} &= \xi_{i\pm 1/2} + \frac{\Delta t}{2} u_{i\pm 1/2}^{*,n}.\end{aligned}$$

## 2.3 The projection step

As part of the Lagrange-Projection algorithm, after solving the Lagrangian step, we shall project the result back onto Eulerian coordinates. In other words, we must project the

piece-wise constant approximations of  $L\bar{U}(\xi, t)$  obtained after the Lagrangian step, onto the Eulerian cells  $(x_{i-1/2}, x_{i+1/2})$ . This will always be done explicitly.

For doing so, we need to compute

$$U_i(t) = \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U(x, t) dx.$$

For any time  $T \geq 0$ , let us define  $\hat{\xi}_{i+1/2}(T)$  as the origin of the curve  $x(\hat{\xi}_{i+1/2}, t)$  which at time  $t = T$  passes through  $x_{i+1/2}$  (see Figure 2.1), that is, given  $t \geq 0$  we define  $\hat{\xi}_{i+1/2}(t)$  such that

$$x(\hat{\xi}_{i+1/2}(t), t) = x_{i+1/2}.$$

Using this notation, we may write

$$U_i(t) = \frac{1}{\Delta x} \int_{x(\hat{\xi}_{i-1/2}(t), t)}^{x(\hat{\xi}_{i+1/2}(t), t)} U(x, t) dx = \frac{1}{\Delta x} \int_{\hat{\xi}_{i-1/2}(t)}^{\hat{\xi}_{i+1/2}(t)} L(\xi, t) \bar{U}(\xi, t) d\xi,$$

and we can split the integral as follows

$$\begin{aligned} U_i(t) &= \frac{1}{\Delta x} \int_{\hat{\xi}_{i-1/2}(t)}^{\xi_{i-1/2}(t)} L(\xi, t) \bar{U}(\xi, t) d\xi + \frac{1}{\Delta x} \int_{\xi_{i-1/2}(t)}^{\xi_{i+1/2}(t)} L(\xi, t) \bar{U}(\xi, t) d\xi \\ &\quad + \frac{1}{\Delta x} \int_{\xi_{i+1/2}(t)}^{\hat{\xi}_{i+1/2}(t)} L(\xi, t) \bar{U}(\xi, t) d\xi. \end{aligned}$$

Remark that the central integral corresponds to  $(L\bar{U})_i(t)$ , which is already known from the previous Lagrangian step. Therefore,

$$U_i^{n+1} = (L\bar{U})_i^{n+1} + \frac{1}{\Delta x} \int_{\hat{\xi}_{i-1/2}}^{\xi_{i-1/2}} L(\xi, t^{n+1}) \bar{U}(\xi, t^{n+1}) d\xi + \frac{1}{\Delta x} \int_{\xi_{i+1/2}}^{\hat{\xi}_{i+1/2}} L(\xi, t^{n+1}) \bar{U}(\xi, t^{n+1}) d\xi. \quad (2.3.1)$$

The evaluation of the integrals in previous expressions is now required. To do so, first and second order approaches will be presented.

### 2.3.1 First order projection scheme

The previous integrals can be approximated in the following way:

$$\frac{1}{\Delta x} \int_{\hat{\xi}_{i+1/2}}^{\xi_{i+1/2}} L(\xi, t^{n+1}) \bar{U}(\xi, t^{n+1}) d\xi = \frac{\hat{\xi}_{i+1/2} - \xi_{i+1/2}}{\Delta x} (L\bar{U})_{i+1/2}^{n+1},$$

where

$$(L\bar{U})_{i+1/2}^{n+1} = \begin{cases} (L\bar{U})_i^{n+1} & \text{for } \xi_{i+1/2} > \hat{\xi}_{i+1/2}, \\ (L\bar{U})_{i+1}^{n+1} & \text{for } \xi_{i+1/2} \leq \hat{\xi}_{i+1/2}. \end{cases}$$

Moreover, using the approximation

$$\hat{\xi}_{i+1/2} = x_{i+1/2} - \Delta t u_{i+1/2}^*,$$

from (2.3.1) we get

$$U_i^{n+1} = (L\bar{U})_i^{n+1} - \frac{\Delta t}{\Delta x} \left( u_{i+1/2}^* (L\bar{U})_{i+1/2}^{n+1} - u_{i-1/2}^* (L\bar{U})_{i-1/2}^{n+1} \right).$$

### 2.3.2 Second order projection scheme

For the second order case, the integrals (2.3.1) need to be computed with second order accuracy. To do so, let us consider a piecewise linear reconstruction of averages of  $(L\bar{U})_i^{n+1}$  and, taking into account that the velocities  $u_{i+1/2}^{*,n+1}$  are continuously defined at the intercells, we use a continuous piecewise linear interpolation at the intercells.

It can be seen that we can write (2.3.1) again as

$$U_i^{n+1} = (L\bar{U})_i^{n+1} - \frac{\Delta t}{\Delta x} \left( u_{i+1/2}^* (L\bar{U})_{i+1/2}^{n+1} - u_{i-1/2}^* (L\bar{U})_{i-1/2}^{n+1} \right),$$

where

$$(L\bar{U})_{i+1/2}^{n+1} = \begin{cases} (L\bar{U})_{i+1/2-}^{n+1} & \text{for } u_{i+1/2}^* > 0, \\ (L\bar{U})_{i+1/2+}^{n+1} & \text{for } u_{i+1/2}^* \leq 0, \end{cases}$$

and

$$\begin{aligned} (L\bar{U})_{i+1/2-}^{n+1} &= \frac{1}{L_i^{n+1}} \left( (L\bar{U})_i^{n+1} + \frac{1}{2} (\delta L\bar{U})_i^{n+1} \left( \Delta x - \frac{\Delta t}{L_i^{n+1}} u_{i+1/2}^{*,n+1} \right) \right), \\ (L\bar{U})_{i+1/2+}^{n+1} &= \frac{1}{L_{i+1}^{n+1}} \left( (L\bar{U})_{i+1}^{n+1} + \frac{1}{2} (\delta L\bar{U})_{i+1}^{n+1} \left( -\Delta x - \frac{\Delta t}{L_{i+1}^{n+1}} u_{i+1/2}^{*,n+1} \right) \right). \end{aligned}$$

In the previous expressions,  $(\delta L\bar{U})_i^{n+1}$  and  $(\delta L\bar{U})_{i+1}^{n+1}$  are approximations of the derivatives of  $L\bar{U}$  at time  $t^{n+1}$  at  $x_i$  and  $x_{i+1}$ , respectively, that are computed by means of an avg limiter (see (1.1.43)).  $(L\bar{U})_{i-1/2}^{n+1}$  is defined in a similar way.

The CFL condition associated with the transport step reads

$$\Delta t \max_j \{ (u_{j-1/2}^*)^+ - (u_{j+1/2}^*)^- \} \leq \Delta x. \quad (2.3.2)$$

Let us remark that this stability restriction always remains as this step, contrary to the Lagrangian one, is performed explicitly.

Observe that this projection step does not destroy the exactly well-balanced character of the scheme as we focus only on steady states corresponding to water at rest, where  $u = 0$ , and we have shown that  $L_i = 1$ , what makes the projection step trivial in that particular case. A more interesting situation occurs when moving equilibria are considered. In that case, the procedure described in Section 3 could be extended to moving equilibria following [25], but the projection step must be modified in order to preserve those steady states. The difficulty in this case is due to the fact that when  $u \neq 0$ , smooth stationary solutions depend on time in the Lagrangian framework. The idea presented in [72] for the design of fully well-balanced schemes for the transport step consists on properly defining the fluctuations and equilibria of  $LU_i(t)$ , denoting by  $LU_i^e(t)$  the solution of the Lagrangian system (2.2.2) applied to the stationary solution  $U_i^e(x)$  and the fluctuation  $L\bar{U}_i^f(t)$ , which is given by

$$L\bar{U}_i^f(t) = L\bar{U}_i(t) - L\bar{U}_i^e(t).$$

## 2.4 Numerical results

We now intend to test and compare the different numerical schemes introduced in this chapter. In what follows, we shall use the following notation for the used schemes. First and second order explicit schemes, which are described in 2.2.3, are denoted by EXP O1 and EXP O2 respectively. The first order implicit scheme given by (2.2.34)-(2.2.41) is denoted by IMP-NL O1 while its second order extension (2.2.46)-(2.2.51) is denoted by IMP-NL O2. Finally the first order (2.2.43)-(2.2.45) implicit-explicit scheme and its second order extension by means of (1.1.76) are represented by IMP-IMEX O1 and IMP-IMEX O2 respectively. In order to make a better assessment of the performance of our second order schemes, in some of the tests we will also include the results obtained with a DIRK (Diagonal Implicit Runge-Kutta) scheme [76] that we will represent by IMP-DIRK O2.

In what follows, CFL condition refers to the restriction needed to satisfy for the stability of the explicit schemes, that is,  $\Delta t$  is the minimum that grants the stability conditions given by (2.2.14) and (2.3.2). Note that both the fully implicit and the IMEX schemes do not need the restriction (2.2.14) and they are only limited by (2.3.2). This means that, in situations where velocity is small compared to the sound speed, the stability condition for the Lagrangian step (2.2.14) is very limiting when compared to the projection step stability condition (2.3.2), which justifies the use of the implicit approach. Therefore, we will see in the following test cases that CFL greater than 1 (when compared to the stability criterion for the explicit scheme) may be chosen for the implicit schemes.

### 2.4.1 Exactly well-balanced property test

This first test aims at assessing the well balanced property of the schemes. Let us propose the water at rest solution given by:

$$u = 0, \quad h + z = C, \quad \text{with } C \in \mathbb{R}$$

where  $C = 0$  and the bottom topography is a gaussian bump given by

$$z(x) = -1 + \frac{1}{2}e^{-x^2}, \quad x \in [-5, 5]. \quad (2.4.1)$$

The interval  $[-5, 5]$  is discretized using 200 cells and final time is set to  $t = 5$ . In Table 2.1 we show the  $L^1$  errors obtained at final time using CFL 0.5 for the explicit schemes and CFL 2 for the implicit ones. As expected, all the exactly well-balanced schemes are able to preserve the water at rest steady state.

Order 1			Order 2		
EXP	IMP-IMEX	IMP-NL	EXP	IMP-IMEX	IMP-NL
1.95e-14	8.57e-14	5.36e-14	1.67e-13	5.56e-14	6.66e-14

Table 2.1:  $L^1$  errors between the numerical solution at initial and final time  $t = 5$  for water at rest initial condition

### 2.4.2 Computational time vs error

The objective of this test is to show the better performance of IMEX schemes against the fully implicit approach. To do so, we study the computational efficiency of these schemes.

Inspired by [4], we have considered a test case consisting in a channel of length  $L = 14000 m$  with a bottom topography defined by

$$z(x) = - \left( 50.5 - 40 \frac{L-x}{L} + 10 \sin \left( \pi \left( 4 \frac{L-x}{L} - \frac{1}{2} \right) \right) \right).$$

We then simulate a tidal wave of  $0.5 m$  amplitude by imposing the following initial and boundary conditions:

$$\begin{aligned} h(x, 0) &= -z(x) + 1, \\ q(x, 0) &= 0, \\ h(L, t) &= -z(L) + \frac{1}{2} + \frac{1}{2} \sin \left( \pi \left( \frac{4t}{86400} + \frac{1}{2} \right) \right), \\ q(0, t) &= 0. \end{aligned}$$

Let us recall that the computational advantage of the numerical approaches presented in this work make sense especially in the case of low Froude number. Indeed, in such situations the restriction of the usual CFL condition is mainly driven by the acoustic waves, so that implicit schemes allow to avoid it. This is indeed the case of this tidal wave test.

In Figures 2.2 and 2.3 we show the computational time needed to perform the simulation and the error for an increasing number of cells for the different type of schemes. The error has been obtained by computing a reference solution for a very fine mesh.

In all the numerical tests, a CFL value of 0.5 has been used in the case of explicit schemes. As far as the implicit schemes are concerned, we consider two different approaches: first, the simulations are performed using the same CFL restriction 0.5 as in the explicit case (left-hand side pictures); second a CFL corresponding to 100 with respect to the usual restriction is used (right-hand side pictures).

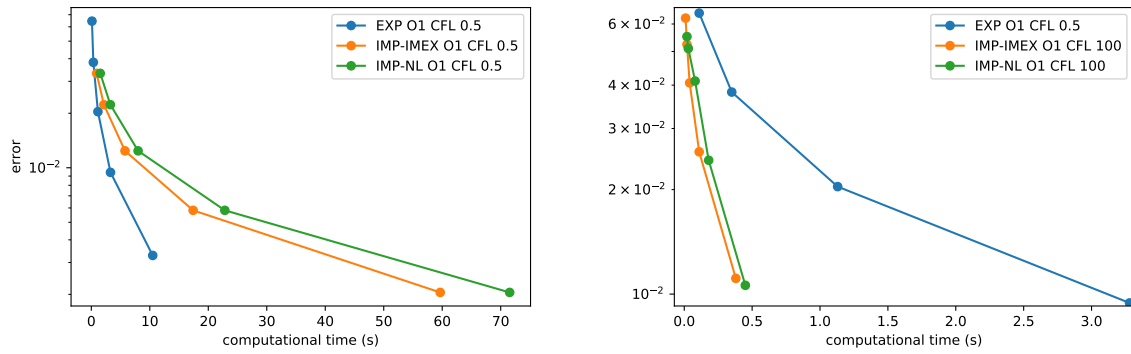
We can observe that for the low CFL value case, the explicit scheme is more efficient than the implicit ones, as expected. It is in the case when large CFL numbers are used that the implicit schemes outperform the explicit one. Indeed we see an important reduction of computational time needed for the implicit schemes when compared to the explicit one in order to obtain a similar error. For example, in the first order case, we see a 95% reduction in computational time for the IMP-IMEX scheme compared to the EXP one for comparable errors. Similarly, in the second order case, the reduction found is approximately 60%. This improvement in the efficiency is expected to be even bigger for 2D problems.

Moreover, when comparing IMP-IMEX and IMP-NL, the former seems to be more efficient than the latter. This is due to the iterations of the fixed point scheme that are needed. In Tables 2.2 and 2.3 we show the maximum number of iterations needed to solve the nonlinear systems during the computation and the global number of fixed point iterations required for the whole simulation. This is done for CFL equal to 0.5 and 100, respectively, for the IMP-NL O2 and the IMP-DIRK O2 schemes. We can also remark that the computational time required by the DIRK scheme is larger than in the IMP-NL O2 case when the CFL is set to 0.5 but it is a bit shorter when the CFL is 100. This can also be explained by checking the total number of iterations in Tables 2.2 and 2.3. However, the IMP-IMEX O2 scheme is faster than the IMP-NL and the IMP-DIRK schemes, since only two linear schemes have to be solved.

Lastly, we could think that even though the implicit schemes are more efficient, the errors should be greater than for the explicit schemes. However, with this initial condition we are considering long waves, so the errors are not too big. To show this, in Figure 2.4 we have plotted the free surface for the different first and second order schemes using both CFL 0.5 and 100 for the implicit ones. In the first order case, we observe that the implicit schemes are more diffusive than the explicit one, but this is to be expected and the size of the errors is not too large. When we analyse the results obtained for the second order schemes we see that both the explicit and the implicit schemes with CFL 0.5 give almost

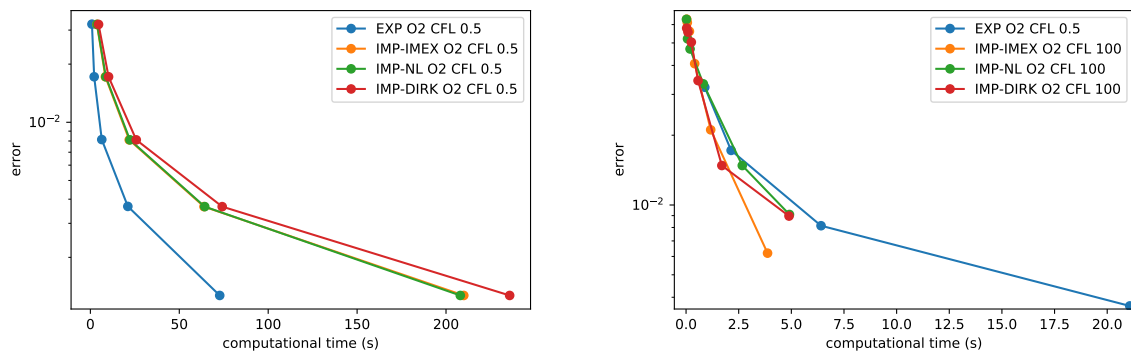
## 48 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

identical results. Of course, when for the implicit schemes the CFL is increased to 100, they are again more diffusive but not too big differences are seen.



(a) Using a CFL value of 0.5 for the explicit and implicit schemes (b) Using a CFL value of 100 for the implicit schemes and 0.5 for the explicit scheme

Figure 2.2: Computational time vs. error for an increasing number of cells using first order schemes



(a) Using a CFL value of 0.5 for the explicit and implicit schemes (b) Using a CFL value of 100 for the implicit schemes and 0.5 for the explicit scheme

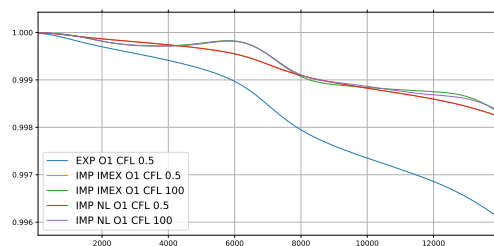
Figure 2.3: Computational time vs. error for an increasing number of cells using second order schemes

CFL 0.5	IMP-NL O2		IMP-DIRK O2	
No. of cells	Max Iter	Total Iter	Max Iter	Total Iter
25	3	496	3	717
50	3	981	2	1342
100	3	1897	2	1682
200	3	3681	2	5366
400	3	6942	2	10736

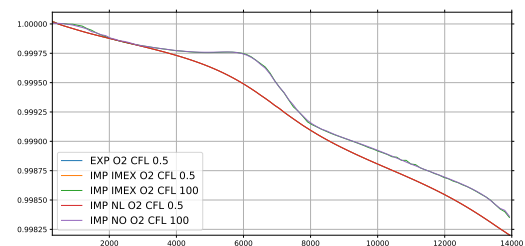
Table 2.2: Maximum number of fixed point iterations to solve the nonlinear systems and total number of iterations performed for different number of cells for the second order nonlinear scheme and the DIRK scheme using a CFL equal to 0.5

CFL 100	IMP-NL O2		IMP-DIRK O2	
No. of cells	Max Iter	Total Iter	Max Iter	Total Iter
25	2	2	1	1
50	7	9	6	7
100	7	21	7	18
200	7	39	6	34
400	7	82	6	73

Table 2.3: Maximum number of fixed point iterations to solve the nonlinear systems and total number of iterations performed for different number of cells for the second order nonlinear scheme and the DIRK scheme using a CFL equal to 100



(a) First order schemes



(b) Second order schemes

Figure 2.4: Free surface for the different first and second order schemes

### 2.4.3 Order test

Let us now check the order of the schemes. In order to do so, we consider as initial condition a small perturbation flowing over a gaussian bump in a domain with length

## 50 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

$L = 14000 m$ . More explicitly, the bottom topography is given by

$$z(x) = - \left( 50 - \exp \left( - \frac{(x - 7000)^2}{1000000} \right) \right),$$

and the initial condition writes as  $q(x, 0) = 0$  and

$$h(x, 0) = -z(x) + \begin{cases} 0.05 \left( 1 + \cos \left( \frac{2\pi(x-4750)}{3500} \right) \right) & \text{if } 3000 < x < 6500 \\ 0.05 \left( - \left( 1 + \cos \left( \frac{2\pi(x-9250)}{3500} \right) \right) \right) & \text{if } 7500 < x < 11000 \\ 0 & \text{otherwise.} \end{cases}$$

Free surface corresponding to this initial condition is shown in Figure 2.5. The errors in  $L^1$  obtained for the different schemes are then shown in Tables 2.4, 2.5 and 2.6, where we see that the expected accuracy is obtained.

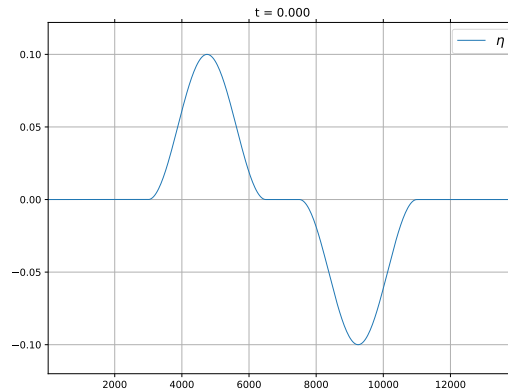


Figure 2.5: Free surface corresponding to the initial condition for the order test case

No. of cells	EXP - Order 1				EXP - Order 2			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
25	3.49e-1	0.00	4.50e0	0.00	2.17e-1	0.00	1.58e0	0.00
50	1.88e-1	0.89	2.69e0	0.74	6.03e-2	1.85	6.89e-1	1.20
100	9.39e-2	1.00	1.35e0	0.99	1.68e-2	1.84	2.01e-1	1.77
200	4.21e-2	1.16	6.07e-1	1.16	4.17e-3	2.01	5.42e-2	1.89
400	1.44e-2	1.54	2.07e-2	1.54	8.50e-4	2.30	1.17e-2	2.21

Table 2.4: Dimensionless errors in  $L^1$  norm and convergence rates for the explicit LP schemes with CFL value 0.5

No. of cells	IMP-IMEX - Order 1				IMP-IMEX - Order 2			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
25	7.03e-1	0.00	1.37e1	0.00	6.53e-1	0.00	3.60e0	0.00
50	5.25e-1	0.42	9.39e0	0.55	1.66e-1	1.97	3.55e0	0.02
100	3.78e-1	0.47	5.85e0	0.68	4.78e-2	1.80	1.16e0	1.61
200	2.18e-1	0.79	3.01e0	0.96	1.30e-2	1.87	3.21e-1	1.86
400	8.89e-2	1.30	1.16e0	1.37	3.22e-3	2.02	7.21e-2	2.16

Table 2.5: Dimensionless errors in  $L^1$  norm and convergence rates for the implicit IMEX LP schemes with CFL value 3

No. of cells	IMP-NL - Order 1				IMP-NL - Order 2			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
25	7.03e-1	0.00	1.37e1	0.00	4.57e-1	0.00	5.61e0	0.00
50	5.25e-1	0.42	9.39e0	0.55	1.60e-1	1.51	3.89e0	0.52
100	3.78e-1	0.47	5.85e0	0.68	5.46e-2	1.55	1.19e0	1.30
200	2.18e-1	0.79	3.01e0	0.96	1.68e-2	1.70	4.53e-1	1.80
400	8.89e-2	1.30	1.16e0	1.37	4.14e-3	2.02	1.06e-1	2.09

Table 2.6: Dimensionless errors in  $L^1$  norm and convergence rates for the implicit nonlinear LP schemes with CFL value 3

The previous initial condition corresponds to a Froude number of order  $10^{-9}$ . By changing the factor 0.05 in  $h$  to bigger values we obtain similar initial conditions with different Froude numbers. We have then tested the errors and convergence rates of the second order implicit schemes when other small Froude numbers are considered. These errors and convergence rates are shown in Tables 2.7 and 2.8, showing second order convergence as expected.

## 52 Implicit and implicit-explicit LP schemes exactly well-balanced for SWE

No. of cells	IMP-IMEX - Order 2				IMP-NL - Order 2			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
25	8.00e-3	0.00	1.37e-2	0.00	7.77e-03	0.00	1.26e-02	0.00
50	1.96e-3	2.02	4.95e-3	1.47	1.97e-03	1.98	8.72e-02	0.53
100	5.12e-4	1.94	1.65e-3	1.64	5.19e-04	1.93	3.18e-03	1.45
200	1.30e-4	1.98	3.81e-4	2.06	1.30e-04	1.99	7.76e-04	2.04
400	3.23e-5	2.00	9.85e-5	1.96	3.19e-05	2.03	1.78e-04	2.12

Table 2.7: Dimensionless errors in  $L^1$  norm and convergence rates for the implicit second order IMEX and nonlinear LP schemes with CFL value 3 for the initial condition with Froude number of order  $10^{-7}$

No. of cells	IMP-IMEX - Order 2				IMP-NL - Order 2			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
25	6.92e-02	0.00	1.36e-01	0.00	6.80e-02	0.00	1.21e-01	0.00
50	1.54e-02	2.17	4.94e-02	1.46	1.55e-02	2.13	8.49e-02	0.52
100	4.13e-03	1.90	1.58e-02	1.64	4.15e-03	1.90	3.07e-02	1.46
200	1.05e-03	1.98	3.60e-03	2.13	1.06e-03	1.97	7.36e-03	2.06
400	2.61e-04	2.00	8.04e-04	2.16	2.69e-04	1.98	1.36e-03	2.44

Table 2.8: Dimensionless errors in  $L^1$  norm and convergence rates for the implicit second order IMEX and nonlinear LP schemes with CFL value 3 for the initial condition with Froude number of order  $10^{-5}$

### 2.4.4 Perturbation of water at rest

Let us consider  $z(x)$  given by (2.4.1) and the following initial condition, which is a perturbation of the water at rest:

$$h(x) = -z(x) + 0.1e^{-x^2}, \quad u(x) = 0.$$

This initial condition is shown in Figure 2.6.

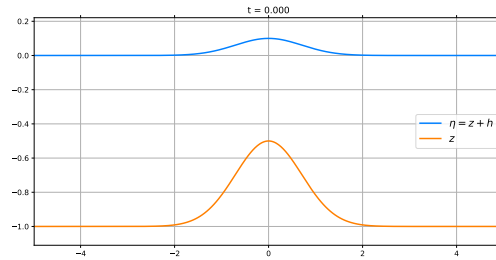


Figure 2.6: Perturbation of water at rest initial condition

In Figures 2.7, 2.8, 2.9 and 2.10 we can check the solutions obtained with the different schemes at time  $t = 0.5$  and  $t = 1$  with 200 cells in the interval  $[-5, 5]$  for the free surface  $\eta = h + z$  and the discharge  $q$ . We include a reference solution that has been computed with the first order explicit scheme using 1600 cells.

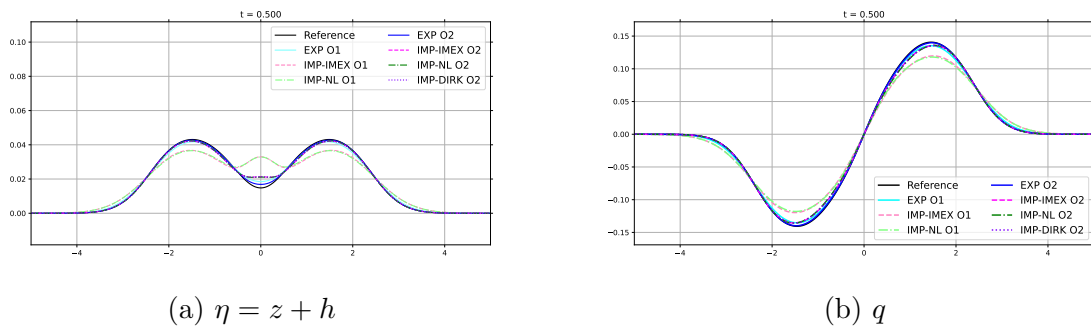


Figure 2.7: Solution for  $\eta$  and  $q$  at  $t = 0.5$  with 200 cells. Explicit: CFL=0.5, Implicit: CFL=2

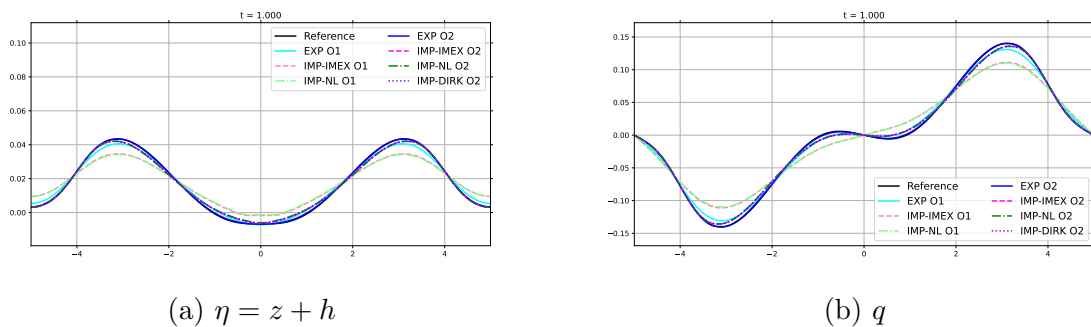


Figure 2.8: Solution for  $\eta$  and  $q$  at  $t = 1$  with 200 cells. Explicit: CFL=0.5, Implicit: CFL=2

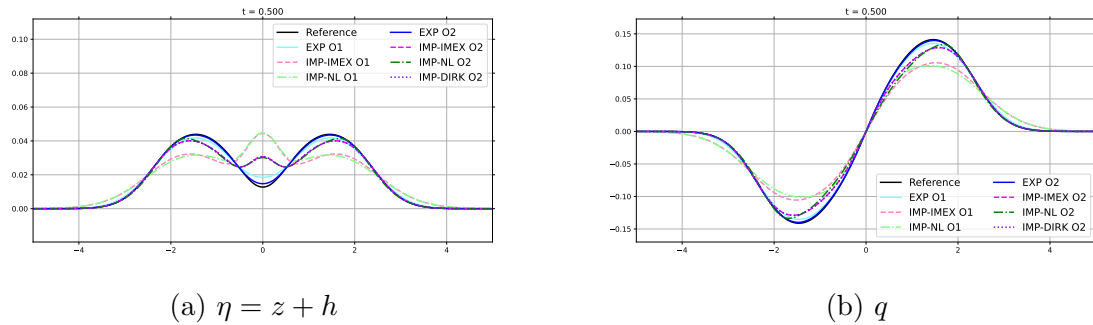


Figure 2.9: Solution for  $\eta$  and  $q$  at  $t = 0.5$  with 200 cells. Explicit: CFL=0.5, Implicit: CFL=5

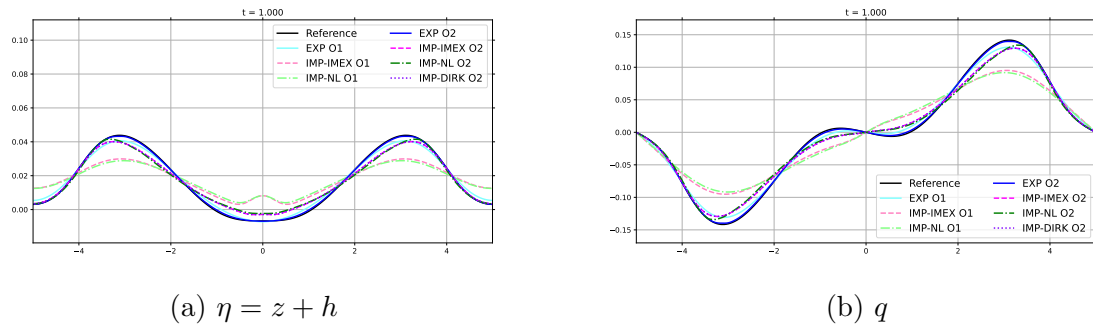


Figure 2.10: Solution for  $\eta$  and  $q$  at  $t = 1$  with 200 cells. Explicit: CFL=0.5, Implicit: CFL=5

In general we see a good agreement of the numerical solutions when compared to the reference solution. The first order implicit schemes are the most diffusive, which is expected. Also, the second order implicit schemes present more diffusion when the CFL is increased. Nevertheless, the greater CFL allowed by implicit schemes make them more efficient.

### 2.4.5 Perturbed water at rest with shock waves

We now intend to study the behavior of the schemes, specially for second order, in the presence of shocks. In such situation, the limiter will play an important role. We consider  $z(x)$  given by (2.4.1) in  $[-5, 5]$  and we set the initial condition:

$$h(x) = \begin{cases} -z(x) & \text{if } |x| \geq 1 \\ -z(x) + 0.1 & \text{if } |x| < 1 \end{cases}, \quad u(x) = 0.$$



This initial condition, shown in Figure 2.11, presents two discontinuities at the interface that generate two shock waves travelling in different directions. We consider a mesh consisting on 200 cells and we perform the simulation up to time  $t = 1$ . For the explicit schemes we have set the CFL value to be 0.5, and for the implicit ones it is set to be 2.

Figures 2.12, 2.13, 2.14 and 2.15 show the numerical solution obtained at times  $t = 0.1$  and  $t = 1$  setting the CFL value to be 2 and 5 for the implicit schemes. As in the previous test, the reference solution has been computed by using the first order explicit scheme with 1600 cells. We see that the schemes are able to correctly handle the initial condition, although very small spurious oscillations are seen for the second order schemes at the advancing front of the shock waves. These spurious oscillations are more pronounced for the second order IMP-NL than for the IMP-IMEX. Nevertheless, in the CFL 2 case, they are mostly suppressed thanks to the presence of the slope limiter. These oscillations become bigger for larger CFL values. This drawback has been pointed out at least in [84] or [59]. Some recent strategies that could be considered in order to reduce such oscillations can be found in [78], [58] or [71].

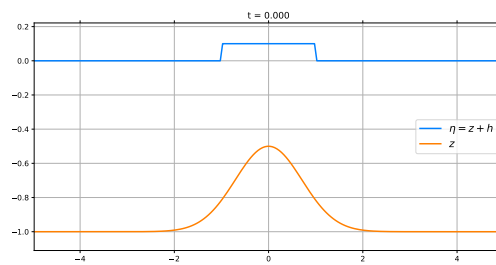


Figure 2.11: water at rest solution with a discontinuous perturbation on the surface

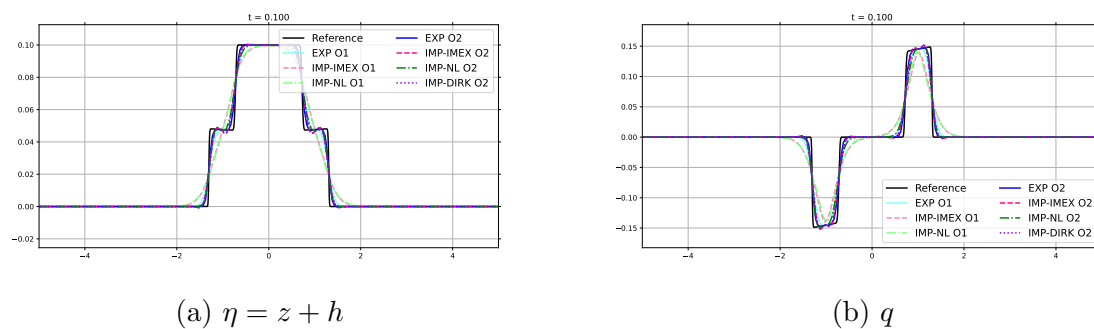


Figure 2.12: Solution at  $t = 0.1$  with 200 cells for a discontinuous perturbation of a water at rest. Explicit: CFL=0.5, Implicit: CFL=2

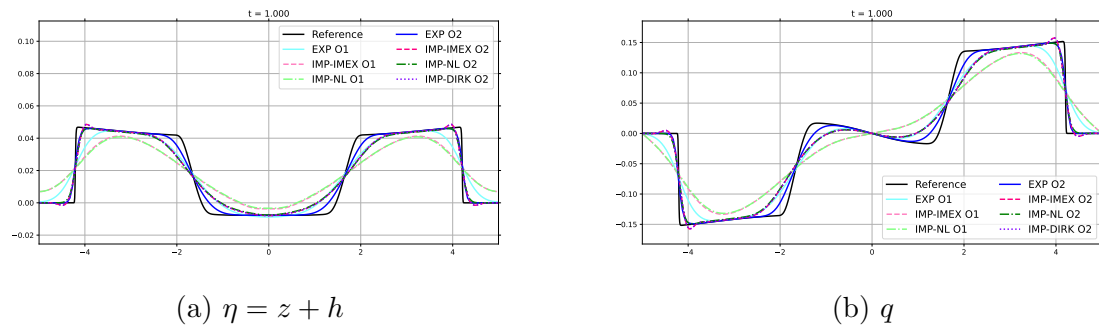


Figure 2.13: Solution at  $t = 1$  with 200 cells for a discontinuous perturbation of a water at rest. Explicit: CFL=0.5, Implicit: CFL=2

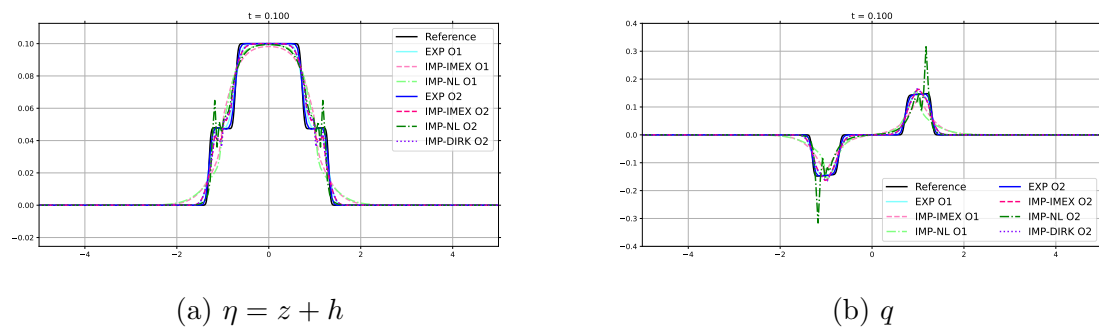


Figure 2.14: Solution at  $t = 0.1$  with 200 cells for a discontinuous perturbation of a water at rest. Explicit: CFL=0.5, Implicit: CFL=5

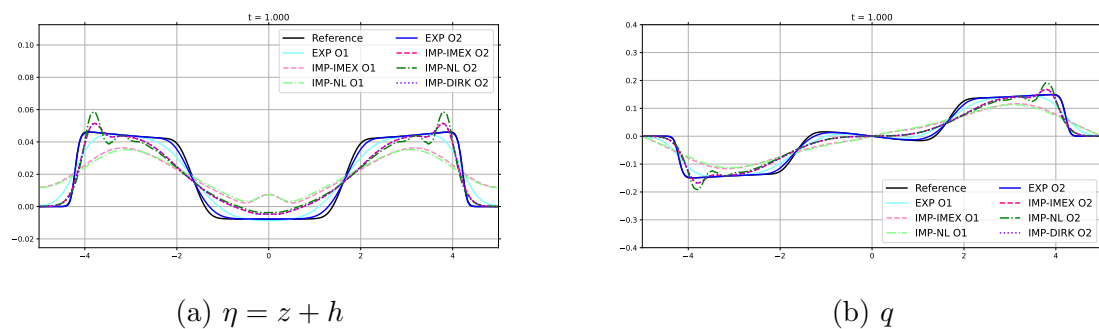


Figure 2.15: Solution at  $t = 1$  with 200 cells for a discontinuous perturbation of a water at rest. Explicit: CFL=0.5, Implicit: CFL=5.



### 2.4.6 Generation of subcritical steady state

Let us consider again the bottom topography (2.4.1) and a water at rest as initial condition in the interval  $[-5, 5]$ :

$$\eta = h + z = 0, u = 0.$$

At the left boundary, we impose a discharge given by  $q(x = -5, t) = 0.5$  and at the right boundary we fix the water height to  $h(x = 5, t) = 1$ , following the test proposed in [60]. This boundary conditions are set using a ghost-cell technique. The results obtained at time  $t = 2$ ,  $t = 50$  and  $t = 100$  for the different methods with 200 cells, with CFL value 0.5 for the explicit schemes and CFL value 2 and for the implicit ones, can be seen in Figures 2.16, 2.17, 2.18, 2.19, 2.20 and 2.21. Due to the imposed discharge on the left boundary, a shock wave enters the domain and travels over the bump. The solution evolves until a subcritical steady state is reached. We remark that the front of the advancing shock is too diffusive for the first order schemes. Moreover, the location of the depression on the surface at times  $t = 50$  and  $t = 100$  is slightly misplaced in the case of the implicit first order schemes. The left water level is not the same as well. This is due to the higher diffusion of the schemes, which is overcome when using a second order scheme.

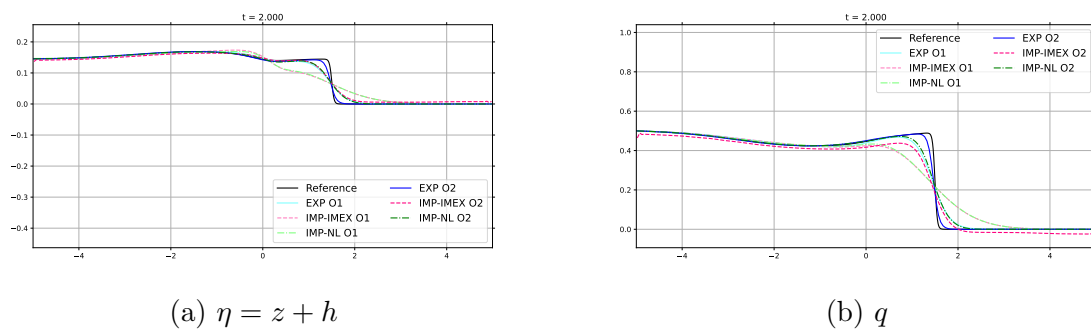


Figure 2.16: Generation of subcritical steady state a time  $t = 2$ . Explicit: CFL=0.5, Implicit: CFL=2

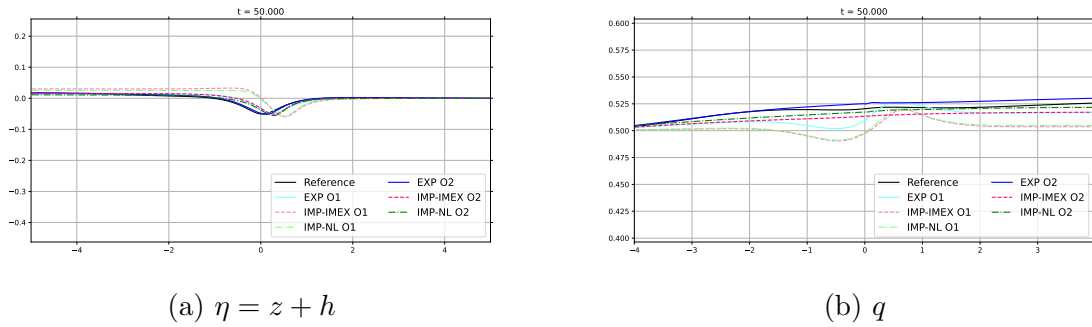


Figure 2.17: Generation of subcritical steady state a time  $t = 50$ , Explicit: CFL=0.5, Implicit: CFL=2

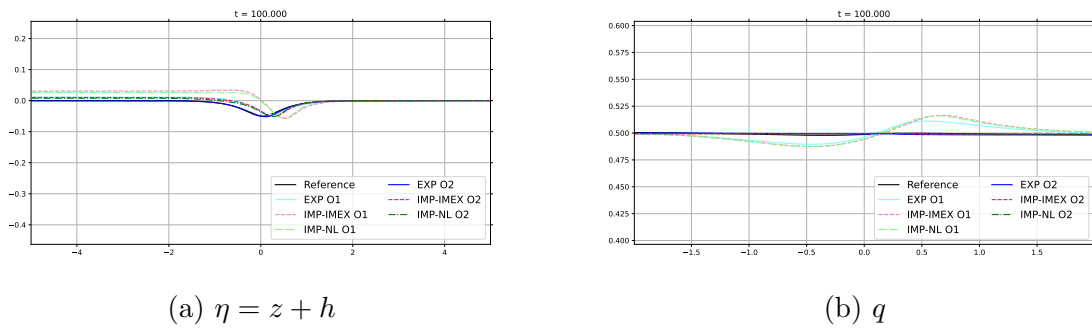


Figure 2.18: Generation of subcritical steady state a time  $t = 100$ . Explicit: CFL=0.5, Implicit: CFL=2

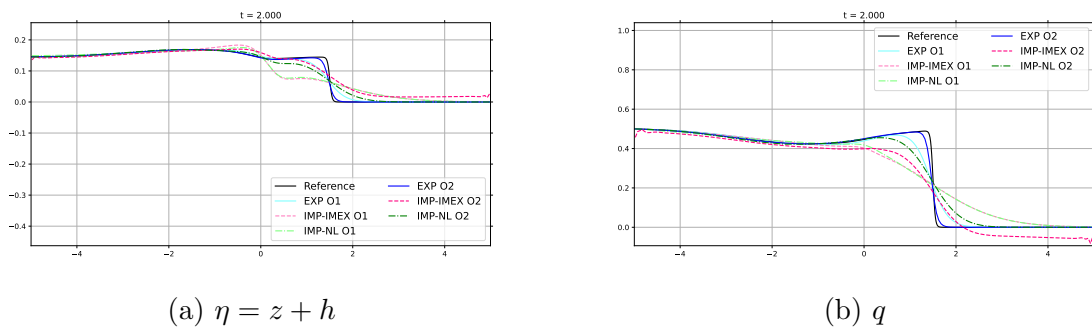
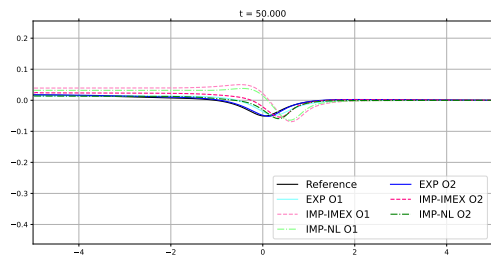
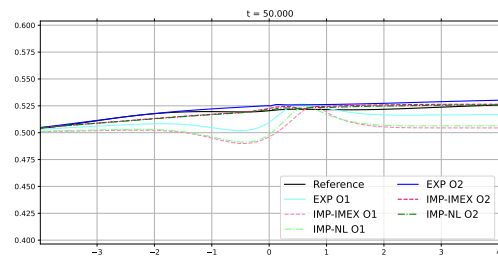


Figure 2.19: Generation of subcritical steady state a time  $t = 2$ . Explicit: CFL=0.5, Implicit: CFL=5



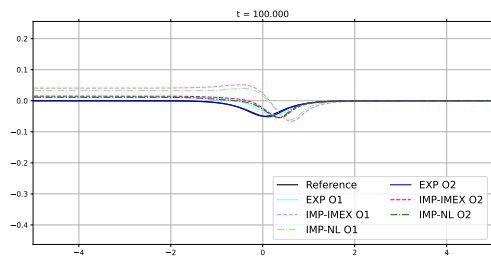


(a)  $\eta = z + h$

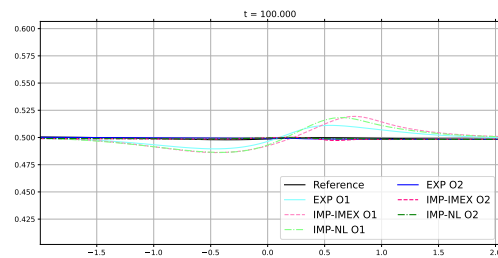


(b)  $q$

Figure 2.20: Generation of subcritical steady state a time  $t = 50$ , Explicit: CFL=0.5, Implicit: CFL=5



(a)  $\eta = z + h$



(b)  $q$

Figure 2.21: Generation of subcritical steady state a time  $t = 100$ . Explicit: CFL=0.5, Implicit: CFL=5



# Chapter 3

## Implicit Lagrange-projection well-balanced finite volume scheme for the Ripa model

### 3.1 Introduction

The Ripa system given by (1.1.22) simulates shallow water flows in situations where variations in temperature play a crucial role. This system was introduced in [79] and [80] for modeling ocean currents. The derivation of the system is based on considering multilayered ocean models, and vertically integrating the density, horizontal pressure gradient and velocity fields in each layer. The model incorporates the horizontal temperature gradients, which results in the variations in the fluid density within each layer. The simple case is the one considered here, where only one single layer is considered. The steady states of the 1D Ripa system satisfy (1.1.23). Here we are interested in the hydrostatic ones, that is those corresponding to  $u = 0$ .

In the following we propose a strategy to numerically solve the Ripa model by applying the Lagrange-projection approach as done in the previous chapter, following what has been done previously in [23, 72, 43, 33, 42, 22].

Some well-balanced numerical methods for this system are available in the literature, such as the central-upwind scheme in [37], in which steady states of the type (1.1.25) are preserved, or the HLLC type scheme in [82], which preserves both (1.1.25) and (1.1.26), and fits within the path-conservative schemes introduced in [75]. In this chapter we design first order finite volume schemes that are well-balanced for every hydrostatic steady state, that is, any steady states such that  $u = 0$ . These steady states satisfy (1.1.24).

We now consider the Lagrangian coordinates as introduced in Chapter 1.

In the case of smooth solutions, from (1.1.22) we have

$$\begin{cases} \partial_t h + u \partial_x h + h \partial_x u = 0, \\ \partial_t(hu) + u \partial_x(hu) + hu \partial_x u + \partial_x p + gh\theta \partial_x z = 0, \\ \partial_t(h\theta) + u \partial_x(h\theta) + h\theta \partial_x u = 0. \end{cases} \quad (3.1.1)$$

Mimicking what was done in Section 1.3 for Euler equations, we obtain the following system for Ripa model in Lagrangian coordinates:

$$\begin{cases} \partial_t(L\bar{h}) = 0, \\ \partial_t(L\bar{h}u) + \partial_\xi \bar{p} + g\bar{h}\theta \partial_\xi \bar{z} = 0, \\ \partial_t(L\bar{h}\theta) = 0. \end{cases} \quad (3.1.2)$$

The Lagrangian step can also be written (see [23]) in the following way:

$$\begin{cases} \partial_t \tau - \tau_0 \partial_\xi u = 0, \\ \partial_t u + \tau_0 \partial_\xi p + g \frac{\tau_0}{\tau} \theta \partial_\xi z = 0, \\ \partial_t \theta = 0, \end{cases} \quad (3.1.3)$$

where  $\tau = 1/h$  is the covolume and  $\tau_0 = \tau|_{t=0}$ .

Now, using a relaxation approach for (3.1.3), we obtain the system

$$\begin{cases} \partial_t \tau - \tau_0 \partial_\xi u = 0, \\ \partial_t u + \tau_0 \partial_\xi \pi + g \frac{\tau_0}{\tau} \theta \partial_\xi z = 0, \\ \partial_t \pi + a^2 \tau_0 \partial_\xi u = 0, \\ \partial_t \theta = 0, \end{cases} \quad (3.1.4)$$

where  $a$  is a constant that satisfies the subcharacteristic condition (1.1.65) and the variable  $\pi$  corresponds to the relaxation of the pressure  $p = \frac{g}{2} h^2 \theta$ .

Proceeding similarly as in Chapter 2 and using the variables  $\vec{w} = \pi + au$  and  $\overleftarrow{w} = \pi - au$ , we obtain

$$\begin{cases} \partial_t \tau - \tau_0 \partial_\xi u = 0, \\ \partial_t \vec{w} + a \tau_0 \partial_\xi \vec{w} + ga \frac{\tau_0}{\tau} \theta \partial_\xi z = 0, \\ \partial_t \overleftarrow{w} - a \tau_0 \partial_\xi \overleftarrow{w} + ga \frac{\tau_0}{\tau} \theta \partial_\xi z = 0, \\ \partial_t \theta = 0. \end{cases} \quad (3.1.5)$$

Let us recall that  $\pi$  and  $u$  can easily be recovered from  $\vec{w}$  and  $\overleftarrow{w}$  using (2.1.3).

We consider a similar numerical discretization as in Chapter 2 to use again the Lagrange-Projection approach, which we recall consists in, given discrete states  $U_i^n = (h, hu, h\theta)_i^n$ ,  $i \in \mathbb{Z}$ , corresponding to instant  $t^n$ , computing the approximations corresponding to time  $t^{n+1}$  in two steps:

1. Approximating the solution  $\bar{U}_i^{n+1}$  of system (3.1.2) (Lagrangian step).
2. Going back to the Eulerian coordinates to obtain the values  $U_i^{n+1}$  (projection step).

## 3.2 The Lagrangian step

In order to solve system (3.1.2), we rewrite the source term as we did in (2.2.1), obtaining

$$g\bar{h}\theta\partial_{\xi}\bar{z} = gL\bar{h}\theta\partial_x\bar{z} = g(h\theta)(0)\partial_x\bar{z},$$

with

$$\partial_x\bar{z}(\xi_i, t) = \partial_x z(x(\xi_i, t)) = z'(x_i(t)).$$

Using this, we can write a semi-discrete scheme for the Lagrangian step as follows:

$$\begin{cases} (L\bar{h})'_i(t) = 0, \\ (L\bar{h}u)'_i(t) = -\frac{1}{\Delta\xi} (\pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t)) - g(h\theta)_i(0)z'(x_i(t)), \\ (L\bar{h}\theta)'_i(t) = 0, \end{cases} \quad (3.2.1)$$

where  $\pi_{i\pm 1/2}^*(t) \approx \bar{\pi}(\xi_{i\pm 1/2}, t)$ .

Similarly, a semi-discrete scheme for the second and third equations of the relaxed system (3.1.5) can be considered:

$$\begin{cases} \vec{w}'_i(t) = -\frac{a}{h_i(0)\Delta\xi} (\vec{w}_{i+1/2}(t) - \vec{w}_{i-1/2}(t)) - ga\theta_i(0)z'(x_i(t)), \\ \overleftarrow{w}'_i(t) = \frac{a}{h_i(0)\Delta\xi} (\overleftarrow{w}_{i+1/2}(t) - \overleftarrow{w}_{i-1/2}(t)) + ga\theta_i(0)z'(x_i(t)). \end{cases} \quad (3.2.2)$$

where  $\vec{w}_{i+1/2}(t)$  and  $\overleftarrow{w}_{i+1/2}(t)$  are upwind numerical fluxes as those in (2.2.3).

Let us now focus on the well-balanced character of the scheme. In the case of a general hydrostatic steady state, unlike what happens with the shallow water equations for zero velocity steady states for which the steady states are explicitly determined by (1.1.19), there are infinitely many possibilities. We then have to fix a profile for one of the variables and obtain the other ones from it. For example, we can start by choosing a profile for  $h$ . A natural choice for first order schemes would be

$$h_i^e(x) = h_i^n, \quad (3.2.3)$$

or

$$h_i^e(x) = h_i^n + z_i - z(x). \quad (3.2.4)$$

In our case, we will work with (3.2.4). Then, we approximate  $\pi$  integrating (1.1.24) using a collocation method as in [56]. Therefore, we can compute the values  $\pi_i^{e,n}(x_{i\pm 1/2})$  as

$$\pi_i^{e,n}(x_{i\pm 1/2}) = \pi_i^{e,n}(x_i) \mp \Delta x \frac{\pi_i^{e,n}(x_i)}{h_i^{e,n}(x_i)} z'(x_i) \quad (3.2.5)$$

being

$$\pi_i^{e,n}(x_i) = \frac{g}{2} h_i^{e,n}(x_i) (h\theta)_i^{e,n}(x_i). \quad (3.2.6)$$

Afterwards, we are able to recover the values of  $\theta$  from the relation:

$$\theta(x) = \frac{2}{g} h^{-2}(x) \pi(x). \quad (3.2.7)$$

This could have also be done similarly by fixing a profile for  $\theta$  instead of doing it for  $h$ .

The fact that we are choosing a discrete profile for  $h$  and computing the other variables from it using a collocation method implies that our scheme will be well-balanced and not exactly well-balanced. That is, the steady states that are preserved are discrete approximations of the exact ones.

As already done in the previous chapter,  $\pi_{i+1/2}^{*,n+1}$  and  $u_{i+1/2}^{*,n+1}$  are computed as:

$$\pi_{i+1/2}^{*,n+1} = \frac{\overrightarrow{w}_{i+1/2-}^{n+1} + \overleftarrow{w}_{i+1/2+}^{n+1}}{2},$$

$$u_{i+1/2}^{*,n+1} = \frac{\overrightarrow{w}_{i+1/2-}^{n+1} - \overleftarrow{w}_{i+1/2+}^{n+1}}{2a},$$

where

$$\overrightarrow{w}_{i+1/2-}^{n+1} = \overrightarrow{w}_i^{n+1} - \pi_i^{e,n}(\xi_i) + \pi_i^{e,n}(\xi_{i+1/2}),$$

$$\overleftarrow{w}_{i+1/2+}^{n+1} = \overleftarrow{w}_{i+1}^{n+1} - \pi_{i+1}^{e,n}(\xi_i) + \pi_{i+1}^{e,n}(\xi_{i+1/2}).$$

Concerning system (3.1.2), note that, for stationary solutions, the second equation reads

$$\partial_\xi \overline{\pi}^e + g \overline{h} \theta^e \partial_\xi \overline{z} = 0,$$

or equivalently,

$$\partial_\xi \overline{\pi}^e + g L \overline{h} \theta^e \overline{\partial_x z} = 0. \quad (3.2.8)$$

Using this equality, the second equation of system (3.1.2) for a general solution can be written as

$$\partial_t (L \overline{h} u) + \partial_\xi (\overline{\pi} - \overline{\pi}^e) + g ((h\theta)(0) - L \overline{h} \theta^e) \overline{\partial_x z} = 0. \quad (3.2.9)$$

The same idea is used to rewrite the equations for  $\overrightarrow{w}$  and  $\overleftarrow{w}$ .

So, following the procedure used in Chapter 2 we consider the following semi-discrete formulation:

$$(L \overline{h} u)_i'(t) = -\mathcal{L}_i(t) - \mathcal{G}_i(t),$$

where

$$\mathcal{L}_i(t) = \frac{1}{\Delta \xi} (\pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t)),$$

$$\mathcal{G}_i(t) = -\frac{1}{\Delta \xi} (\pi_i^{e,n}(x_{i+1/2}(t)) - \pi_i^{e,n}(x_{i-1/2}(t)))$$

$$+ g (P_{i,(h\theta)_0}(\xi_i) - L_i(t) (h\theta)_i^{e,n}(x(\xi_i, t))) z'(x(\xi_i, t)),$$

and

$$\begin{cases} \vec{w}'_i(t) = -(\mathcal{L}_{\vec{w}})_i(t) - \frac{a}{h_i(0)}\mathcal{G}_i(t), \\ \overleftarrow{w}'_i(t) = -(\mathcal{L}_{\overleftarrow{w}})_i(t) + \frac{a}{h_i(0)}\mathcal{G}_i(t), \end{cases}$$

being

$$\begin{aligned} (\mathcal{L}_{\vec{w}})_i(t) &= \frac{a}{h_i(0)\Delta\xi} (\vec{w}_{i+1/2}(t) - \vec{w}_{i-1/2}(t)), \\ (\mathcal{L}_{\overleftarrow{w}})_i(t) &= -\frac{a}{h_i(0)\Delta\xi} (\overleftarrow{w}_{i+1/2}(t) - \overleftarrow{w}_{i-1/2}(t)). \end{aligned}$$

We will consider two different types of implicit schemes depending on how we treat functions  $\mathcal{L}$  and  $\mathcal{G}$ :

- Nonlinear implicit schemes: both  $\mathcal{L}$  and  $\mathcal{G}$  are treated implicitly.
- Implicit-explicit schemes:  $\mathcal{L}$  is treated implicitly while  $\mathcal{G}$  is treated explicitly.

### 3.2.1 First order nonlinear implicit well-balanced Lagrangian scheme

In the case in which we treat functions  $\mathcal{L}$  and  $\mathcal{G}$  implicitly, we consider

$$(\overline{Lhu})_i^{n+1} = (hu)_i^n - \Delta t (\mathcal{L}_i^{n+1} + \mathcal{G}_i^{n+1}),$$

with

$$\begin{aligned} \mathcal{L}_i^{n+1} &= \frac{1}{\Delta\xi} (\pi_{i+1/2}^{*,n+1} - \pi_{i-1/2}^{*,n+1}), \\ \mathcal{G}_i^{n+1} &= -\frac{1}{\Delta\xi} (\pi_i^{e,n}(x_{i+1/2}^{*,n+1}) - \pi_i^{e,n}(x_{i-1/2}^{*,n+1})) \\ &\quad + g((h\theta)_i^n - L_i^{n+1}(h\theta)_i^{e,n}(x_i^{*,n+1})) z'(x_i^{*,n+1}). \end{aligned}$$

being

$$x_{i\pm 1/2}^{*,n+1} = \xi_{i\pm 1/2} + \Delta t u_{i\pm 1/2}^{*,n+1},$$

and

$$L_i^{n+1} = 1 + \frac{\Delta t}{\Delta\xi} (u_{i+1/2}^{*,n+1} - u_{i-1/2}^{*,n+1}).$$

Moreover, in order to obtain  $\vec{w}_i^{n+1}$  and  $\overleftarrow{w}_i^{n+1}$  we need to solve the non-linear systems defined by

$$\begin{cases} \vec{w}_i^{n+1} = \vec{w}_i^n - \Delta t \left( (\mathcal{L}_{\vec{w}})_i^{n+1} + \frac{a}{h_i^n} \mathcal{G}_i^{n+1} \right), \\ \overleftarrow{w}_i^{n+1} = \overleftarrow{w}_i^n - \Delta t \left( (\mathcal{L}_{\overleftarrow{w}})_i^{n+1} - \frac{a}{h_i^n} \mathcal{G}_i^{n+1} \right), \end{cases}$$

where

$$\begin{aligned} (\mathcal{L}_{\vec{w}})_i^{n+1} &= \frac{a}{h_i^n \Delta \xi} \left( \vec{w}_{i+1/2-}^{n+1} - \vec{w}_{i-1/2-}^{n+1} \right), \\ (\mathcal{L}_{\overleftarrow{w}})_i^{n+1} &= -\frac{a}{h_i^n \Delta \xi} \left( \overleftarrow{w}_{i+1/2+}^{n+1} - \overleftarrow{w}_{i-1/2+}^{n+1} \right). \end{aligned}$$

These systems are solved using a fixed point scheme as in Chapter 2. However, the results obtained there show that implicit-explicit schemes have the advantage of being more efficient as only linear systems need to be solved. Therefore from now on we focus on implicit-explicit schemes.

### 3.2.2 First order implicit-explicit well-balanced Lagrangian scheme

In order to treat  $\mathcal{L}$  implicitly and  $\mathcal{G}$  explicitly, we consider the scheme

$$(L\overline{hu})_i^{n+1} = (hu)_i^n - \Delta t (\mathcal{L}_i^{n+1} + \mathcal{G}_i^n),$$

where

$$\begin{aligned} \mathcal{L}_i^{n+1} &= \frac{1}{\Delta \xi} \left( \pi_{i+1/2}^{*,n+1} - \pi_{i-1/2}^{*,n+1} \right), \\ \mathcal{G}_i^n &= -\frac{1}{\Delta \xi} \left( \pi_i^{e,n}(\xi_{i+1/2}) - \pi_i^{e,n}(\xi_{i-1/2}) \right). \end{aligned}$$

In this case,  $\vec{w}_i^{n+1}$  and  $\overleftarrow{w}_i^{n+1}$  are the solution of the linear systems

$$\begin{cases} \vec{w}_i^{n+1} = \vec{w}_i^n - \Delta t \left( (\mathcal{L}_{\vec{w}})_i^{n+1} + \frac{a}{h_i^n} \mathcal{G}_i^n \right), \\ \overleftarrow{w}_i^{n+1} = \overleftarrow{w}_i^n - \Delta t \left( (\mathcal{L}_{\overleftarrow{w}})_i^{n+1} - \frac{a}{h_i^n} \mathcal{G}_i^n \right), \end{cases}$$

where

$$\begin{aligned} (\mathcal{L}_{\vec{w}})_i^{n+1} &= \frac{a}{h_i^n \Delta \xi} \left( \vec{w}_{i+1/2-}^{n+1} - \vec{w}_{i-1/2-}^{n+1} \right), \\ (\mathcal{L}_{\overleftarrow{w}})_i^{n+1} &= -\frac{a}{h_i^n \Delta \xi} \left( \overleftarrow{w}_{i+1/2+}^{n+1} - \overleftarrow{w}_{i-1/2+}^{n+1} \right). \end{aligned}$$

Therefore, by using this implicit-explicit strategy we avoid the need to solve a nonlinear system and to use fixed point schemes, obtaining a more efficient scheme similarly to what was seen for the shallow water equations case in Chapter 2.

### 3.3 The projection step

Once the Lagrangian step has been performed, we need to project the solutions back to Eulerian coordinates as it was done in Section 2.3, but considering now the appropriate reconstructions of the variables.

### 3.4 Numerical results

Some tests will now be presented in order to check the behaviour of the schemes presented in this chapter. We will consider three different Lagrange-Projection schemes:

- A non well-balanced scheme obtained by using an HLL scheme for the Lagrangian step. This scheme should perform adequately in the isobaric steady state case since when applying the Lagrange-Projection approach, the Lagrangian step (3.1.2) with flat topography and constant pressure reduces to

$$\begin{cases} \partial_t(L\bar{h}) = 0, \\ \partial_t(L\bar{h}u) = 0, \\ \partial_t(L\bar{h}\theta) = 0, \end{cases}$$

and the projection step presents no issues, so any consistent scheme should do it right. However, it should present difficulties with other type of steady states.

- An exactly well-balanced scheme that exactly preserves water at rest steady states (1.1.25) and isobaric steady states (1.1.26). This one is also obtained by using an HLL scheme for the Lagrangian step. Of course, it should not be able to preserve general hydrostatic steady states. However, if a discrete steady state is considered, it could present some errors of the order of the discrete approximation.
- Our implicit-explicit well-balanced scheme that is well-balanced for all the hydrostatic steady states. Since it is not exactly well-balanced it could present some errors when considering a exact steady state as initial condition, but it should converge to the steady state when the mesh is refined. Moreover, if the discrete approximation of a hydrostatic steady state is considered, it should exactly preserve it.

#### 3.4.1 Isobaric steady state

In this first test we consider a flat topography

$$z(x) = 0.$$

and the following initial conditions:

$$\begin{aligned} h(x, 0) &= 1 - 0.5e^{-x^2}, \\ u(x, 0) &= 0, \\ \theta(x, 0) &= (1 - 0.5e^{-x^2})^{-2}. \end{aligned}$$

The ones corresponding to  $h$  and  $\theta$  are plotted in Figure 3.1. The computational domain is the interval  $[-5, 5]$  and the final time is  $T = 2$ .

The aim of this test is to check that our scheme preserves isobaric steady states, which are the ones satisfying (1.1.26).

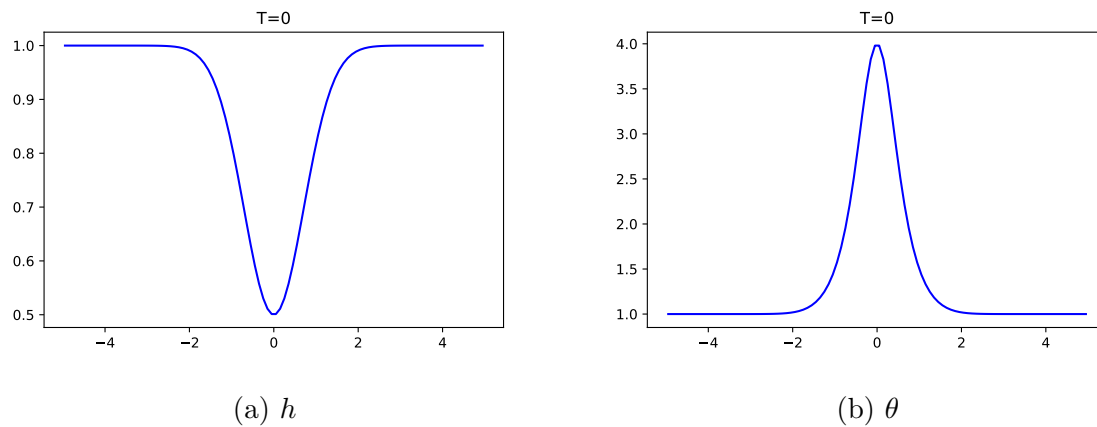


Figure 3.1: Initial condition for an isobaric steady state

In Tables 3.1, 3.2 and 3.2 we find the errors for  $h$ ,  $hu$  and  $h\theta$  obtained with the non well-balanced scheme, the exactly well-balanced one and our implicit-explicit well-balanced scheme, respectively. As expected, in the absence of source term, the three schemes are able to preserve the steady state.

No. of cells	$h$	$hu$	$h\theta$
100	2.77e-15	1.43e-15	3.73e-15
200	3.53e-15	1.35e-15	4.99e-15
400	5.83e-15	1.58e-15	7.12e-15

Table 3.1: Errors in  $L^1$  norm for the non well-balanced scheme for an isobaric steady state

No. of cells	$h$	$hu$	$h\theta$
100	8.21e-16	8.42e-16	4.44e-16
200	6.55e-16	7.41e-16	1.77e-16
400	1.09e-15	8.71e-16	7.54e-16

Table 3.2: Errors in  $L^1$  norm for the exactly well-balanced scheme for an isobaric steady state

No. of cells	$h$	$hu$	$h\theta$
100	8.65e-16	1.40e-15	5.32e-16
200	6.88e-16	7.19e-16	1.77e-16
400	1.15e-15	8.96e-16	8.43e-16

Table 3.3: Errors in  $L^1$  norm for the well-balanced scheme for an isobaric steady state

### 3.4.2 Water at rest case

We consider the non-flat bottom topography

$$z(x) = 0.5e^{-x^2}$$

and the initial conditions

$$\begin{aligned} h(x, 0) &= 1 - 0.5e^{-x^2}, \\ u(x, 0) &= 0, \\ \theta(x, 0) &= 1. \end{aligned}$$

So, now the free surface  $h + z$  is constant and this is a water at rest type steady state, satisfying (1.1.25). We will check if our scheme is well-balanced for this type of stationary solution.

Again, the computational domain is  $[-5, 5]$  and the final time is  $T = 2$ .

On the one hand, the profiles of the free surface obtained with a non well balanced scheme and with the proposed well balanced scheme are similar, as we can see in Figure 3.2. Moreover, no big differences are seen for the temperature either (see Figure 3.4). On the other hand, the velocity obtained with the non well-balanced scheme is very far away from the one that we should obtain, while for the well-balanced scheme we obtain errors of order  $10^{-13}$  (see Figure (3.3)). This behaviour can also be observed in Tables 3.4, 3.5 and 3.6, where we show the errors obtained with the non well-balanced scheme, the exactly well-balanced one and our well-balanced scheme. The well-balanced scheme is, of course, non exact, but the errors decrease when the mesh is refined, so it is expected to converge to the stationary solution. In fact, the convergence rate is 2, as it can be deduced from Table 3.7.

Of course, the exactly well-balanced scheme perfectly fits this situation, while the well-balanced scheme is not exact but it is better than the non well-balanced one, in particular for approximating the discharge. In fact, if we had started from a discrete stationary solution instead of an exact solution, the well-balanced scheme would have been exact and we would have seen some small errors in the exactly well-balanced one.

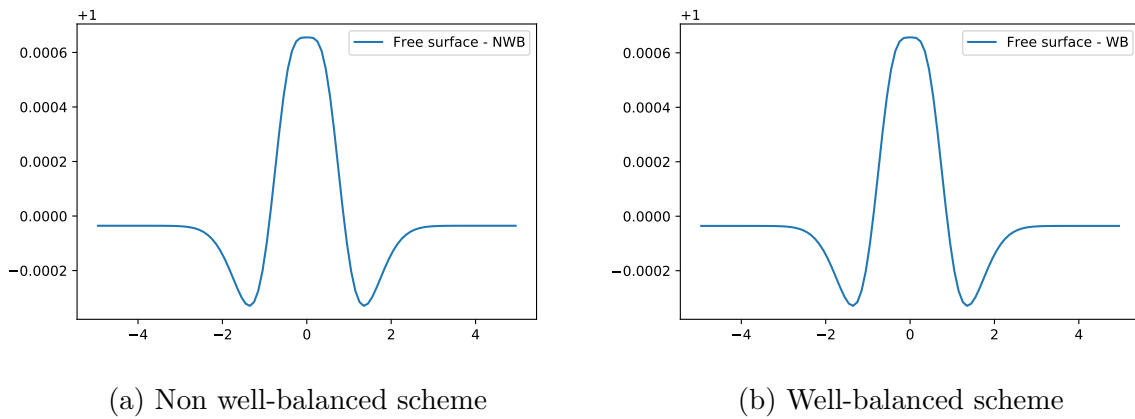


Figure 3.2: Free surface computed with well-balanced and non well-balanced schemes for the water at rest steady state

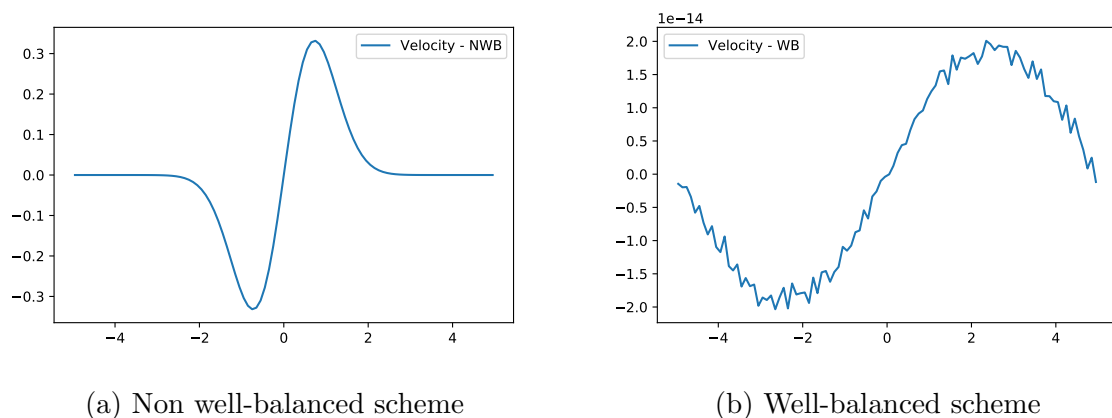
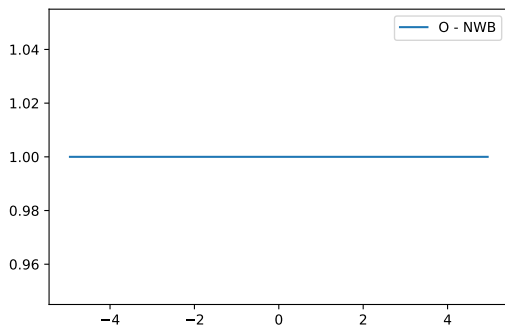
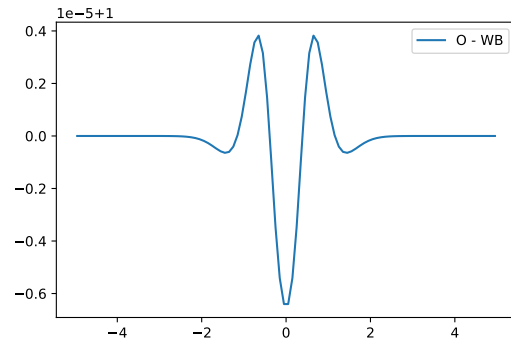


Figure 3.3: Velocity computed with well-balanced and non well-balanced schemes for the water at rest steady state



(a) Non well-balanced scheme



(b) Well-balanced scheme

Figure 3.4: Temperature computed with well-balanced and non well-balanced schemes for the water at rest steady state

No. of cells	$h$	$hu$	$h\theta$
50	6.64e-03	8.67e-01	6.64e-03
100	1.16e-03	5.80e-01	1.16e-03
200	4.04e-04	3.99e-01	4.04e-04
400	1.01e-04	2.81e-01	1.01e-04

Table 3.4: Errors in  $L^1$  norm for the non well-balanced scheme for the water at rest steady state

No. of cells	$h$	$hu$	$h\theta$
50	7.54e-16	3.67e-15	7.54e-16
100	1.99e-15	5.14e-15	1.99e-15
200	3.63e-15	7.59e-15	3.63e-15
400	3.36e-15	4.41e-15	3.36e-15

Table 3.5: Errors in  $L^1$  norm for the exactly well-balanced scheme for the water at rest steady state

No. of cells	$h$	$hu$	$h\theta$
50	6.46e-03	2.23e-13	6.46e-03
100	1.16e-03	2.99e-13	1.16e-03
200	4.04e-04	1.55e-13	4.04e-04
400	1.01e-04	3.12e-13	1.01e-04

Table 3.6: Errors in  $L^1$  norm for the well-balanced scheme for the water at rest steady state

No. of cells	Order for $h$	Order for $h\theta$
100	2.47	2.47
200	1.52	1.52
400	2.00	2.00

Table 3.7: Convergence rates for the well-balanced scheme for the water at rest steady state

### 3.4.3 A general hydrostatic steady state case

The topography considered now is

$$z(x) = 1 - 0.5e^{-x^2}.$$

From the initial conditions

$$h(x, 0) = 1 - 0.5e^{-x^2},$$

$$u(x, 0) = 0,$$

$$\pi(x, 0) = (1 - 0.5e^{-x^2})^{-2}$$

we deduce that the temperature  $\theta$  is no longer constant, and this is a general steady state, not a water at rest steady state or an isobaric one. Let us check if this general steady state is preserved by our scheme.

Although we know the exact solution, we start from the discrete solution in order to test our scheme (see Figure 3.5 and Table 3.8 to compare them). We use the same space interval and final time of previous examples.

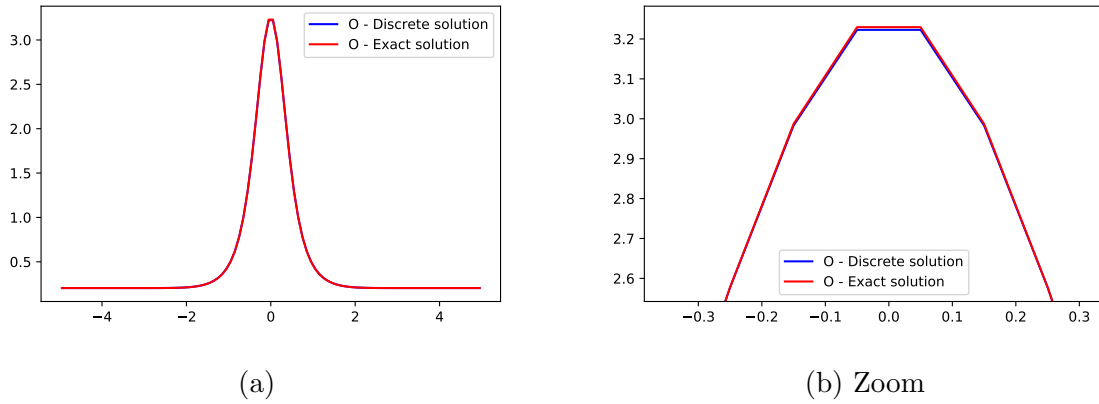
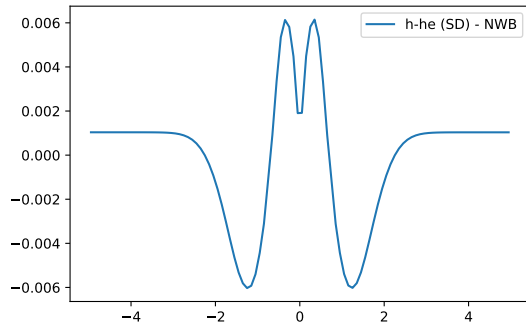


Figure 3.5: Exact and discrete solution for  $\theta$  for a general hydrostatic steady state case

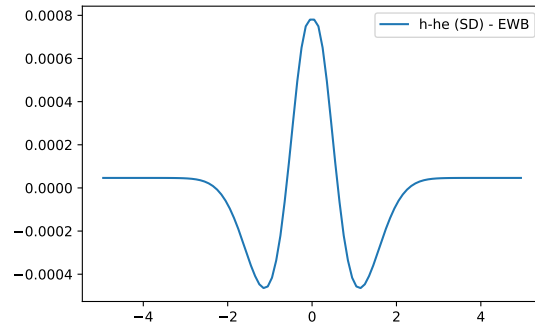
No. of cells	Error	Order
50	1.61e-02	
100	3.99e-03	2.01
200	1.00e-03	1.99
400	2.50e-04	2.00
800	6.25e-05	2.00

Table 3.8: Error and convergence rates in the approximation of the global discrete solution for a general hydrostatic steady state case

In Figures 3.6, 3.7 and 3.8 we have the final states obtained for  $h$ ,  $u$  and  $\theta$ , and in Tables 3.9, 3.10 and 3.11 the errors for the three considered schemes. Of course, the errors in the exactly well-balanced case are smaller than in the non well-balanced one, especially in the case of the discharge. However, our well-balanced scheme performs better, preserving the discrete stationary solution up to machine precision.

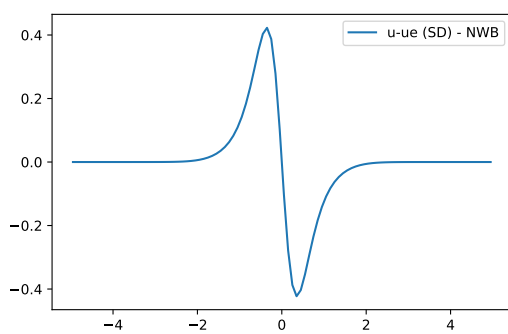


(a) Non well-balanced

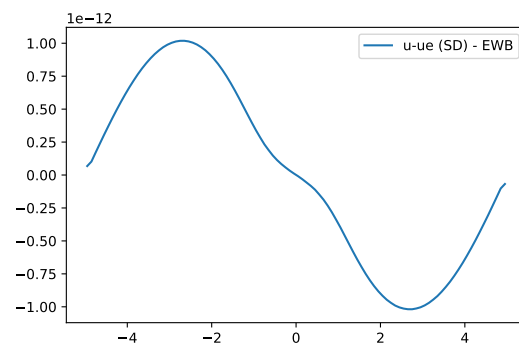


(b) Exactly well-balanced

Figure 3.6: Error in  $h$  for the non well-balanced and exactly well-balanced schemes for a general hydrostatic steady state case



(a) Non well-balanced



(b) Exactly well-balanced

Figure 3.7: Error in  $u$  for the non well-balanced and exactly well-balanced schemes for a general hydrostatic steady state case

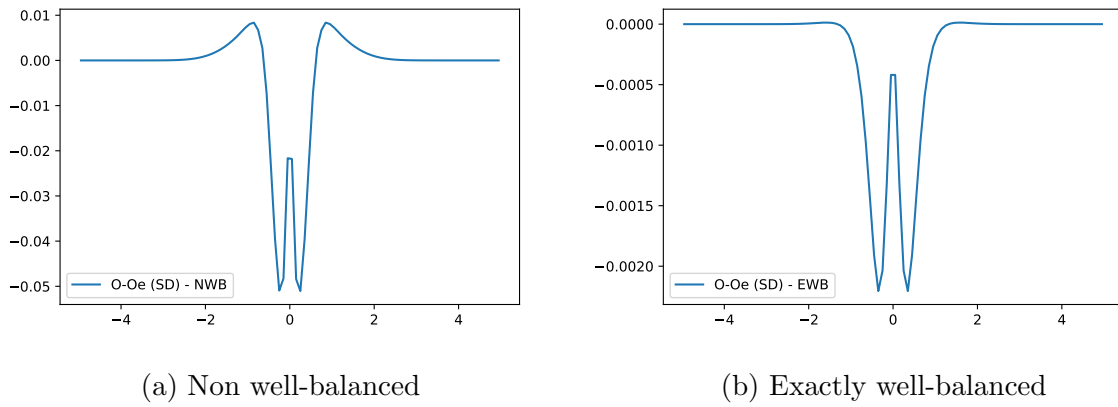


Figure 3.8: Error in  $\theta$  for the non well-balanced and exactly well-balanced schemes for a general hydrostatic steady state case

No. of cells	$h$	$hu$	$h\theta$
50	4.93e-02	6.52e-01	4.11e-02
100	2.15e-02	4.20e-01	1.72e-02
200	9.29e-03	2.82e-01	7.21e-03
400	6.53e-03	1.91e-01	6.49e-03

Table 3.9: Errors in  $L^1$  norm for the non well-balanced scheme for a general hydrostatic steady state case

No. of cells	$h$	$hu$	$h\theta$
50	7.33e-03	1.37e-13	6.84e-03
100	1.69e-03	1.49e-13	1.18e-03
200	4.07e-04	1.00e-13	4.81e-04
400	4.97e-05	1.10e-13	1.22e-04

Table 3.10: Errors in  $L^1$  norm for the exactly well-balanced scheme for a general hydrostatic steady state case

No. of cells	$h$	$hu$	$h\theta$
50	4.47e-17	1.04e-15	1.22e-16
100	2.22e-16	1.04e-15	9.99e-17
200	1.22e-16	1.20e-15	1.24e-16
400	4.99e-17	1.94e-15	5.41e-17

Table 3.11: Errors in  $L^1$  norm for the well-balanced scheme for a general hydrostatic steady state case

### 3.4.4 A general steady hydrostatic state with a perturbation

In this last test we add a small perturbation in  $h$  to the discrete solution of the previous test. In particular, we consider the initial conditions

$$\begin{aligned} h(x_i, 0) &= h^e(x_i) + 0.1e^{-16(x_i-2)^2}, \\ (hu)(x_i, 0) &= 0, \\ (h\theta)(x_i, 0) &= (h\theta^e)(x_i), \end{aligned}$$

where the superscript  $e$  refers to those discrete stationary solutions of the previous test.

We consider open boundary conditions and let the simulation run until a stationary solution is reached.

In Figure 3.9 we have plotted  $h$  at different times obtained with the well-balanced scheme. Moreover, in Figures 3.10 and 3.11 we observe the final steady state obtained with this scheme for variables  $h$  and  $\theta$ , respectively. Lastly, in Figures 3.12 and 3.13 we have plotted the differences between the non well-balanced and the well-balanced scheme and the differences between the exactly well-balanced and the well-balanced scheme for the variables  $u$  and  $\theta$ .

On one hand, the well-balanced scheme generates from that initial condition another hydrostatic solution different from the perturbed starting one, which is expected, since there are infinitely many of that type and it goes to any one of those. On the other hand, the exactly well-balanced scheme is only able to generate a stationary solution very similar to the initial one, as it is not well balanced for the discrete stationary solution. Finally, the non well-balanced scheme does not obtain a solution with zero velocity, but generates a large perturbation in  $u$ , as expected.

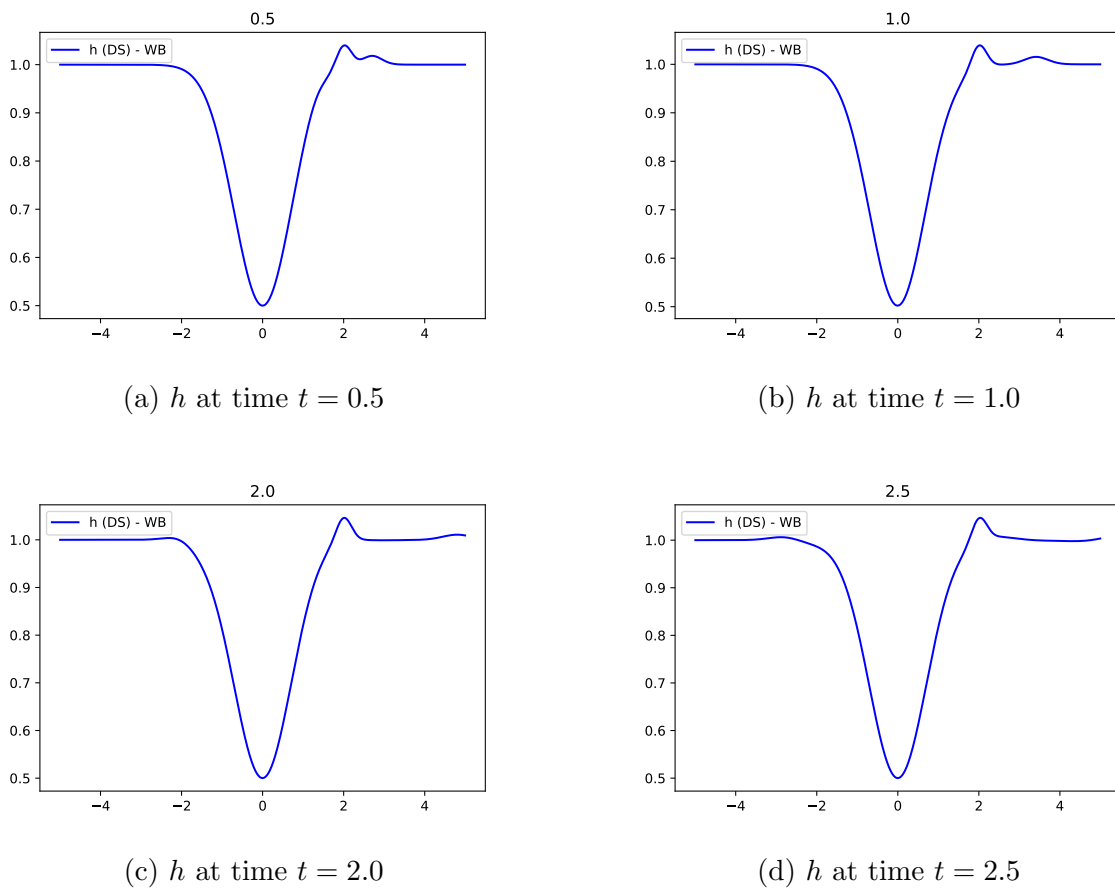


Figure 3.9:  $h$  at different times for the well-balanced scheme for a general hydrostatic steady state case with a perturbation

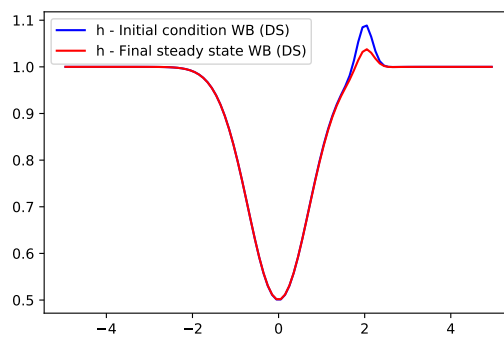


Figure 3.10: Final steady state for  $h$  for the well-balanced scheme for a general hydrostatic steady state case with a perturbation

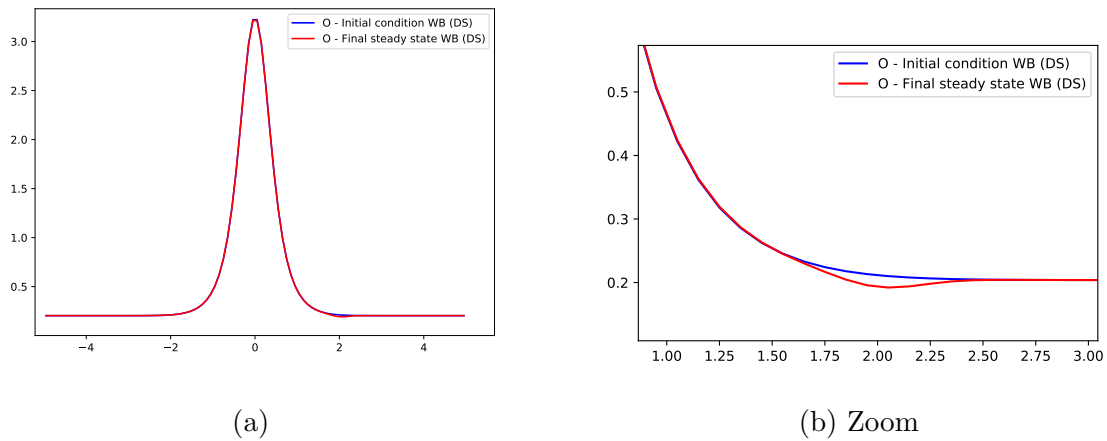
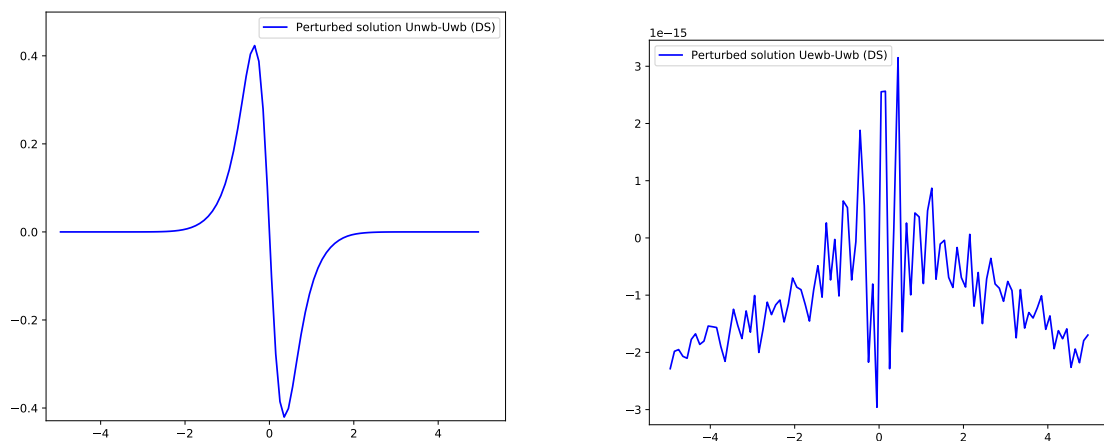
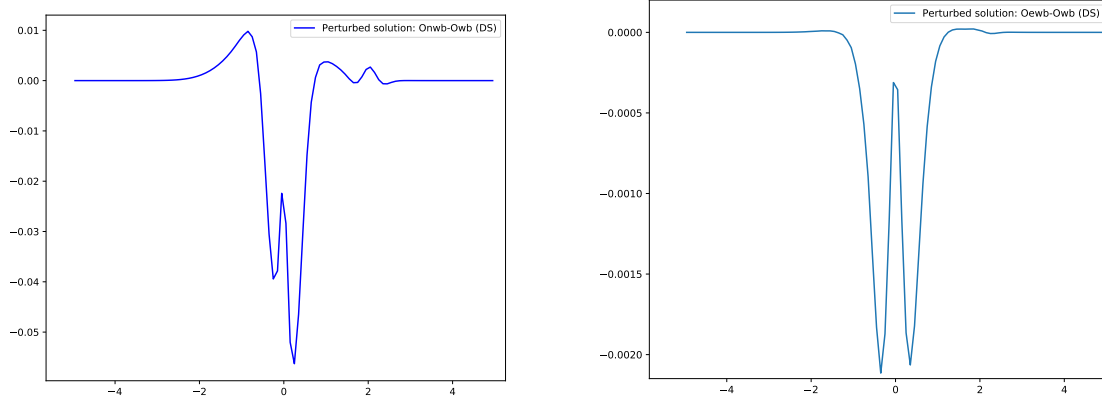


Figure 3.11: Final steady state for  $\theta$  for the well-balanced scheme for a general hydrostatic steady state case with a perturbation



(a) Difference between the non well-balanced and the well-balanced scheme for  $u$ . (b) Difference between the exactly well-balanced and the well-balanced scheme for  $u$ .

Figure 3.12: Difference for  $u$



(a) Difference between the non well-balanced and the well-balanced scheme for  $\theta$ . (b) Difference between the exactly well-balanced and the well-balanced scheme for  $\theta$ .

Figure 3.13: Difference for  $\theta$



# Chapter 4

## Semi-implicit fully exactly well-balanced finite volume schemes for the 1D shallow water system

### 4.1 Introduction

In Chapter 2 we presented first and second order implicit and implicit-explicit numerical schemes that preserved water at rest steady states of the SWE. We now aim to design semi-implicit schemes that are fully exactly well-balanced for the 1D shallow water equations, that is, schemes that exactly preserve all the smooth steady states of the system.

Different techniques for the design of implicit or semi-implicit schemes for the shallow water models have been developed since [26]. These approaches include the application of finite volume methods in studies such as [9, 31, 91], a discontinuous Galerkin (DG) approach as seen in the works [46, 53, 89, 64, 88], finite difference methods as explored in Casulli's earlier work [26, 29], and hybrid strategies, as demonstrated in studies like those by [11, 21, 17]. In general, the idea consists on performing a splitting that allows to separate the fast waves from the slow ones, and on combining explicit and implicit schemes. However, to the best of our knowledge, no previous work has been presented in which a semi-implicit scheme preserves all the steady states of the one-dimensional shallow water equations.

In Chapter 2, the Lagrangian formalism was used in order to define semi-implicit schemes that preserve water at rest stationary solutions for the shallow water equations. The Lagrangian formalism was also used in [23] to define an explicit first order fully well balanced scheme, while in [72] an explicit high order well-balanced scheme was presented. Nevertheless, the use of the Lagrangian formalism complicates in excess the task of defining high-order, semi-implicit and fully well-balanced schemes. Even in [23] particular care had to be taken in the projection step in order to obtain a fully well-balanced scheme, as a steady-state in Eulerian coordinates is not necessarily a steady-state in Lagrangian

coordinates. Moreover, extending this technique to 2D cases could be cumbersome. Therefore, we propose here a different approach that overcomes these difficulties.

We aim then to define a semi-implicit fully well-balanced scheme, specially adapted for small Froude situations. To do so, a splitting strategy inspired on [47] combined with a relaxation technique will be used, which will be described in Section 4.2, where we will also be concerned with the design of schemes that are well balanced. In Sections 4.3 and 4.4, the proposed first and second order schemes are presented, respectively. Finally, several numerical experiments are shown in Section 4.5 in order to test the accuracy and efficiency of the schemes presented.

## 4.2 Splitting and relaxation techniques

We will start by performing a splitting of the shallow water system (1.1.15). In order to set the basic ideas that will be detailed afterwards, let us rewrite the system as

$$\partial_t U = S_P(U, z) + S_T(U),$$

where  $U = (h, hu)^T$ ,

$$S_P(U, z) = \begin{bmatrix} 0 \\ -\partial_x \left( \frac{1}{2} gh^2 \right) - gh z' \end{bmatrix} \quad (4.2.1)$$

and

$$S_T(U) = \begin{bmatrix} -\partial_x(hu) \\ -\partial_x(hu^2) \end{bmatrix}. \quad (4.2.2)$$

The splitting strategy consists in solving each of the two systems

$$\partial_t U = S_P(U, z), \quad (4.2.3)$$

and

$$\partial_t U = S_T(U). \quad (4.2.4)$$

sequentially.

System (4.2.3) will be referred to as the pressure system, while system (4.2.4) will be referred to as the transport system.

We could either solve the system defined by  $S_P$  first, followed by the one defined by  $S_T$  or vice versa.

The advantage of applying this splitting approach is that it decouples the acoustic and the transport phenomena. This way, the pressure system can be solved both explicitly or implicitly, while the transport system will always be solved explicitly. Performing the pressure step implicitly allows us to consider larger time steps since it involves a less restrictive CFL condition. Indeed, for small Froude numbers, the main restriction on the

time step is driven by the pressure term. The reason why this happens is that system (4.2.3) has eigenvalues

$$\lambda = \pm\sqrt{gh}, \quad (4.2.5)$$

while system (4.2.4) has a double eigenvalue given by

$$\lambda = u. \quad (4.2.6)$$

Furthermore, this strategy offers advantages compared to the Lagrangian-Projection strategy, as it exclusively considers Eulerian coordinates, eliminating the necessity to deal with steady states dependent on time.

Let us now describe in detail how each of the systems will be solved. First, we consider the general framework of finite volume schemes, already described in the previous chapters: the computational domain is discretized in a set of cells  $[x_{i-1/2}, x_{i+1/2})$  for  $i \in \mathbb{Z}$  using, for the sake of simplicity, a constant volume length  $\Delta x = x_{i+1/2} - x_{i-1/2}$ . The values  $x_{i+1/2}$  correspond to the intercells, while the centers of the volume cells will be denoted by  $x_i = (x_{i-1/2} + x_{i+1/2})/2$ . The time variable will be kept continuous for now and the approximation of the cell averages will be denoted by

$$(h_i(t), (hu)_i(t))^T = U_i(t) \approx \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} U(x, t) dx.$$

A semi-discrete finite volume scheme for systems (4.2.3)-(4.2.4) can be written as

$$\begin{cases} h'_i(t) = 0, \\ (hu)'_i(t) = -\frac{1}{\Delta x} \left( \pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t) \right) + S_i(t), \end{cases} \quad (4.2.7)$$

$$\begin{cases} h'_i(t) = -\frac{1}{\Delta x} \left( h_{i+1/2}^*(t) u_{i+1/2}^*(t) - h_{i-1/2}^*(t) u_{i-1/2}^*(t) \right), \\ (hu)'_i(t) = -\frac{1}{\Delta x} \left( (hu)_{i+1/2}^*(t) u_{i+1/2}^*(t) - (hu)_{i-1/2}^*(t) u_{i-1/2}^*(t) \right), \end{cases} \quad (4.2.8)$$

or in compact form as

$$\frac{d}{dt} U_i(t) = S_{P_i}(t), \quad (4.2.9)$$

$$\frac{d}{dt} U_i(t) = S_{T_i}(t), \quad (4.2.10)$$

where

$$S_{P_i}(t) = \begin{pmatrix} 0 \\ -\frac{1}{\Delta x} \left( \pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t) \right) + S_i(t) \end{pmatrix},$$

$$S_{Ti}(t) = -\frac{1}{\Delta x} (U_{i+1/2}^*(t)u_{i+1/2}^*(t) - U_{i-1/2}^*(t)u_{i-1/2}^*(t)).$$

Here,  $U_{i+1/2}^*(t) = (h_{i\pm 1/2}^*(t), (hu)_{i\pm 1/2}^*(t))^T$ , being  $h_{i\pm 1/2}^*(t)$  and  $(hu)_{i\pm 1/2}^*(t)$  approximations at the interface of the height  $h$  and the discharge  $hu$ , that is,  $U_{i+1/2}^*(t) \approx (h(x_{i\pm 1/2}, t), (hu)(x_{i\pm 1/2}, t))^T$ , while  $\pi_{i\pm 1/2}^*(t)$  and  $u_{i\pm 1/2}^*(t)$  are approximations at the interface of the pressure  $\pi = \frac{g}{2}h^2$  and the velocity  $u$ , respectively:  $\pi_{i\pm 1/2}^*(t) \approx \pi(x_{i\pm 1/2}, t)$  and  $u_{i\pm 1/2}^*(t) \approx u(x_{i\pm 1/2}, t)$ . The approximation of the source term is denoted by  $S_i(t)$ , that is,

$$S_i(t) \approx \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} gh(x, t)z'(x)dx.$$

In order to appropriately define the values  $\pi_{i\pm 1/2}^*(t)$  and  $u_{i\pm 1/2}^*(t)$  and solve system defined (4.2.9), we will make use of relaxation techniques following the ideas presented in 1.1.3.2 as was done in Chapters 2 and 3 for the Lagrangian systems. More explicitly, we propose a relaxed system for (4.2.3):

$$\begin{cases} \partial_t h = 0, \\ \partial_t(hu) + \partial_x \pi = -gh z', \\ \partial_t(h\pi) + a^2 \partial_x u = 0, \end{cases} \quad (4.2.11)$$

where  $a$  is a constant satisfying the subcharacteristic condition (1.1.65). In practice, this means that

$$a \geq h\sqrt{gh}. \quad (4.2.12)$$

Now, using again the framework of finite volume methods and keeping the time variable continuous, we shall define the semi-discrete scheme for (4.2.11):

$$\begin{cases} h_i'(t) = 0, \\ (hu)_i'(t) = -\frac{1}{\Delta x} (\pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t)) - S_i(t), \\ (h\pi)_i'(t) = -\frac{a^2}{\Delta x} (u_{i+1/2}^*(t) - u_{i-1/2}^*(t)). \end{cases} \quad (4.2.13)$$

To do so, as  $h_i$  is constant through time,  $h_i$  is frozen at  $t = t_0$ . In practice it will be frozen at the corresponding time step, every time the pressure system is solved. Now, if first and second order finite volume schemes are considered and focusing on the equations for  $u$  and  $\pi$ , we could write (4.2.13) as follows:

$$\begin{cases} u_i'(t) = -\frac{1}{h_i(t_0)\Delta x} (\pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t)) - \frac{1}{h_i(t_0)} S_i(t), \\ \pi_i'(t) = -\frac{a^2}{h_i(t_0)\Delta x} (u_{i+1/2}^*(t) - u_{i-1/2}^*(t)), \end{cases} \quad (4.2.14)$$

where  $S_i(t)$  is now such that

$$S_i(t) \approx \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} g P_{i,h}(t, x) z'(x) dx,$$

with  $P_{i,h}$  a reconstruction operator within the cell  $i$  for the variable  $h(t, x)$ .

As  $h$  does not change in time in this step,  $S_i(t)$  is no longer time dependent, so we could consider

$$S_i \approx \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} g P_{i,h_0}(x) z'(x) dx,$$

where  $P_{i,h_0}(x)$  is a reconstruction of  $h|_{t=t_0}(x)$ , to be determined.

In practice, (4.2.14) can be rewritten in terms of the Riemann invariants  $\vec{w} = \pi + au$  and  $\overleftarrow{w} = \pi - au$  as follows:

$$\begin{cases} \vec{w}'_i(t) = -\frac{a}{h_i(t_0)\Delta x} (\vec{w}_{i+1/2}(t) - \vec{w}_{i-1/2}(t)) - \frac{a}{h_i(t_0)} S_i, \\ \overleftarrow{w}'_i(t) = \frac{a}{h_i(t_0)\Delta x} (\overleftarrow{w}_{i+1/2}(t) - \overleftarrow{w}_{i-1/2}(t)) + \frac{a}{h_i(t_0)} S_i. \end{cases} \quad (4.2.15)$$

Here, the approximations at the intercells  $\vec{w}_{i+1/2}(t) \approx \vec{w}(x_{i+1/2}, t)$  and  $\overleftarrow{w}_{i+1/2}(t) \approx \overleftarrow{w}(x_{i+1/2}, t)$  will be computed using a reconstruction operator. The advantage of using the Riemann invariants is that we manage to decouple the two equations, obtaining thus two transport equations with source terms.

Let us recall that  $\pi$  and  $u$  can be easily recovered from the values of  $\vec{w}$  and  $\overleftarrow{w}$  using (2.1.3).

Therefore, once we have solved (4.2.15), we can define  $\pi_{i+1/2}^*(t)$  and  $u_{i+1/2}^*(t)$  as

$$\pi_{i+1/2}^*(t) = \frac{P_{i,\vec{w}}(x_{i+1/2}, t) + P_{i+1,\overleftarrow{w}}(x_{i+1/2}, t)}{2}, \quad (4.2.16)$$

$$u_{i+1/2}^*(t) = \frac{P_{i,\vec{w}}(x_{i+1/2}, t) - P_{i+1,\overleftarrow{w}}(x_{i+1/2}, t)}{2a}, \quad (4.2.17)$$

where  $P_{i,\vec{w}}$  and  $P_{i,\overleftarrow{w}}$  correspond to some reconstruction operators within the cell. Using (4.2.16)-(4.2.17), we can compute the values  $\pi_{i+1/2}^*(t)$  and  $u_{i+1/2}^*(t)$  in (4.2.7) and (4.2.8).

Once the pressure system is approximated, we will use an upwind scheme in order to solve the transport ODE system defined in (4.2.8). That is, the values  $U_{i+1/2}^*(t)$  are defined as

$$U_{i+1/2}^*(t) = \begin{cases} P_{U,i}(x_{i+1/2}, t) & \text{if } u_{i+1/2}^*(t) \geq 0, \\ P_{U,i+1}(x_{i+1/2}, t) & \text{if } u_{i+1/2}^*(t) < 0, \end{cases} \quad (4.2.18)$$

where  $P_{U,i}$  denotes the reconstruction operator corresponding to  $U = (h, hu)$ .

Although up to now the time variable is kept continuous, the time steps will be solved afterward by means of an explicit or implicit scheme. In practice, the transport system

will always be solved explicitly. Remark that, for the sake of simplicity, we have not dealt yet with the well-balancing issue.

In order to achieve the exactly well-balanced character of the scheme we will follow again the ideas described in [25], where the main ingredients are: a fully exactly well-balanced reconstruction operator, a quadrature formula and a proper approximation of the source term  $S_i$ , that guarantees the exactly well-balanced character of the numerical scheme. We will describe our choices for the first and second-order exactly well-balanced reconstruction operators in Section 4.2.1. As we are interested in first and second order numerical schemes, we will use the mid-point rule as quadrature formula. Finally, the approximation of the source term  $S_i$  is also done following the ideas described in [25].

More precisely, given a time  $t = t_0$ , at every cell we compute the steady state  $(h_i^{e,t_0}, u_i^{e,t_0})$  that satisfies (1.1.17) such that  $h_i^{e,t_0}(x_i) = h_i^{t_0}$  and  $u_i^{e,t_0}(x_i) = u_i^{t_0}$ , or equivalently, the solution of (1.1.18) with  $C_{1,i}^{t_0} = (hu)_i^{t_0}$  and

$$C_{2,i}^{t_0} = \frac{(u_i^{t_0})^2}{2} + g(gh_i^{t_0} + z(x_i)).$$

Now, integrating (1.1.17) over the cell  $[x_{i-1/2}, x_{i+1/2}]$  we have that

$$\begin{aligned} & \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} gh_i^{e,t_0}(x)z'(x) dx \\ &= \frac{1}{\Delta x} (\pi_i^{e,t_0}(x_{i+1/2}) - \pi_i^{e,t_0}(x_{i-1/2}) + (hu)_i^{e,t_0} (u_i^{e,t_0}(x_{i+1/2}) - u_i^{e,t_0}(x_{i-1/2}))), \end{aligned}$$

where  $\pi_i^{e,t_0}(x) = \frac{g}{2}(h_i^{e,t_0})^2(x)$ . Taking into account the splitting procedure that we consider here, we could rewrite systems (4.2.13) and (4.2.8) as

$$\left\{ \begin{aligned} h_i'(t) &= 0, \\ (hu)_i'(t) &= -\frac{1}{\Delta x} (\pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t) - \pi_i^{e,t_0}(x_{i+1/2}) + \pi_i^{e,t_0}(x_{i-1/2})) \\ &\quad - \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} g(P_{i,h_0}(x) - h_i^{e,t_0}(x))z'(x)dx, \\ (h\pi)_i'(t) &= -\frac{a^2}{\Delta x} (u_{i+1/2}^*(t) - u_{i-1/2}^*(t)), \end{aligned} \right. \quad (4.2.19)$$

and

$$\left\{ \begin{aligned} h_i'(t) &= -\frac{1}{\Delta x} (h_{i+1/2}^*(t)u_{i+1/2}^*(t) - h_{i-1/2}^*(t)u_{i-1/2}^*(t)), \\ (hu)_i'(t) &= -\frac{1}{\Delta x} ((hu)_{i+1/2}^*(t)u_{i+1/2}^*(t) - (hu)_{i-1/2}^*(t)u_{i-1/2}^*(t)) \\ &\quad - \frac{1}{\Delta x} (-(hu)_i^{e,t_0}(u_i^{e,t_0}(x_{i+1/2}) - u_i^{e,t_0}(x_{i-1/2}))). \end{aligned} \right. \quad (4.2.20)$$

However, given the previous semi-discrete systems, if we considered a steady state initial condition, the third equation in (4.2.19) would not guarantee obtaining  $(h\pi)'_i(t) = 0$ . The way to deal with this issue is to consider a modified relaxed system, that could be seen as the relaxed system of the fluctuations:

$$\begin{cases} \partial_t h = 0, \\ \partial_t(hu) + \partial_x(\pi - \pi^e) = -g(h - h^e) z', \\ \partial_t(h\pi) + a^2 \partial_x(u - u^e) = 0. \end{cases} \quad (4.2.21)$$

Now, applying the splitting we obtain

$$\begin{cases} h'_i(t) = 0, \\ (hu)'_i(t) = -\frac{1}{\Delta x} \left( \pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t) - \pi_i^{e,t_0}(x_{i+1/2}) + \pi_i^{e,t_0}(x_{i-1/2}) \right) \\ - \frac{1}{\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} g(P_{i,h_0}(x) - h_i^{e,t_0}(x)) z'(x) dx, \\ (h\pi)'_i(t) = -\frac{a^2}{\Delta x} \left( u_{i+1/2}^*(t) - u_{i-1/2}^*(t) - u_i^{e,t_0}(x_{i+1/2}) + u_i^{e,t_0}(x_{i-1/2}) \right), \end{cases} \quad (4.2.22)$$

and (4.2.20).

System (4.2.22) can be written analogously in terms of the Riemann invariants:

$$\begin{cases} \vec{w}'_i(t) = -\frac{a}{h_i(t_0)\Delta x} \left( \vec{w}_{i+1/2}(t) - \vec{w}_{i-1/2}(t) - \vec{w}_i^{e,t_0}(x_{i+1/2}) + \vec{w}_i^{e,t_0}(x_{i-1/2}) \right) \\ - \frac{a}{h_i(t_0)\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} g(P_{i,h_0}(x) - h_i^{e,t_0}(x)) z'(x) dx, \\ \overleftarrow{w}'_i(t) = \frac{a}{h_i(t_0)\Delta x} \left( \overleftarrow{w}_{i+1/2}(t) - \overleftarrow{w}_{i-1/2}(t) - \overleftarrow{w}_i^{e,t_0}(x_{i+1/2}) + \overleftarrow{w}_i^{e,t_0}(x_{i-1/2}) \right) \\ + \frac{a}{h_i(t_0)\Delta x} \int_{x_{i-1/2}}^{x_{i+1/2}} g(P_{i,h_0}(x) - h_i^{e,t_0}(x)) z'(x) dx. \end{cases} \quad (4.2.23)$$

However, since we are considering schemes up to second order, we may apply the mid-point rule and considering that the averages correspond to the values at the center of the cells. Therefore, the source terms in (4.2.22) and (4.2.23) vanish, giving systems

$$\begin{cases} h'_i(t) = 0, \\ (hu)'_i(t) = -\frac{1}{\Delta x} \left( \pi_{i+1/2}^*(t) - \pi_{i-1/2}^*(t) - \pi_i^{e,t_0}(x_{i+1/2}) + \pi_i^{e,t_0}(x_{i-1/2}) \right) \\ (h\pi)'_i(t) = -\frac{a^2}{\Delta x} \left( u_{i+1/2}^*(t) - u_{i-1/2}^*(t) - u_i^{e,t_0}(x_{i+1/2}) + u_i^{e,t_0}(x_{i-1/2}) \right), \end{cases} \quad (4.2.24)$$

and

$$\begin{cases} \vec{w}'_i(t) = -\frac{a}{h_i(t_0)\Delta x} \left( \vec{w}_{i+1/2}(t) - \vec{w}_{i-1/2}(t) - \vec{w}_i^{e,t_0}(x_{i+1/2}) + \vec{w}_i^{e,t_0}(x_{i-1/2}) \right) \\ \overleftarrow{w}'_i(t) = \frac{a}{h_i(t_0)\Delta x} \left( \overleftarrow{w}_{i+1/2}(t) - \overleftarrow{w}_{i-1/2}(t) - \overleftarrow{w}_i^{e,t_0}(x_{i+1/2}) + \overleftarrow{w}_i^{e,t_0}(x_{i-1/2}) \right). \end{cases} \quad (4.2.25)$$

### 4.2.1 Well-balanced variable reconstructions

Let us now focus on the reconstruction of our variables. To do so, we need to keep in mind the well-balanced property and to adapt the general strategy presented in [25], combined with the ideas introduced in Chapter 2.

Note that, in practice a quadrature formula will be used. As only first and second order reconstruction operators are considered here, the integrals are approximated by the mid-point rule.

We will now show how the reconstruction of variables is done, given in a general form for a variable  $X$  that can be either  $\vec{w}$ ,  $\overleftarrow{w}$ ,  $h$  or  $q$ .

**First order reconstruction.** The first order reconstruction can be written as

$$P_{i,X}^{o1}(x,t) = X_i^{e,t_0}(x) + X_i(t) - X_i^{e,t_0}(x_i), \quad (4.2.26)$$

with  $Q_{i,X}(x)$  being

$$Q_{i,X}(x) = X_i(t) - X_i^{e,t_0}(x_i).$$

**Second order reconstruction.** For the second order schemes we consider the following reconstruction:

$$P_{i,X}^{o2}(x,t) = X_i^{e,t_0}(x) + X_i(t) - X_i^{e,t_0}(x_i) + \Delta X_i^{t,f}(x - x_i) + \Delta X_i^{t_0,f}(x - x_i), \quad (4.2.27)$$

where

$$\Delta X_i^{t,f} = \frac{1}{\Delta x} \left( \tilde{\phi}_{i+}^{t_0}(X_i^{t,f} - X_{i-1}^{t,f}) + \tilde{\phi}_{i-}^{t_0}(X_{i+1}^{t,f} - X_i^{t,f}) \right)$$

with  $X_i^{t,f} = X_i(t) - X_i^{t_0}$  and

$$\tilde{\phi}_{i-}^{t_0} = \begin{cases} \frac{|d_{i-}|}{|d_{i-}| + |d_{i+}|} & \text{if } |d_{i-}| + |d_{i+}| > 0, \\ 0 & \text{otherwise,} \end{cases}$$

and

$$\tilde{\phi}_{i+}^{t_0} = \begin{cases} \frac{|d_{i+}|}{|d_{i-}| + |d_{i+}|} & \text{if } |d_{i-}| + |d_{i+}| > 0, \\ 0 & \text{otherwise,} \end{cases}$$

where  $d_{i-} = X_i^{t,f} - X_{i-1}^{t,f}$  and  $d_{i+} = X_{i+1}^{t,f} - X_i^{t,f}$ , and

$$\Delta X_i^{t_0,f} = \frac{1}{\Delta x} \left( \phi_{i+}^{t_0}(X_i^{t_0,f} - X_{i-1}^{t_0,f}) + \phi_{i-}^{t_0}(X_{i+1}^{t_0,f} - X_i^{t_0,f}) \right),$$

being  $\phi_{i\pm}^{t_0} = \tilde{\phi}_{i\pm}^{t_0}$  and

$$X_j^{t_0,f} = X_j^{t_0} - X_i^{e,t_0}(x_j)$$

for a given cell  $i$ .

Note that now  $Q_{i,X}(x)$  is given by

$$Q_{i,X}(x) = X_i(t) - X_i^{e,t_0}(x_i) + \Delta X_i^{t,f}(x - x_i) + \Delta X_i^{t_0,f}(x - x_i).$$

**Theorem 4.2.1.** *The schemes that result after considering the semi-discrete schemes (4.2.24), (4.2.25), (4.2.20) and the previous reconstructions (4.2.26) and (4.2.27) are fully well-balanced.*

*Proof.* Let us suppose that the initial condition is stationary. Then, in the first order case

$$\begin{aligned} \vec{w}_{i+1/2}(t) &= P_{i,\vec{w}}(x_{i+1/2}, t) = \vec{w}_i(t) - \pi_i^{e,t_0}(x_i) - au_i^{e,t_0}(x_i) + \pi_i^{e,t_0}(x_{i+1/2}) + au_i^{e,t_0}(x_{i+1/2}) \\ &= \vec{w}_i^{e,t_0}(x_i) - \pi_i^{e,t_0}(x_i) - au_i^{e,t_0}(x_i) + \pi_i^{e,t_0}(x_{i+1/2}) + au_i^{e,t_0}(x_{i+1/2}) \\ &= \pi_i^{e,t_0}(x_{i+1/2}) + au_i^{e,t_0}(x_{i+1/2}) = \vec{w}_i^{e,t_0}(x_{i+1/2}), \end{aligned}$$

where in the second line we are considering that the averages correspond to the values at the center of the cells. The same result holds for the other variables as well as for the second order case, since  $\Delta \vec{w}_i^{t_0,f} = 0$  and  $\Delta \vec{w}_i^{t,f} = 0$ . The same would happen with the reconstruction of  $\overleftarrow{w}$ ,  $h$  and  $q$ .

Then, from (4.2.25) we obtain

$$\begin{cases} \vec{w}'_i(t) &= 0, \\ \overleftarrow{w}'_i(t) &= 0, \end{cases}$$

and therefore, the pressure and transport semi-discrete systems (4.2.24) and (4.2.20) are trivial, being both

$$\begin{cases} h'_i(t) &= 0, \\ (hu)'_i(t) &= 0. \end{cases}$$

Therefore, the stationary solution is preserved.  $\square$

### 4.3 First order scheme

In this section we shall describe an exactly fully well-balanced first order scheme. Two approaches will be considered: an explicit version, where both (4.2.3) and (4.2.4) are solved explicitly, and a semi-implicit approach, where (4.2.3) is solved implicitly. As said previously, the semi-implicit scheme will allow us to have less restrictive CFL condition in subcritical regimes where velocity terms are smaller than the pressure terms.

The time-stepping will be done as follows: given a set of cell averages at time  $t^n$ ,  $U_i^n$ , we solve first system (4.2.3) in the time interval  $[t^n, t^{n+1}]$  obtaining the cell averages at time  $t^{n+1}$  and denoted by superindex  $n+1-$ . Then, starting for these cell averages, we solve system (4.2.4) in the time interval  $[t^n, t^{n+1}]$ , obtaining the approximation of the solution at the next time step,  $t^{n+1}$ , denoted by superindex  $n+1$ .

The stationary solutions will be computed as discussed in Section 1.1.2 (see equation (1.1.20)), by computing the constants  $C_1$  and  $C_2$  with the values at the center of the cells, which in this case correspond to the averages. In the case in which two solutions of the equation exist, we keep the subcritical one or the supercritical one so that it matches the character of the cell.

### 4.3.1 Explicit scheme

In view of the semi-discrete scheme (4.2.24), we propose the following first order explicit scheme for the pressure system:

$$\begin{aligned} h_i^{n+1-} &= h_i^n, \\ (hu)_i^{n+1-} &= (hu)_i^n - \frac{\Delta t}{\Delta x} \left( \pi_{i+1/2}^{*,n} - \pi_{i-1/2}^{*,n} - \pi_i^{e,n}(x_{i+1/2}) + \pi_i^{e,n}(x_{i-1/2}) \right). \end{aligned} \quad (4.3.1)$$

where the values  $\pi_{i+1/2}^{*,n}$  are computed by means of a fully well-balanced reconstruction operator for  $\vec{w}$  and  $\vec{w}$ , given by (4.2.26). That is:

$$\pi_{i+1/2}^{*,n} = \frac{P_{i,\vec{w}}^{o1}(x_{i+1/2}, t^n) + P_{i+1,\vec{w}}^{o1}(x_{i+1/2}, t^n)}{2},$$

Then, the transport system (4.2.20) is solved using  $(h_i^{n+1-}, (hu)_i^{n+1-})$  as initial condition.

$$\begin{aligned} h_i^{n+1} &= h_i^{n+1-} - \frac{\Delta t}{\Delta x} \left( h_{i+1/2}^{*,n+1-} u_{i+1/2}^{*,n+1-} - h_{i-1/2}^{*,n+1-} u_{i-1/2}^{*,n+1-} \right), \\ (hu)_i^{n+1} &= (hu)_i^{n+1-} - \frac{\Delta t}{\Delta x} \left( (hu)_{i+1/2}^{*,n+1-} u_{i+1/2}^{*,n+1-} - (hu)_{i-1/2}^{*,n+1-} u_{i-1/2}^{*,n+1-} \right) \\ &\quad + \frac{\Delta t}{\Delta x} \left( (hu)_i^{n+1-} \left( u_{i,i+1/2}^{e,t_{n+1-}} - u_{i,i-1/2}^{e,t_{n+1-}} \right) \right), \end{aligned} \quad (4.3.2)$$

where  $u_{i,i\pm 1/2}^{e,t_{n+1-}} = u_i^{e,t_{n+1-}}(x_{i\pm 1/2})$ . Note that now, the values  $u_{i+1/2}^{*,n+1-}$  and  $h_{i\pm 1/2}^{*,n+1-}$  and  $(hu)_{i\pm 1/2}^{*,n+1-}$  must be determined. The values  $u_{i+1/2}^{*,n+1-}$  are computed at each intercell by means of a fully well-balanced first order reconstruction operator as follows:

$$u_{i+1/2}^{*,n+1-} = \frac{P_{i,\vec{w}}^{o1}(x_{i+1/2}, t^{n+1-}) - P_{i+1,\vec{w}}^{o1}(x_{i+1/2}, t^{n+1-})}{2a}.$$

Finally the values  $h_{i\pm 1/2}^{*,n+1-}$  and  $(hu)_{i\pm 1/2}^{*,n+1-}$  are also computed using again the fully well-balanced first order reconstruction operator and the upwind scheme, that is

$$X_{i+1/2}^{*,n+1-} = \begin{cases} P_{i,X}^{ol}(x_{i+1/2}) & \text{if } u_{i+1/2}^{*,n+1-} \geq 0, \\ P_{i+1,X}^{ol}(x_{i+1/2}) & \text{if } u_{i+1/2}^{*,n+1-} < 0, \end{cases}$$

where  $X = h, hu$ .

Remark that here we have used a splitting technique by solving first the system defined by  $S_P$  and then the one defined by  $S_T$ . Nevertheless, nothing obliges to do the splitting in that order and one could consider a variant of this explicit scheme by solving first (4.2.4) and then (4.2.3).

### 4.3.2 Semi-implicit scheme

As previously said, in subcritical regimes where  $u^2 \ll gh$ , the main restriction of the CFL condition comes from the pressure terms. Therefore, in view of the eigenvalues of (4.2.3) (see (4.2.5)) we consider an implicit version of the pressure system.

$$\begin{aligned} h_i^{n+1-} &= h_i^n, \\ (hu)_i^{n+1-} &= (hu)_i^n - \frac{\Delta t}{\Delta x} \left( \pi_{i+1/2}^{*,n+1-} - \pi_{i-1/2}^{*,n+1-} - \pi_i^{e,n}(x_{i+1/2}) + \pi_i^{e,n}(x_{i-1/2}) \right). \end{aligned} \quad (4.3.3)$$

It can be seen that  $(hu)_i^{n+1-}$  could be also obtained as

$$(hu)_i^{n+1-} = h_i^n u_i^{n+1-} = h_i^n \cdot \frac{\vec{w}_i^{n+1-} - \overleftarrow{w}_i^{n+1-}}{2a},$$

where  $\vec{w}_i^{n+1-}$  and  $\overleftarrow{w}_i^{n+1-}$  are given by

$$\vec{w}_i^{n+1-} = \vec{w}_i^n - \frac{a\Delta t}{h_i^n \Delta x} \left( \vec{w}_{i+1/2}^{n+1} - \vec{w}_{i-1/2}^{n+1} - \vec{w}_{i+1/2}^{e,n} + \vec{w}_{i-1/2}^{e,n} \right),$$

and

$$\overleftarrow{w}_i^{n+1-} = \overleftarrow{w}_i^n + \frac{a\Delta t}{h_i^n \Delta x} \left( \overleftarrow{w}_{i+1/2}^{n+1} - \overleftarrow{w}_{i-1/2}^{n+1} - \overleftarrow{w}_{i+1/2}^{e,n} + \overleftarrow{w}_{i-1/2}^{e,n} \right),$$

where  $\vec{w}_{i+1/2}^{n+1}$  is given by

$$\vec{w}_{i+1/2}^{n+1} = P_{i,\vec{w}}^{ol}(x_{i+1/2}, t^{n+1}) = \vec{w}_i^{n+1} - \pi_i^{e,n}(x_i) - au_i^{e,n}(x_i) + \pi_i^{e,n}(x_{i+1/2}) + au_i^{e,n}(x_{i+1/2}),$$

$\overleftarrow{w}_{i+1/2}^{n+1}$  is defined similarly, and

$$\vec{w}_{i\pm 1/2}^{e,n} = \frac{1}{2}g(h_i^e)^2(x_{i\pm 1/2}) + au_i^e(x_{i\pm 1/2}).$$

Next, the transport step is computed as in the explicit case.

As we have said for the explicit case, we may reverse the order of the splitting so that we may first solve system (4.2.4) and then system (4.2.3). Of course, system (4.2.4) would be solved explicitly and (4.2.3), implicitly.

## 4.4 Second order scheme

In order to obtain second order accuracy, the time-stepping will be done using a Strang splitting method (see [85, 68, 69]), described in Section 1.1.3.3. More explicitly:

1. Perform a step of the first system with time step  $\Delta t/2$ , obtaining an approximation  $\tilde{h}_i^{n+1}$  and  $(\widetilde{hu})_i^{n+1}$ .
2. Perform a step of the second system with time step  $\Delta t$ , obtaining the approximation denoted by  $\widehat{h}_i^{n+1}$  and  $(\widehat{hu})_i^{n+1}$ .
3. Perform a final step of the first system with time step  $\Delta t/2$ , obtaining the approximations  $h_i^{n+1}$  and  $(hu)_i^{n+1}$  at time  $t^{n+1}$ .

Let us remark that there is no a priori restriction on which of the systems (4.2.3) or (4.2.4) should go first.

This may be summarize in a compact form as follows:

Denote by  $S_P^\tau$ ,  $S_T^\tau$  the approximate solution operators in the interval  $[t, t + \tau]$  of the corresponding exact solution operators to the pressure system  $S_P$  and transport system  $S_T$  respectively. Then, the first version of the scheme corresponds to

$$U(x, t + \Delta t) = S_P^{\frac{\Delta t}{2}} \circ S_T^{\Delta t} \circ S_P^{\frac{\Delta t}{2}}(U(x, t)), \quad (4.4.1)$$

while the second version corresponds to

$$U(x, t + \Delta t) = S_T^{\frac{\Delta t}{2}} \circ S_P^{\Delta t} \circ S_T^{\frac{\Delta t}{2}}(U(x, t)), \quad (4.4.2)$$

Remark that in each of the steps we need to consider second order approximations in space, while the time stepping is just first order within the step, the second order in time being obtained thanks to Strang method.

In this second order case, the stationary solutions are again computed by applying the discussion given in Section 1.1.2 and we keep the subcritical or the supercritical one depending on the character of the cell in which we are computing the value, in the case in which there are two solutions of the cubic equation (1.1.20).

### 4.4.1 Explicit scheme

We shall describe the case corresponding to (4.4.1). The second version given by (4.4.2) is analogous.

In this explicit case, the solution of the first step is obtained by applying (4.3.1) with time step  $\Delta t/2$ . Afterwards, we solve the transport system using (4.3.2), and finally the pressure system is solved again applying (4.3.1) with time step  $\Delta t/2$ . Of course, in the previous schemes, second order approximations in space are considered.

### 4.4.2 Semi-implicit scheme

As before, we will describe the case corresponding to (4.4.1). In this case, similarly as done for the first order scheme, the steps corresponding to the operator  $S_P$  are performed implicitly. Therefore, we use the same procedure as in the second order explicit case but now using for the pressure system the implicit scheme (4.3.3) instead of (4.3.1).

Remark that the semi-implicit second order scheme requires to solve linear systems for the pressure step. Therefore, accounting for the computational cost, in this case it is especially interesting to consider the second version of the scheme, where only one step corresponds to the pressure system.

## 4.5 Numerical experiments

In this section, we consider a wide range of numerical experiments in order to test the performance of the different schemes proposed here. We will denote by EXP the results obtained by the fully explicit schemes and by IMP the semi-implicit ones. The accuracy will be indicated as O1 or O2 for the first or second order respectively. Moreover, the different versions of the schemes will be denoted by PT, TP, PTP and TPT, which indicate the order in which the pressure (P) and the transport (T) system have been solved.

### 4.5.1 Fully well-balanced property

In order to check that the fully well-balanced property is satisfied, we consider the spatial domain  $[-5, 5]$  and define as bottom topography a gaussian bump given by

$$z(x) = 0.5 \exp(-x^2). \quad (4.5.1)$$

Then, a subcritical steady state is computed by setting  $C_1 = hu = 0.1$  and constant energy level

$$C_2 = \frac{(0.1)^2}{2} + g(1 + z(-5)),$$

which corresponds to the value obtained by imposing  $h = 1$  at the left boundary. This subcritical steady state is considered as initial condition for the schemes. Tables 4.1, 4.2,

4.3 and 4.4 show the difference in  $L^1$  norm between the initial condition and the solution obtained with the schemes at time  $t = 1$  with  $N = 100$  cells in the domain and setting the CFL value to 1 for the explicit schemes and 5 for the implicit ones. As expected, the steady state is preserved, obtaining errors of order  $10^{-14}$ .

EXP O1 PT		EXP O1 TP	
$h$	$hu$	$h$	$hu$
$5.83 \cdot 10^{-14}$	$5.80 \cdot 10^{-14}$	$4.54 \cdot 10^{-14}$	$6.57 \cdot 10^{-14}$

Table 4.1: Difference in  $L^1$  norm between the initial condition and the solution obtained at time  $t = 1$  with each of the explicit first order schemes

IMP O1 PT		IMP O1 TP	
$h$	$hu$	$h$	$hu$
$8.14 \cdot 10^{-14}$	$9.97 \cdot 10^{-14}$	$6.82 \cdot 10^{-14}$	$7.58 \cdot 10^{-14}$

Table 4.2: Difference in  $L^1$  norm between the initial condition and the solution obtained at time  $t = 1$  with each of the implicit first order schemes

EXP O2 PTP		EXP O2 TPT	
$h$	$hu$	$h$	$hu$
$4.01 \cdot 10^{-14}$	$6.87 \cdot 10^{-14}$	$3.44 \cdot 10^{-14}$	$5.74 \cdot 10^{-14}$

Table 4.3: Difference in  $L^1$  norm between the initial condition and the solution obtained at time  $t = 1$  with each of the explicit second order schemes

IMP O2 PTP		IMP O2 TPT	
$h$	$hu$	$h$	$hu$
$5.41 \cdot 10^{-14}$	$8.41 \cdot 10^{-14}$	$4.80 \cdot 10^{-14}$	$8.13 \cdot 10^{-14}$

Table 4.4: Difference in  $L^1$  norm between the initial condition and the solution obtained at time  $t = 1$  with each of the implicit second order schemes

### 4.5.2 Accuracy test

Let us now check the order of the schemes by performing an accuracy test. To do so, we consider as initial condition a small perturbation of a water at rest steady state:

$$h(x, 0) = \begin{cases} 0.05 \left( 1 + \cos \left( \frac{2\pi(x-4750)}{3500} \right) \right) & \text{if } 3000 < x < 6500 \\ 0.05 \left( - \left( 1 + \cos \left( \frac{2\pi(x-9250)}{3500} \right) \right) \right) & \text{if } 7500 < x < 11000 \\ 0 & \text{otherwise} \end{cases}$$

with  $q(x, 0) = 0$  and bottom topography

$$z(x) = - \left( 50 - \exp \left( - \frac{(x - 7000)^2}{1000000} \right) \right).$$

The spatial domain corresponds to  $[0, 14000]$  and the final time is  $t = 0.5$ . Periodic boundary conditions are considered. The considered reference solution has 6400 cells.

The errors are shown in Tables 4.6, 4.7 and 4.8. The expected order is reached either for the explicit or semi-implicit version. We remark that, concerning the order of convergence, no major differences are observed either if we begin with the pressure or the transport system. Therefore, focusing exclusively in the order of accuracy, it would make sense to start with the transport system for the second order semi-implicit scheme, since the computational cost would be less.

No. of cells	EXP O1 PT (CFL 1)				EXP O1 TP (CFL 1)			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
100	2.82e+00		1.38e+00		2.82e+00		1.44e+00	
200	7.50e-01	1.91	4.22e-01	1.71	7.50e-01	1.91	4.55e-01	1.66
400	2.28e-01	1.72	1.39e-01	1.60	2.28e-01	1.72	1.56e-01	1.55
800	7.93e-02	1.53	4.91e-02	1.50	7.93e-02	1.53	5.73e-02	1.44
1600	2.91e-02	1.45	1.75e-02	1.49	2.91e-02	1.45	2.12e-02	1.43

Table 4.5: Errors in  $L^1$  norm and convergence rates for the first order explicit schemes

### 4.5.3 Perturbation of water at rest

We propose now to closely study the behavior of the different schemes. Let us consider the bottom topography given by (4.5.1) in the domain  $[-5, 5]$ . The following perturbation of water at rest is considered as initial condition

$$h(x) = -z(x) + 0.1e^{-x^2}, \quad u(x) = 0.$$

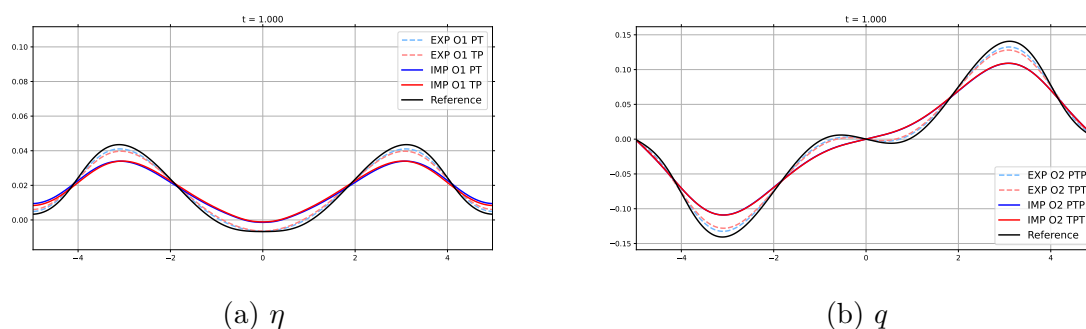
No. of cells	IMP O1 PT (CFL 5)				IMP O1 TP (CFL 5)			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
100	2.82e+00		1.51e+00		2.82e+00		1.57e+00	
200	7.51e-01	1.91	4.89e-01	1.62	7.51e-01	1.91	5.23e-01	1.59
400	2.29e-01	1.71	1.73e-01	1.50	2.29e-01	1.71	1.90e-01	1.46
800	7.98e-02	1.52	6.54e-02	1.40	7.98e-02	1.52	7.32e-02	1.37
1600	2.96e-02	1.43	2.47e-02	1.40	2.96e-02	1.43	2.81e-02	1.38

Table 4.6: Errors in  $L^1$  norm and convergence rates for the first order implicit schemes

No. of cells	EXP O2 PTP (CFL 1)				EXP O2 TPT (CFL 1)			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
100	3.05e+00		2.30e+00		3.05e+00		2.31e+00	
200	7.64e-01	2.00	5.43e-01	2.08	7.64e-01	2.00	5.43e-01	2.09
400	1.91e-01	2.00	1.26e-01	2.11	1.91e-01	2.00	1.25e-01	2.12
800	4.73e-02	2.01	3.09e-02	2.02	4.73e-02	2.01	3.06e-02	2.03
1600	1.13e-02	2.06	7.80e-03	1.99	1.13e-02	2.06	7.97e-03	1.94

Table 4.7: Errors in  $L^1$  norm and convergence rates for the second order explicit schemes

In Figures 4.1 and 4.2 we can see the solution obtained with the first and second order schemes at time  $t = 1$  using 200 cells and CFL 0.8 for the explicit schemes and 2 for the implicit ones. In both figures we have also plotted a reference solution that has been computed using the EXP O1 PT scheme with 1600 cells. Again, periodic boundary conditions have been considered.

Figure 4.1: Solution at time  $t=1$  obtained with the first order schemes using 200 cells

As expected, the implicit schemes are more diffusive than the explicit ones. However, this diffuseness is reduced when we consider second order schemes.

No. of cells	IMP O2 PTP (CFL 5)				IMP O2 TPT (CFL 5)			
	$h$		$q$		$h$		$q$	
	Error	Order	Error	Order	Error	Order	Error	Order
100	3.05e+00		2.29e+00		3.05e+00		2.28e+00	
200	7.64e-01	2.00	5.44e-01	2.07	7.64e-01	2.00	5.44e-01	2.07
400	1.91e-01	2.00	1.26e-01	2.11	1.91e-01	2.00	1.26e-01	2.11
800	4.73e-02	2.01	3.10e-02	2.03	4.73e-02	2.01	3.06e-02	2.04
1600	1.13e-02	2.06	7.44e-03	2.06	1.13e-02	2.07	7.26e-03	2.07

Table 4.8: Errors in  $L^1$  norm and convergence rates for the second order implicit schemes

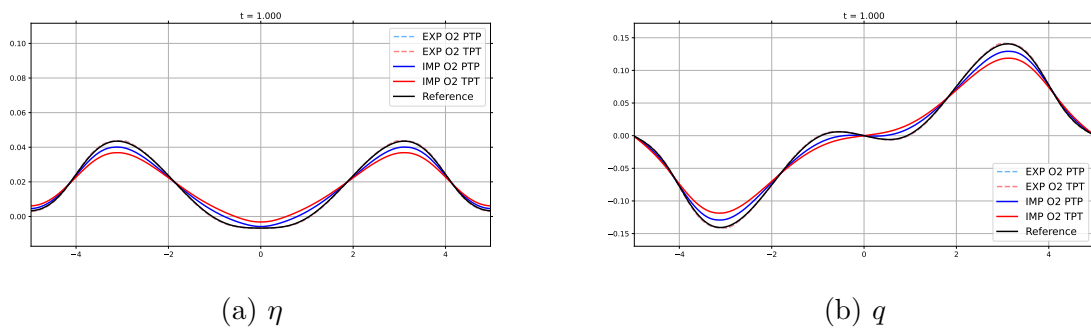


Figure 4.2: Solution at time  $t=1$  obtained with the second order schemes using 200 cells

Let us now consider a bigger CFL number for the implicit schemes, which is shown in Figure 4.3. We now remark a major difference whether we begin with the pressure or transport system. As shown on the left-hand side image in Figure 4.3, with  $CFL=5$ , the IMP O1 PT scheme performs better in terms of stability than the IMP O1 TP. Conversely, on the right-hand side for the second order case, the IMP O2 TPT scheme shows better performance than the IMP O2 PTP in terms of stability. Therefore, from now on, we will just consider the IMP O1 PT and IMP O2 TPT versions of the semi-implicit schemes.

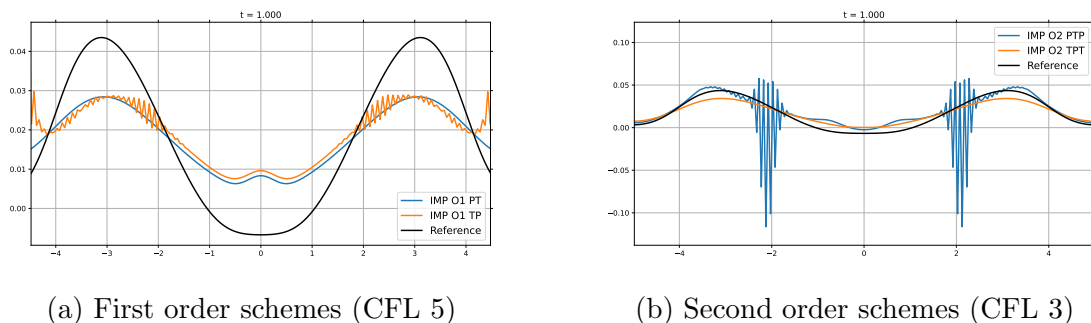


Figure 4.3: Solution for  $\eta$  at time  $t=1$  obtained using 200 cells when increasing the CFL

#### 4.5.4 Perturbation of water at rest with shock waves

Up to this point, the numerical tests considered corresponded to smooth solutions. We want now to check the performance of the schemes in the presence of shocks. In order to do so we consider the same bottom topography as in the previous test case, that is the topography given by (4.5.1), and define the following initial condition in  $[-5, 5]$ :

$$h(x) = \begin{cases} -z(x) & \text{if } |x| \geq 1 \\ -z(x) + 0.1 & \text{if } |x| < 1 \end{cases}, \quad u(x) = 0,$$

Periodic boundary conditions will be used.

Figures 4.4 and 4.5 show the solutions obtained at time  $t = 1$  for the first and second order schemes respectively. For the explicit schemes the CFL value has been set to 0.8 and for the implicit ones we have considered two cases: solution with CFL 2 and with CFL 3. In order to compare the results, we include the reference solution computed with the EXP O1 PT scheme and 1600 cells.

We observe that the schemes successfully handle the initial condition, not observing important spurious oscillations for the second order schemes. This might be thanks to the use of slope limiters in the reconstruction operators.

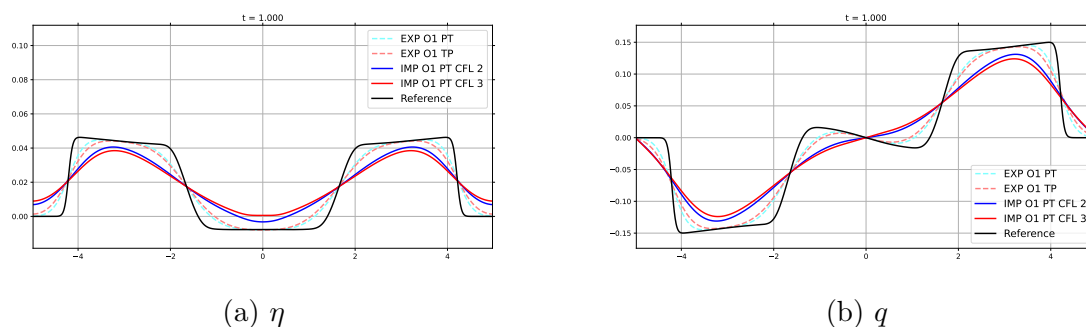


Figure 4.4: Solution at time  $t=1$  obtained with the first order schemes using 200 cells

#### 4.5.5 Perturbation of a subcritical solution

We will now perform a test proposed in [56], in which a perturbation of a smooth subcritical stationary solution is considered as initial condition. The initial condition will be given by  $U_0(x) = (h_0(x), q_0(x))^t$  for  $x \in [0, 3]$ , where

$$h_0(x) = \begin{cases} h^*(x) + 0.02 & \text{if } 0.7 \leq x \leq 1, \\ h^*(x) & \text{otherwise,} \end{cases}$$

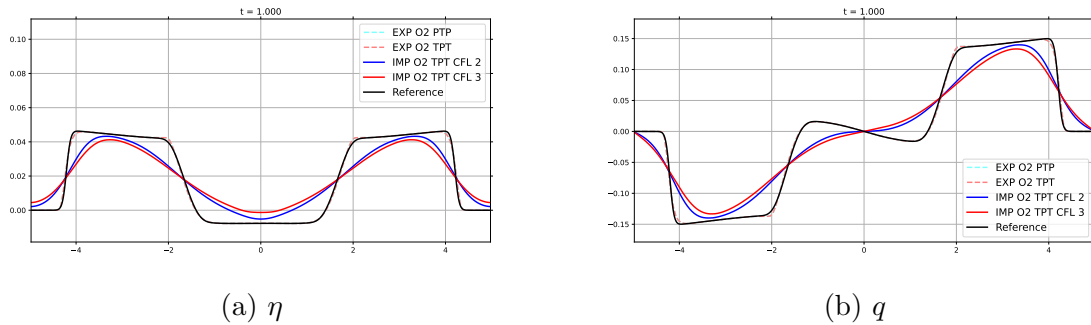


Figure 4.5: Solution at time  $t=1$  obtained with the second order schemes using 200 cells

and  $q_0(x) = q^*(x)$ , being  $U^*(x) = (h^*(x), q^*(x))^t$ , the solution of the following Cauchy problem (1.1.17) with initial condition  $h(0) = 2, q(0) = 3.5$ . Moreover, the bottom topography is given by

$$z(x) = \begin{cases} 0.25(1 + \cos(5\pi(x + 0.5))) & \text{if } 1.3 \leq x \leq 1.7, \\ 0 & \text{otherwise.} \end{cases} \quad (4.5.2)$$

This initial condition is plotted in Figure 4.6.

As boundary conditions, we impose the value of  $q$  on the left and the one of  $h$  on the right.

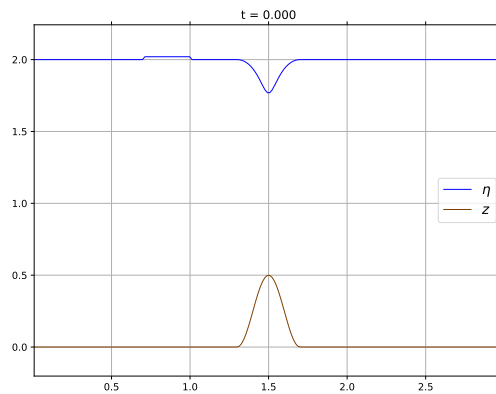


Figure 4.6: Perturbation of a subcritical solution initial condition

In Figures 4.7 we have plotted the difference between the result of the scheme and the steady state at time  $t = 0.1$  using  $N = 200$  cells for the first and second order schemes for the variable  $h$ . A reference solution has been computed using the first order explicit scheme with 1600 cells. For the implicit schemes, the CFL value has been set to 5.

We clearly observe the well-balanced character of the schemes, since they preserve the stationary solution in the areas where the perturbation has not arrived yet.

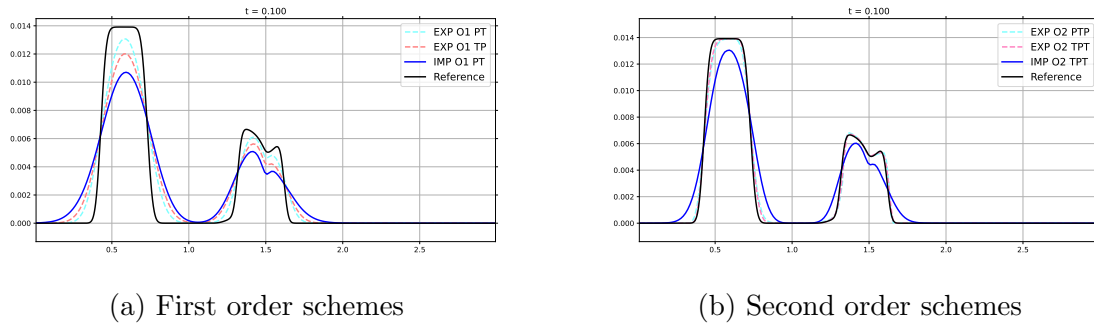


Figure 4.7: Difference between the result of the scheme and the steady state at time  $t=0.1$  using  $N = 200$  cells for the variable  $h$

Moreover, in Table 4.9, we show the errors for  $h$  and the CPU times at time  $t = 10$  for the different schemes using 100 cells. In the first order case, the semi-implicit scheme takes some more seconds than the explicit one but the error is also lower than the others. However, in the second order case we observe errors of the same magnitude and the CPU time needed by the semi-implicit scheme is approximately 25% lower than the explicit ones. Of course, if we increase the final time and the perturbation leave the domain, we capture the steady state, as shown in Table 4.10.

Scheme	Error for $h$	CPU time
EXP O1 PT	$1.31 \cdot 10^{-3}$	40.79
EXP O2 TP	$1.12 \cdot 10^{-3}$	41.29
IMP O1 PT	$7.75 \cdot 10^{-4}$	45.47
EXP O2 PTP	$1.62 \cdot 10^{-3}$	126.82
EXP O2 TPT	$1.64 \cdot 10^{-3}$	129.70
IMP O2 TPT	$1.22 \cdot 10^{-3}$	96.62

Table 4.9: Error in  $L^1$  norm and CPU time for the different schemes at time  $t = 10$  using 100 cells

#### 4.5.6 Perturbation of a transcritical smooth solution

For this test, we will once again take into account a test proposed in [24] in which a stationary solution with a transition at  $x_{crit} = 1.5$  which is the solution of (1.1.17) with constants  $C_1 = 2.5$  and  $C_2 = 17.56957396120237$ , and the same depth function as in previous tests, (4.5.2). A small perturbation of size  $\Delta h = 0.02$  is imposed in the interval

Scheme	Error for $h$
EXP O1 PT	$3.94 \cdot 10^{-13}$
EXP O2 TP	$3.73 \cdot 10^{-13}$
IMP O1 PT	$2.86 \cdot 10^{-13}$
EXP O2 PTP	$3.86 \cdot 10^{-13}$
EXP O2 TPT	$2.14 \cdot 10^{-13}$
IMP O2 TPT	$2.93 \cdot 10^{-13}$

Table 4.10: Error in  $L^1$  norm and CPU time for the different schemes at time  $t = 100$  using 100 cells

[1.1, 1.2]. This initial condition is plotted in Figure 4.8. As boundary conditions, we impose the value of  $q$  on the left and leave free boundary conditions on the right.

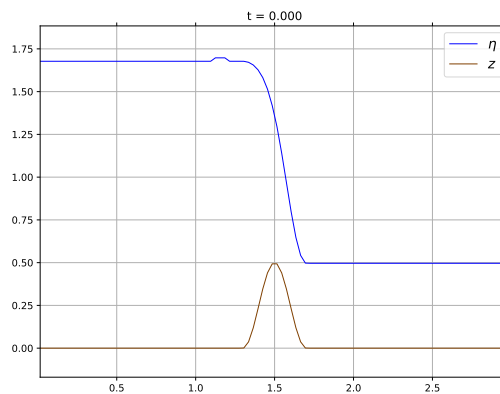


Figure 4.8: Perturbation of a transcritical smooth solution initial condition

In Figures 4.9 and 4.10, we have plotted the difference between the result of the different schemes and the steady state at time  $t = 0.15$  using 200 cells and CFL 5 for the implicit schemes. Again, the reference solution has been computed by using the first order explicit scheme with 1600 cells. It might look like in the implicit case the right wave is shifted to the left with respect to the reference solution, but this is due to diffusion. To be sure about this, for the first order case we have also computed a reference solution with the implicit scheme by increasing the number of cells to 1600 and setting the CFL value to be 1, observing that the solutions converge to reference solution computed with the explicit scheme. Moreover, the peak of the sonic point is not as pronounced as it appears to be. In order to show this we have plotted the free surface for the different schemes in Figure 4.11, where no big differences are observed between the different schemes.

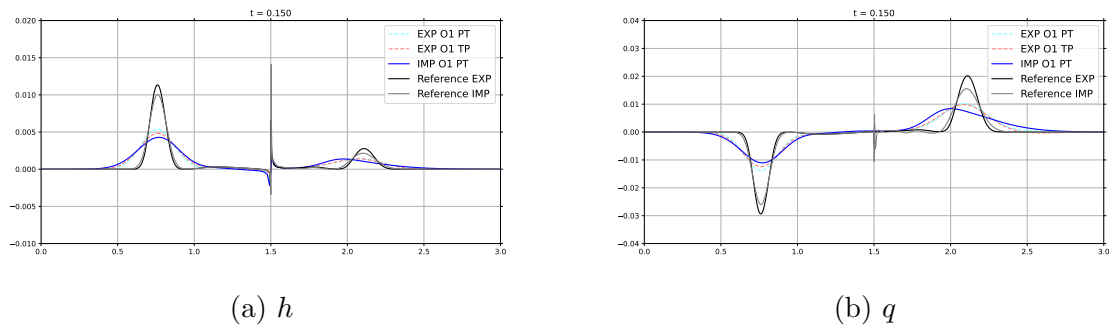


Figure 4.9: Difference between the result of the first order schemes and the steady state at time  $t = 0.15$  using 200 cells

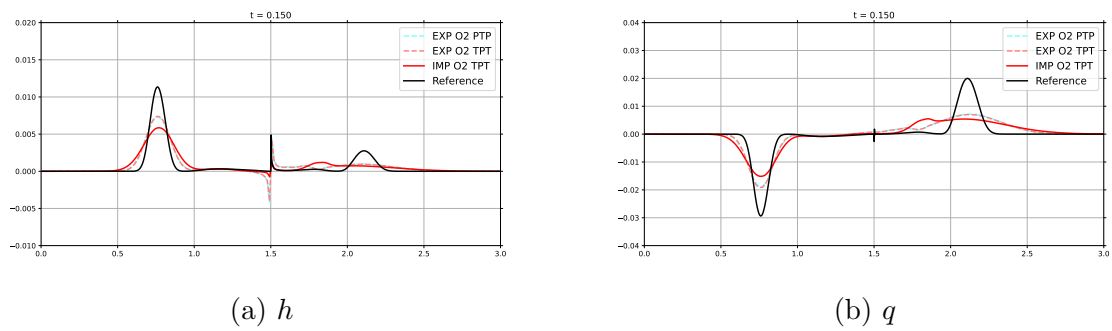


Figure 4.10: Difference between the result of the second order schemes and the steady state at time  $t = 0.15$  using 200 cells

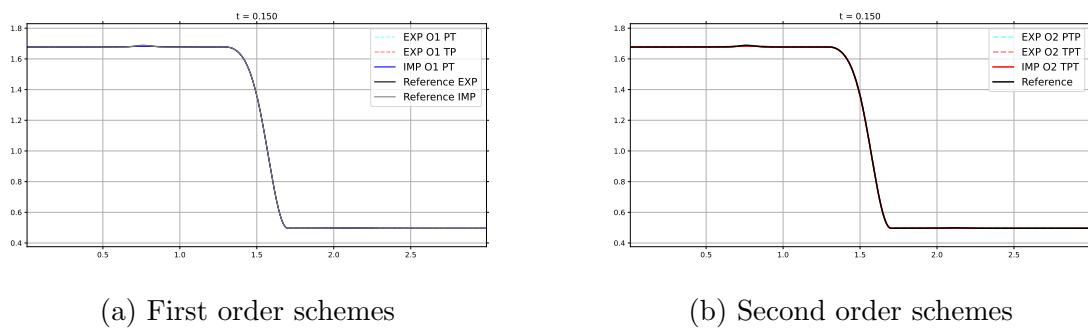


Figure 4.11: Solution for the free surface at time  $t = 0.15$  obtained using 200 cells

# Chapter 5

## Conclusions and future work

The aim of this thesis is the design of well-balanced schemes for the shallow water model and for the Ripa one that allow to decouple the acoustic and the transport waves. This is achieved for one-dimensional problems. The fact that we manage to decouple these waves enables us to consider semi-implicit schemes which allow a bigger time step when compared to explicit schemes, especially in the case of low Froude number, making these schemes more efficient.

In Chapter 2 we propose first and second order semi-implicit schemes that exactly preserve water at rest steady states for the SWE by applying the Lagrange-Projection approach.

By applying this Lagrange-Projection strategy we also manage to obtain a first order semi-implicit scheme for the Ripa system that is well-balanced for the hydrostatic steady states, which is presented in Chapter 3. In this case, we face the additional difficulty that the steady states are not explicitly determined as in the SWE case, so we choose to use a collocation method to define a discrete approximation of them. Therefore, the resulting scheme is well-balanced, but not exactly well-balanced.

The Lagrange-Projection approach has already been applied in multiple works and it has been proven to be an interesting strategy to consider when we want to design semi-implicit schemes. However, the use of Lagrangian coordinates can be challenging in some cases. For example, if we want to design schemes that preserve moving equilibria, since a steady-state in Eulerian coordinates is not necessarily a steady-state in Lagrangian coordinates. This was the problem we faced when we wanted to design fully exactly well-balanced schemes for the shallow water equations. For this reason, another strategy was chosen to carry out this task, consisting on applying an appropriate splitting of the system followed by a relaxation technique. This strategy is presented in Chapter 4.

To summarise, we have successfully managed to design well-balanced semi-implicit schemes for different shallow flows that allow to separate low and fast waves, being especially interesting in the case of subsonic regimes when compared to explicit schemes.

As far as future endeavours are concerned, we aim to extend the schemes presented in this thesis to dimension 2. However, applying the Lagrange-Projection strategy presented in Chapters 2 and 3 for this expansion could pose some challenges, particularly in the projection step. The main reason is that projecting the Lagrangian coordinates into Eulerian coordinates within a two-dimensional mesh introduces complexities that may make the procedure cumbersome.

Considering these facts, an alternative approach emerges from the strategy presented in Chapter 4. We are optimistic that the application of the proposed splitting strategy could give favorable results in the extension of our scheme to two dimensional systems.

Furthermore, our aspirations include the design of schemes of order higher than 2. The main difficulty that we would face is that the schemes proposed in this thesis cannot be directly extended to higher order, since when we make use of Riemann invariants, we are applying a second order approximation when dividing by the water height. We would then need to come up with an alternative strategy that does not involve the use of Riemann invariants.

In order to obtain higher order schemes, it could also be interesting to consider an IMEX formulation of the schemes. This could also reduce the diffusive aspect observed in numerical results when solving the pressure equation implicitly.

Finally, it would also be interesting to apply the different strategies proposed in this thesis to other type of systems that could benefit from the separation of fast and slow waves. This is the case, for example, of the Euler system with gravity. Moreover, we could also think of addressing other interesting scenarios such as turbidity currents or sediment transport.

# Bibliography

- [1] E. Audusse, F. Bouchut, M. O. Bristeau, R. Klein, and B. Perthame. A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows. *SIAM Journal on Scientific Computing*, 25(6):2050–2065, 2004.
- [2] E. Audusse, C. Chalons, and P. Ung. A simple well-balanced and positive numerical scheme for the shallow-water system. *Communications in Mathematical Sciences*, 13:1317–1332, 2015.
- [3] M. Baudin, C. Berthon, F. Coquel, R. Masson, and Q. H. Tran. A relaxation method for two-phase flow models with hydrodynamic closure law. *Numerische Mathematik*, 99(3):411–440, 2004.
- [4] A. Bermudez and M. E. Vazquez. Upwind methods for hyperbolic conservation laws with source terms. *Computers & Fluids*, 23(8):1049–1071, 1994.
- [5] C. Berthon and C. Chalons. A fully well-balanced, positive and entropy-satisfying Godunov-type method for the shallow-water equations. *Mathematics of Computation*, 85(299):1281–1307, 2015.
- [6] C. Berthon, C. Chalons, S. Cornet, and G. Sperone. Fully well-balanced, positive and simple approximate riemann solver for shallow water equations. *Bulletin of the Brazilian Mathematical Society, New Series*, 47(1):117–130, 2016.
- [7] C. Berthon and F. Foucher. Efficient well-balanced hydrostatic upwind schemes for shallow-water equations. *Journal of Computational Physics*, 231(15):4993–5015, 2012.
- [8] C. Berthon and V. Michel-Dansac. A simple fully well-balanced and entropy preserving scheme for the shallow-water equations. *Applied Mathematics Letters*, 86:284–290, 2018.
- [9] G. Bispen, K. R. Arun, M. Lukáčová-Medvid'ová, and S. Noelle. IMEX large time step finite volume methods for low Froude number shallow water flows. *Communications in Computational Physics*, 16(2):307–347, 2014.



- [10] G. Bispen, M. Lukáčová-Medvid'ová, and L. Yelash. Asymptotic preserving IMEX finite volume schemes for low Mach number Euler equations with gravitation. *Journal of Computational Physics*, 335:222–248, April 2017.
- [11] L. Bonaventura, E. D. Fernández-Nieto, J. Garres-Díaz, and G. Narbona-Reina. Multilayer shallow water models with locally variable number of layers and semi-implicit time discretization. *Journal of Computational Physics*, 364:209–234, 2018.
- [12] S. Boscarino, L. Pareschi, and G. Russo. Implicit-explicit Runge-Kutta schemes for hyperbolic systems and kinetic equations in the diffusion limit. *SIAM Journal on Scientific Computing*, 35(1):A22–A51, 2013.
- [13] S. Boscarino, L. Pareschi, and G. Russo. A unified IMEX Runge–Kutta approach for hyperbolic systems with multiscale relaxation. *SIAM Journal on Numerical Analysis*, 55(4):2085–2109, 2017.
- [14] S. Boscarino and G. Russo. On a class of uniformly accurate IMEX Runge–Kutta schemes and applications to hyperbolic systems with relaxation. *SIAM Journal on Scientific Computing*, 31(3):1926–1945, 2009.
- [15] W. Boscheri, M. Dumbser, M. Ioriatti, I. Peshkov, and E. Romenski. A structure-preserving staggered semi-implicit finite volume scheme for continuum mechanics. *Journal of Computational Physics*, 424:109866, January 2021.
- [16] W. Boscheri and M. Tavelli. On the construction of conservative semi-Lagrangian IMEX advection schemes for multiscale time dependent pdes. *Journal of Scientific Computing*, 90(3), 6 2021.
- [17] W. Boscheri, M. Tavelli, and C. E. Castro. An all Froude high order IMEX scheme for the shallow water equations on unstructured Voronoi meshes. *Applied Numerical Mathematics*, 185:311–335, 2023.
- [18] F. Bouchut. Entropy satisfying flux vector splittings and kinetic BGK models. *Numerische Mathematik*, 94(4):623–672, 2003.
- [19] F. Bouchut. *Nonlinear stability of finite volume methods for hyperbolic conservation laws and well-balanced schemes for sources*. Frontiers in Mathematics. Birkhäuser Verlag, Basel, 2004.
- [20] F. Bouchut and T. Morales de Luna. A subsonic-well-balanced reconstruction scheme for shallow water flows. *SIAM Journal on Numerical Analysis*, 48(5):1733–1758, 2010.
- [21] S. Busto and M. Dumbser. A staggered semi-implicit hybrid finite volume / finite element scheme for the shallow water equations at all Froude numbers. *Applied Numerical Mathematics*, 175:108–132, 2022.

- [22] C. Caballero-Cárdenas, M. J. Castro, T. Morales de Luna, and M. L. Muñoz-Ruiz. Implicit and implicit-explicit Lagrange-projection finite volume schemes exactly well-balanced for 1D shallow water system. *Applied Mathematics and Computation*, 443:127784, 2023.
- [23] M. J. Castro, C. Chalons, and T. Morales de Luna. A fully well-balanced Lagrange-projection-type scheme for the shallow-water equations. *SIAM Journal on Numerical Analysis*, 56(5):3071–3098, 2018.
- [24] M. J. Castro, J. A. López-García, and C. Parés. High order exactly well-balanced numerical methods for shallow water systems. *Journal of Computational Physics*, 246:242–264, 2013.
- [25] M. J. Castro and C. Parés. Well-balanced high-order finite volume methods for systems of balance laws. *Journal of Scientific Computing*, 82(2), 2020.
- [26] V. Casulli. Semi-implicit finite difference methods for the two-dimensional shallow water equations. *Journal of Computational Physics*, 86(1):56–74, 1990.
- [27] V. Casulli. Numerical simulation of three-dimensional free surface flow in isopycnal co-ordinates. *International Journal for Numerical Methods in Fluids*, 25(6):645–658, 1997.
- [28] V. Casulli and E. Cattani. Stability, accuracy and efficiency of a semi-implicit method for three-dimensional shallow water flow. *Computers & Mathematics with Applications*, 27(4):99–112, 1994.
- [29] V. Casulli and R. T. Cheng. Semi-implicit finite difference methods for three-dimensional shallow water flow. *International Journal for Numerical Methods in Fluids*, 15(6):629–648, 1992.
- [30] V. Casulli and R. A. Walters. An unstructured grid, three-dimensional model based on the shallow water equations. *International Journal for Numerical Methods in Fluids*, 32(3):331–348, 2000.
- [31] L. Cea and A. López-Núñez. Extension of the two-component pressure approach for modeling mixed free-surface-pressurized flows with the two-dimensional shallow water equations. *International Journal for Numerical Methods in Fluids*, 93(3):628–652, 2020.
- [32] C. Chalons. *Thesis École Polytechnique Palaiseau, France*. 2002.
- [33] C. Chalons and A. del Grosso. A second-order well-balanced Lagrange-projection numerical scheme for shallow water exner equations in 1D and 2D. *Communications in Mathematical Sciences*, 20(7):1839–1873, 2022.

- [34] C. Chalons, M. Girardin, and S. Kokh. Large Time Step and Asymptotic Preserving Numerical Schemes for the Gas Dynamics Equations with Source Terms. *SIAM Journal on Scientific Computing*, 35(6):A2874–A2902, 2013.
- [35] C. Chalons, M. Girardin, and S. Kokh. An All-Regime Lagrange-Projection Like Scheme for the Gas Dynamics Equations on Unstructured Meshes. *Communications in Computational Physics*, 20(1):188–233, 2016.
- [36] C. Chalons, P. Kestener, S. Kokh, and M. Stauffert. A large time-step and well-balanced Lagrange-projection type scheme for the shallow water equations. *Communications in Mathematical Sciences*, 15(3):765–788, 2017.
- [37] A. Chertock, A. Kurganov, and Y. Liu. Central-upwind schemes for the system of shallow water equations with horizontal temperature gradients. *Numerische Mathematik*, 127(4):595–639, 2013.
- [38] A. Chinnayya, A.-Y. Leroux, and N. Seguin. A well-balanced numerical scheme for the approximation of the shallow-water equations with topography: The resonance phenomenon. *International Journal on Finite Volumes*, 1, 2004.
- [39] F. Coquel, E. Godlewski, B. Perthame, A. In, and P. Rascle. *Some New Godunov and Relaxation Methods for Two-Phase Flow Problems*, pages 179–188. Springer US, New York, NY, 2001.
- [40] R. Courant, K. Friedrichs, and H. Lewy. Über die partiellen differenzgleichungen der mathematischen physik. *Mathematische Annalen*, 100(1):32–74, December 1928.
- [41] A.J.C.B. de Saint-Venant. *Théorie du mouvement non permanent des eaux, avec application aux crues des rivières et à l'introduction des marées dans leur lit*. Comptes rendus hebdomadaires des séances de l'Académie des sciences. Gauthier-Villars, 1871.
- [42] A. del Grosso, M. J. Castro, C. Chalons, and T. Morales de Luna. On well-balanced implicit-explicit Lagrange-projection schemes for two-layer shallow water equations. *Applied Mathematics and Computation*, 442:127702, 2023.
- [43] A. del Grosso and C. Chalons. Second-order well-balanced Lagrange-projection schemes for blood flow equations. *Calcolo*, 58(4), 2021.
- [44] G. Dimarco, R. Loubère, V. Michel-Dansac, and M. H. Vignal. Second-order implicit-explicit total variation diminishing schemes for the Euler system in the low Mach regime. *Journal of Computational Physics*, 372:178–201, November 2018.
- [45] R. Donat, M. C. Martí, A. Martínez-Gavara, and P. Mulet. Well-balanced adaptive mesh refinement for shallow water flows. *Journal of Computational Physics*, 257:937–953, January 2014.

- [46] M. Dumbser and V. Casulli. A staggered semi-implicit spectral discontinuous Galerkin scheme for the shallow water equations. *Appl. Math. Comput.*, 219:8057–8077, 2013.
- [47] E. Franck and L. Navoret. Semi-implicit two-speed well-balanced relaxation scheme for Ripa model. In *Springer Proceedings in Mathematics & Statistics*, pages 735–743. Springer International Publishing, 2020.
- [48] E. Gaburro, M. J. Castro, and M. Dumbser. Well-balanced Arbitrary-Lagrangian-Eulerian finite volume schemes on moving nonconforming meshes for the euler equations of gas dynamics with gravity. *Monthly Notices of the Royal Astronomical Society*, 477(2):2251–2275, 2018.
- [49] J. M. Gallardo, C. Parés, and M. J. Castro. On a well-balanced high-order finite volume scheme for shallow water equations with topography and dry areas. *Journal of Computational Physics*, 227(1):574–601, November 2007.
- [50] G. Gallice. Positive and entropy stable Godunov-type schemes for gas dynamics and MHD equations in Lagrangian or Eulerian coordinates. *Numerische Mathematik*, 94(4):673–713, 2002.
- [51] G. Gallice. Solveurs simples positifs et entropiques pour les systèmes hyperboliques avec terme source. *Comptes Rendus Mathématique*, 334(8):713–716, 2002.
- [52] J. Garres-Díaz, E. D. Fernández-Nieto, and G. Narbona-Reina. A semi-implicit approach for sediment transport models with gravitational effects. *Applied Mathematics and Computation*, 421:126938, 2022.
- [53] F. X. Giraldo and M. Restelli. High-order semi-implicit time-integrators for a triangular discontinuous Galerkin oceanic shallow water model. *International Journal for Numerical Methods in Fluids*, pages n/a–n/a, 2009.
- [54] E. Godlewski and P. A. Raviart. *Hyperbolic systems of conservation laws*. Ellipses-Edition Marketing, 1991.
- [55] E. Godlewski and P.A. Raviart. *Numerical Approximation of Hyperbolic Systems of Conservation Laws*. Number 118 in Applied Mathematical Sciences. Springer, 1996.
- [56] I. Gómez-Bueno, M. J. Castro, C. Parés, and G. Russo. Collocation methods for high-order well-balanced methods for systems of balance laws. *Mathematics*, 9(15):1799, 2021.
- [57] L. Gosse. A well-balanced flux-vector splitting scheme designed for hyperbolic systems of conservation laws with source terms. *Computers & Mathematics with Applications*, 39(9-10):135–159, 2000.



- [58] S. Gottlieb, Z. J. Grant, J. Hu, and R. Shu. High order strong stability preserving MultiDerivative implicit and IMEX Runge–Kutta methods with asymptotic preserving properties. *SIAM Journal on Numerical Analysis*, 60(1):423–449, 2022.
- [59] S. Gottlieb, C.-W. Shu, and E. Tadmor. Strong stability-preserving high-order time discretization methods. *SIAM Review*, 43(1):89–112, 2001.
- [60] N. Goutal and F. Maurel. Proceedings of the 2nd workshop on dam-break wave simulation. *Technical report, Groupe Hydraulique Fluviale, Département Laboratoire National d’Hydraulique, Electricité de France*, 1997.
- [61] J. M. Greenberg and A. Y. Leroux. A well-balanced scheme for the numerical processing of source terms in hyperbolic equations. *SIAM Journal on Numerical Analysis*, 33(1):1–16, 1996.
- [62] A. Harten, P. D. Lax, and B. van Leer. *On Upstream Differencing and Godunov-Type Schemes for Hyperbolic Conservation Laws*, pages 53–79. Springer Berlin Heidelberg, Berlin, Heidelberg, 1997.
- [63] S. Jin and Z. Xin. The relaxation schemes for systems of conservation laws in arbitrary space dimensions. *Communications on Pure and Applied Mathematics*, 48(3):235–276, 1995.
- [64] S. Kang, F. X. Giraldo, and T. Bui-Thanh. IMEX HDG-DG: A coupled implicit hybridized discontinuous Galerkin and explicit discontinuous Galerkin approach for shallow water systems. *Journal of Computational Physics*, 401:109010, 2020.
- [65] C. Klingenberg, G. Puppo, and M. Semplice. Arbitrary order finite volume well-balanced schemes for the Euler equations with gravity. *SIAM Journal on Scientific Computing*, 41(2):A695–A721, 2019.
- [66] R. J. LeVeque. *Numerical methods for conservation laws*, volume 132. Springer, 1992.
- [67] M. Lukáčová-Medvid’ová, J. Rosemeier, P. Spichtinger, and B. Wiebe. IMEX finite volume methods for cloud simulation. pages 179–187, 2017.
- [68] G. I. Marchuk. *Metody rasshchepleniya*. Moskva: Nauka, 1988.
- [69] G. I. Marchuk. Splitting and alternating direction methods. pages 197–462, 1990.
- [70] V. Michel-Dansac and A. Thomann. On high-precision  $l^\infty$ -stable IMEX schemes for scalar hyperbolic multi-scale equations. pages 79–94, 2021.
- [71] V. Michel-Dansac and A. Thomann. TVD-MOOD schemes based on implicit-explicit time integration. *Applied Mathematics and Computation*, 433:127397, 2022.

- [72] T. Morales de Luna, M. J. Castro, and C. Chalons. High-order fully well-balanced Lagrange-projection scheme for shallow water. *Communications in Mathematical Sciences*, 18(3):781–807, 2020.
- [73] L. O. Müller, C. Parés, and E. F. Toro. Well-balanced high-order numerical schemes for one-dimensional blood flow in vessels with varying mechanical properties. *Journal of Computational Physics*, 242:53–85, 2013.
- [74] S. Noelle, Y. Xing, and C.-W. Shu. High-order well-balanced finite volume WENO schemes for shallow water equation with moving water. *J. Comput. Phys.*, 226(1):29–58, 2007.
- [75] C. Parés. Numerical methods for nonconservative hyperbolic systems: a theoretical framework. *SIAM Journal on Numerical Analysis*, 44(1):300–321, 2006.
- [76] L. Pareschi and G. Russo. Implicit–Explicit Runge–Kutta schemes and applications to hyperbolic systems with relaxation. *Journal of Scientific Computing*, 25(1):129–155, 2005.
- [77] B. Perthame and C. Simeoni. A kinetic scheme for the Saint-Venant system with a source term. *Calcolo*, 38:201–231, 2001.
- [78] G. Puppo, M. Semplice, and G. Visconti. Quinpi: Integrating conservation laws with CWENO implicit methods. *Communications on Applied Mathematics and Computation*, 2022.
- [79] P. Ripa. Conservation laws for primitive equations models with inhomogeneous layers. *Geophysical & Astrophysical Fluid Dynamics*, 70(1-4):85–111, 1993.
- [80] P. Ripa. On improving a one-layer ocean model with thermodynamics. *Journal of Fluid Mechanics*, 303:169–201, 1995.
- [81] G. Russo and A. Khe. High order well-balanced schemes based on numerical reconstruction of the equilibrium variables. In *Waves and Stability in Continuous Media*, pages 230–241. World Scientific, April 2010.
- [82] C. Sánchez-Linares, T. Morales de Luna, and M. J. Castro. A HLLC scheme for Ripa model. *Applied Mathematics and Computation*, 272:369–384, 2016.
- [83] S. Serna and A. Marquina. Power ENO methods: a fifth-order accurate weighted power ENO method. *Journal of Computational Physics*, 194(2):632–658, 2004.
- [84] M. N. Spijker. Contractivity in the numerical solution of initial value problems. *Numerische Mathematik*, 42(3):271–290, 1983.



- [85] G. Strang. On the construction and comparison of difference schemes. *SIAM Journal on Numerical Analysis*, 5(3):506–517, 1968.
- [86] I. Suliciu. On modelling phase transitions by means of rate-type constitutive equations. shock wave structure. *International Journal of Engineering Science*, 28(8):829–841, 1990.
- [87] I. Suliciu. Some stability-instability problems in phase transitions modelled by piecewise linear elastic or viscoelastic constitutive equations. *International Journal of Engineering Science*, 30(4):483–494, 1992.
- [88] M. Tavelli and M. Dumbser. A high order semi-implicit discontinuous Galerkin method for the two dimensional shallow water equations on staggered unstructured meshes. *Applied Mathematics and Computation*, 234:623–644, 2014.
- [89] G. Tumolo, L. Bonaventura, and M. Restelli. A semi-implicit, semi-Lagrangian, p-adaptive discontinuous Galerkin method for the shallow water equations. *Journal of Computational Physics*, 232(1):46–67, 2013.
- [90] B. Van Leer. Towards the ultimate conservative difference scheme. II. Monotonicity and conservation combined in a second-order scheme. *Journal of Computational Physics*, 14(4):361–370, 1974.
- [91] S. Vater and R. Klein. A semi-implicit multiscale scheme for shallow water flows at low Froude number. *Communications in Applied Mathematics and Computational Science*, 13(2):303–336, 2018.
- [92] Y. Xing. Exactly well-balanced discontinuous Galerkin methods for the shallow water equations with moving water equilibrium. *Journal of Computational Physics*, 257:536–553, 2014.
- [93] Y. Xing, C.-W. Shu, and S. Noelle. On the advantage of well-balanced schemes for moving-water equilibria of the shallow water equations. *Journal of Scientific Computing*, 48(1-3):339–349, 2010.