

WELL-BALANCED AND INVARIANT DOMAIN PRESERVING SCHEMES FOR  
DISPERSIVE SHALLOW WATER FLOWS.

A Dissertation

by

ERIC JOSEPH TOVAR

Submitted to the Graduate and Professional School of  
Texas A&M University  
in partial fulfillment of the requirements for the degree of  
DOCTOR OF PHILOSOPHY

Chair of Committee,	Jean-Luc Guermond
Committee Members,	Andrea Bonito
	Bojan Popov
	Jean Ragusa
Head of Department,	Sarah Witherspoon

December 2021

Major Subject: Mathematics

Copyright 2021 Eric Joseph Tovar

## ABSTRACT

As urbanization encroaches more on flood prone regions and paved surfaces are ever expanding, more catastrophic flash floods occurring in urban environments are expected in the near future. These risks are compounded by global changes in the climate. Mathematics can help better predict and understand these situations through modeling and numerical simulations. The aim of this work is to discuss current mathematical and computational issues in modeling shallow water flows with applications in coastal hydraulics, large-scale oceanography and in-land flooding.

Our mathematical starting points are the systems of partial differential equations known as the (i) Saint-Venant shallow water equations and (ii) dispersive Serre–Green–Naghdi (SGN) equations. The goal of this work is to efficiently solve both mathematical models supplemented with external physical source terms for in-land flooding and large-scale coastal oceanography applications. In particular, the work focuses on introducing a novel technique for solving the Serre–Green–Naghdi equations. We introduce new analytical solutions of the SGN equations with topography that are used to verify the accuracy of numerical methods. Then, we propose a new relaxation technique for solving the SGN equations with topography effects that yields a hyperbolic formulation of the equations. This relaxation technique allows us to circumvent the dispersive time step restriction of the Serre Equations which is a major challenge when solving the equations. This method is then supplemented with a novel continuous finite element approximation that is second-order accurate in space, invariant domain preserving and well-balanced. The method is then verified with academic benchmarks and validated by comparison with laboratory experimental data.

## DEDICATION

To my mother and my late grandparents, who raised me and shaped me into the person I am. To my wife, whose patience, advice, and unwavering support allowed me to complete this project. To Aleah, my sister at heart for whom I aim to be an adequate role model. To my uncle Javier for his support throughout the years. And to my Tia Rosie who taught me the importance of education and always pushed me to be the best I could be.

## ACKNOWLEDGMENTS

I would like to thank my advisor, Dr. Jean-Luc Guermond, for his investment, support and advice throughout my time at Texas A&M University. I would also like to thank my committee members for their patience and guidance: Dr. Bojan Popov, Dr. Andrea Bonito and Dr. Jean Ragusa; your time and feedback was greatly appreciated. I would also like to thank Dr. Chris Kees, my former supervisor at ERDC, for his support, patience and trust throughout my time at ERDC. I would also like to thank Hwai-Ping (Pearce) Cheng, the branch chief at CHL/ERDC, for believing in me and supporting me through my internship. I would also like to thank Dr. Matthias Maier at TAMU for his help and support in the development of the Ryuji software. Lastly, I am also grateful for my friends and colleagues, Bennett C. and Weston B., for their input, discussions, and company throughout my last five years in this program.

## CONTRIBUTORS AND FUNDING SOURCES

### **Contributors**

This work was supported by a dissertation committee consisting of Professor Jean-Luc Guermont (advisor), Professor Bojan Popov, Professor Andrea Bonito of the Department of Mathematics at Texas A&M University and Professor Jean Ragusa of the Department of Nuclear Engineering at Texas A&M University. Part of this work was done in collaboration with Professor Chris Kees at the Department of Civil & Environmental Engineering at Louisiana State University for the completion of an internship program at the U.S. Army Engineer Research and Development Center. All other work conducted for the dissertation was completed by the student independently.

### **Funding Sources**

Graduate study was supported in part by the National Science Foundation grants DMS-1619892 and DMS-1620058, by the Air Force Office of Scientific Research, USAF, under grant/contract number FA9550-18-1-0397, and by the Army Research Office under grant/contract number W911NF-19-1-0431. Part of the graduate study was supported by the Texas Water Resources Institute (TWRI).

## NOMENCLATURE

$\boldsymbol{x}$	Cartesian position vector
$g$	acceleration due to gravity $9.81 \text{ ms}^{-2}$
$h_0$	reference water depth
$\boldsymbol{u}$	conserved solution variable
$h$	water depth
$\boldsymbol{q}$	momentum vector
$\boldsymbol{v}$	velocity vector
SV	Saint-Venant
SWE	Shallow Water Equations
SGN	Serre–Green–Naghdi
ODE	ordinary differential equation
PDE	partial differential equation
TAMU	Texas A&M University
ERDC	U.S. Army Engineer Research and Development Center

## TABLE OF CONTENTS

	Page
ABSTRACT .....	ii
DEDICATION .....	iii
ACKNOWLEDGMENTS .....	iv
CONTRIBUTORS AND FUNDING SOURCES .....	v
NOMENCLATURE .....	vi
TABLE OF CONTENTS .....	vii
LIST OF FIGURES .....	xi
LIST OF TABLES.....	xiv
1. INTRODUCTION.....	1
1.1 Introduction.....	1
1.2 Hyperbolic systems .....	2
1.3 Dispersive partial differential equations .....	4
1.4 Shallow water flows .....	5
1.4.1 Free-surface Euler equations for water waves .....	5
1.4.2 The water depth and flow discharge formulation .....	9
1.4.3 Decomposition of the pressure.....	12
1.4.4 Decomposition of velocity vector .....	14
1.4.5 Non-dimensionalization of the depth and flow discharge formulation.....	16
1.4.6 Inner structure of the velocity and pressure fields .....	19
1.4.7 First order approximation – Saint-Venant Shallow Water Equations .....	25
1.4.8 Weakly non-linear second order approximation – The Boussinesq model ....	26
1.4.9 Strongly non-linear second order approximation – Serre–Green–Naghdi model.....	27
2. SHALLOW WATER MODELS .....	29
2.1 Introduction.....	29
2.2 The Saint-Venant model .....	30
2.2.1 The model problem.....	30
2.2.2 Saint-Venant properties .....	31
2.2.3 Physical drawbacks of the Saint-Venant model .....	32

2.3	The Serre–Green–Naghdi Equations .....	33
2.3.1	The model problem .....	34
2.3.2	Admissible set and Lake-at-rest .....	35
2.3.3	Mathematical reinterpretation of dispersive terms .....	36
2.3.4	Challenges .....	41
2.3.4.1	Dispersive time step restriction .....	41
2.3.4.2	Verification and validation .....	42
2.3.4.3	Boundary conditions .....	42
2.3.5	Properties .....	43
2.3.5.1	Conservation of energy .....	43
2.3.5.2	Analytical steady-state solution with topography .....	47
2.3.5.3	Dispersion relation .....	50
2.4	External physical sources .....	52
2.4.1	Preliminaries .....	52
2.4.2	Gauckler-Manning friction .....	53
2.4.3	Wave generation and absorption .....	53
3.	REFORMULATION OF THE DISPERSIVE SERRE MODEL .....	56
3.1	Introduction .....	56
3.2	Reformulation under two algebraic constraints .....	56
3.3	Literature review .....	61
4.	A HYPERBOLIC RELAXATION TECHNIQUE FOR SOLVING THE SERRE EQUATIONS .....	64
4.1	Introduction .....	64
4.2	Preliminaries .....	65
4.3	Relaxing the $\{q_1 = h^2; q_3 = \mathbf{q} \cdot \nabla z\}$ constraints .....	65
4.3.1	The correlation between the energy functional and relaxed pressure and source terms .....	66
4.4	The generic relaxed model .....	68
4.5	Properties .....	69
4.5.1	Hyperbolicity for the relaxed model .....	69
4.5.2	Defining the $\Gamma$ function .....	70
4.5.3	Derivation of energy inequality .....	74
4.5.4	Dispersion relation .....	79
4.6	Literature review .....	83
4.6.1	Favrie and Gavriluk model .....	83
5.	APPROXIMATION .....	85
5.1	Introduction .....	85
5.2	Finite element setting .....	86
5.2.1	Finite element representations .....	88
5.3	The low-order method .....	89



5.3.1	Numerical flux and hydrostatic pressure/source .....	89
5.3.2	Well-balancing star states .....	91
5.3.3	PDE source .....	92
5.3.4	External physical sources .....	92
5.3.5	Low-order graph-viscosity coefficients.....	93
5.3.6	Defining the time-step .....	95
5.3.7	Generic low-order update for both models .....	96
5.3.8	Higher-order time stepping .....	96
5.4	Well-balancing and invariant domain preserving properties .....	97
5.5	Local auxiliary states and bounds .....	100
5.6	Provisional high-order method .....	103
5.6.1	Wave generation .....	103
5.6.2	Commutator-based entropy viscosity.....	104
5.6.3	Consistent mass matrix.....	106
5.6.4	Loss of positivity .....	107
5.7	Convex limiting with sources.....	108
5.7.1	Quasiconcave functionals and bounds.....	108
5.7.2	Limiting process .....	110
5.7.3	Application to the system (4.4.1) .....	112
5.7.4	Relaxation of the bounds .....	116
6.	NUMERICAL ILLUSTRATIONS .....	119
6.1	Introduction.....	119
6.2	Preliminaries .....	119
6.3	Convergence tests .....	120
6.3.1	Solitary wave solution of Serre model (2.3.1) .....	120
6.3.2	Method of manufactured solutions for Hyperbolic Serre model (4.4.1) .....	122
6.3.3	Steady-state solution with topography .....	124
6.4	Well-balancing tests .....	125
6.5	Riemann problem for Serre Equations .....	127
6.5.1	1D – Dam break over wet bed .....	128
6.5.2	1D – Dam break over dry bed.....	130
6.5.3	2D – Circular dam break .....	130
6.5.4	2D – Square dam break .....	131
6.6	Academic benchmarks .....	132
6.6.1	1D – Interaction of solitary waves.....	133
6.6.2	2D – Dam Break over three obstacles with friction.....	137
6.6.3	2D – Circular Dam Break with divergence-free velocity .....	138
6.7	Laboratory experiments .....	139
6.7.1	Shoaling of solitary waves over sloped beach.....	139
6.7.2	Periodic waves propagation over a submerged bar .....	142
6.7.3	Propagation of periodic waves over an elliptic shoal .....	144
6.7.4	Propagation of periodic waves over semi-circular shoal.....	148
6.7.5	2D Solitary wave run-up over a conical island.....	150

6.7.6 Propagation over a solitary wave over a triangular shelf with conical island.. 152

REFERENCES ..... 158

## LIST OF FIGURES

FIGURE	Page
1.1 Two-dimensional schematic of the free-surface problem.....	7
1.2 A representation of the physical accuracy of each model derived compared to the full free-surface Euler problem. ....	28
2.1 Comparison of smooth solutions with Hyperbolic Serre model and Saint-Venant model.....	33
2.2 1D Steady state solution for the Serre model .....	50
2.3 A plot of the Serre phase velocity (2.3.23) as a function of the wave number $k$ with $H_0 = 1$ m.....	52
2.4 Propagation of solitary wave profile over flat-bottom with roughness coefficients of $n = \{0, 0.1, 0.2, 0.3, 0.4\} \text{m}^{-\frac{1}{3}} \text{s}$ . ....	54
2.5 Generation and absorption zone schematic. ....	55
2.6 Numerical illustration of wave generation/absorption in the space-time domain. ....	55
4.1 A plot of $\Gamma(x)$ . ....	73
4.2 A representation of the physical accuracy of the models derived in Section 1.4 compared to the full free-surface Euler problem now including the hyperbolic relaxed model.....	78
4.3 The error $\frac{c_p^S - c_p^-}{c_p^S}$ as a function of the wave number $k$ for different water depths. ....	81
4.4 Comparison of slow phase velocity of hyperbolic relaxed model (4.4.1) with $\epsilon = \{1, 2, 4, 8\}$ and Serre phase velocity (2.3.23). ....	82
4.5 Comparison of rapid phase velocity of hyperbolic relaxed model (4.4.1) with $\epsilon = \{1, 2, 4, 8\}$ and Serre phase velocity (2.3.23). ....	82
6.1 Computational solitary wave solution at $T = 50$ s with $I = 200$ . ....	122
6.2 Error tables for well-balancing tests using conical island topography. ....	126
6.3 Figures for well-balancing tests with conical island topography.....	127

6.4	A plot of $\delta_\infty(t)$ as a function of time for $t \in [0, 200 \text{ s}]$ for the Hyperbolic Serre model and the Saint-Venant model for the low and high-order schemes. ....	128
6.5	Numerical solution to 1D dam-break problem with dry bed with $\Delta h = 1.8 \text{ m}$ . ....	129
6.6	Numerical solution to 1D dam-break problem with a wet bed with $\Delta h = 0.4 \text{ m}$ . ....	129
6.7	Numerical solution to 1D dam-break problem with a dry bed with $\Delta h = 1.8 \text{ m}$ . ....	130
6.8	Initial profiles for circular and square Riemann Problems. ....	131
6.9	Circular dam break – Comparison with Hyperbolic Serre (top) and Saint-Venant (bottom) at $t = \{10, 30, 50\} \text{ s}$ . ....	132
6.10	Square dam break – Comparison with Hyperbolic Serre (top) and Saint-Venant (bottom) at $t = \{10, 30, 50\} \text{ s}$ . ....	133
6.11	Two solitary waves – quasi-elastic collision .....	135
6.12	Two solitary waves – inelastic collision.....	135
6.13	Four solitary waves – quasi-elastic collision.....	136
6.14	Four solitary waves – inelastic collision .....	136
6.15	Collision of Eight solitary waves .....	137
6.16	Dam break with bumps – Surface plot of the water elevation $h + z$ at several time snapshots.....	138
6.17	Dam break with bumps – Comparison with Hyperbolic Serre (top) and Saint-Venant (bottom) at $t = \{1, 7.8, 15\}$ .....	139
6.18	Circular dam break with divergent-free velocity field – Comparison with Hyperbolic Serre (top) and Saint-Venant (bottom) at $t = \{1, 5, 10\} \text{ s}$ . ....	140
6.19	Comparison of numerical results with experimental data for solitary wave shoaling experiments of [35]......	141
6.20	Submerged bar set up with gauge locations. ....	142
6.21	Illustration of <i>period-folding</i> with experimental wave gauge 1 with SL case (left) and SH case (right). ....	143
6.22	SL Case. Water elevation at seven gauges. Numerical results using three meshes, $h = \{0.05 \text{ m}, 0.025 \text{ m}, 0.0125 \text{ m}\}$ (solid lines). Experimental data (red points). ....	145
6.23	SH Case. Water elevation at seven gauges. Numerical results using three meshes, $h = \{0.05 \text{ m}, 0.025 \text{ m}, 0.0125 \text{ m}\}$ (solid lines). experimental data (red points). ....	146

6.24	Elliptic shoal experiments .....	147
6.25	Elliptic shoal – Comparison of numerical results (3 mesh refinements) with the experimental data along the 8 sections. Experimental data: red triangles. ....	149
6.26	Whalin semi-circular shoal results .....	150
6.27	Experiment 4 – Surface plot of the water elevation $h + z$ at several times for Case C. The thin grey cylinders represent the wave gauges WG3, WG6, WG9, WG16, WG22 (left to right). ....	151
6.28	Experiment 4 – Temporal series over the period $t \in [0, 12s]$ of the free surface elevation $h + z$ in meters at the four WGs (blue solid) compared to the experimental data (red circles) for Case B (on the left) and Case C (on the right). ....	153
6.29	(a) Coordinates of the wave gauges and ADVs in meters; (b) Overview of their respective locations on the bathymetry. ....	155
6.30	(a) Temporal series over the period $t \in [0, 40s]$ of the free surface elevation $h + z$ compared to the experimental data (red dashed). The TAMU code results are in blue (solid) and Proteus code results in black (solid). (b) Temporal series over the period $t \in [0, 40s]$ of velocity $v$ (blue solid) and experimental ADVs (red dashed). ..	156
6.31	Surface plot of the water elevation $h + z$ at several times. ....	157

## LIST OF TABLES

TABLE	Page
6.1 Convergence rates for solitary wave solution $T = 50$ s, CFL = 0.05. ....	121
6.2 Convergence table using $\ h - h_h\ _{L^1}/\ h\ _{L^1}$ for solitary wave solution of Serre model (2.3.1). $T = 50$ s, CFL= 0.05. ....	122
6.3 Convergence rates using manufactured solution. $T = 50$ s, CFL = 0.05 .....	124
6.4 Convergence rates table for steady state solution with topography. ....	125
6.5 Solitary wave shoaling experiment [35] – configuration values .....	141

# 1. INTRODUCTION

## 1.1 Introduction

In this chapter, we provide an introduction to the thesis and discuss the background and motivation for the research. In particular, we introduce the concepts of hyperbolic systems of conservation laws and dispersive partial differential equations which are crucial for this work. We then introduce the notion of shallow water flows by deriving the mathematical models of interest from the free-surface Euler equations for water waves: (i) the Saint-Venant shallow water equations and (ii) the Serre–Green–Naghdi equations.

In collaboration with the Coastal and Hydraulics Laboratory at the U.S. Army Engineer Research and Development Center, the goal of this present research is to provide accurate and rapid modeling of shallow water flows with an emphasis on dispersive water waves for applications in coastal hydrodynamics. In particular, we are interested in the prediction and modeling of inland and coastal flooding hazards due to events such as dam breaks, tsunami waves and hurricane surges. In this work, dispersive water waves are defined as waves whose phase speeds are dependent on their wavelength. This phenomena is paramount for modeling smooth periodic waves (and superposition of such waves) and their non-linear interactions in the near-shore region. Such interactions lead to physical processes such as wave shoaling, wave run-up and wave breaking. Dispersion is also an important property for the propagation of tsunami waves in the deep ocean and in the near-shore region when topography effects are important. The mathematical equations describing such physical phenomena are strongly non-linear and pose numerous challenges when one tries to solve them numerically.

Our starting point are the systems of partial differential equations known as the (i) Saint-Venant Shallow Water Equations (SWE) and (ii) dispersive Serre–Green–Naghdi (SGN) equations. The goal of this work is to efficiently solve both mathematical models supplemented with external physical source terms for in-land flooding and large-scale coastal oceanography applications. In

particular, the work focuses on introducing a novel technique for solving Serre–Green–Naghdi equations. That is to say, we propose a new relaxation technique for solving the SGN equations with topography effects that yields a hyperbolic formulation of the equations. This relaxation technique allows us to circumvent the dispersive time step restriction of the Serre Equations which is a major challenge when solving the equations. This method is then supplemented with a continuous finite element approximation that is second-order accurate in space, invariant domain preserving and well-balanced. The method is then verified with academic benchmarks and validated by comparison with laboratory experimental data.

## 1.2 Hyperbolic systems

The mathematical models of interest in this work are those with dominant *hyperbolic* features and are essential for applications in computational fluid dynamics. These mathematical models propagate “wave-like” profiles and exhibit special solutions such as shock and expansion waves. The theory of hyperbolic systems of conservation laws dates back to the 1950s in the seminal work of Lax, Hopf, Glimm and many others. For an overview on conservation laws and hyperbolic systems, we refer the reader to the books of Ern and Guermond [18, Chap. 80] and Godlewski and Raviart [24] and references therein. The following material in this section is an overview of [18, Chap. 80.1.2].

Let  $D$  be a polyhedral domain in  $\mathbb{R}^d$  (where  $d$  is the spatial dimension). Let  $\mathbf{x} \in \mathbb{R}^d$  denote the usual Cartesian position vector. Let  $\mathcal{A}$  be a subset of the space  $\mathbb{R}^m$  (where  $m \in \mathbb{N} \setminus \{0\}$ ) henceforth referred to as the *admissible set of states*. Let  $\mathbf{u}$  be a vector of conserved quantities (we call this the conserved variable) which takes values from the set  $\mathcal{A}$ . Let  $\mathbb{f} \in \text{Lip}(\mathcal{A}; \mathbb{R}^{m \times d})$  be a given flux. For a generic state  $\mathbf{v} \in \mathcal{A}$ , the flux  $\mathbb{f}$  is a matrix with entries  $\mathbb{f}_{il}(\mathbf{v})$  for all  $i \in \{1 : m\}$  and  $l \in \{1 : d\}$ . In this thesis, we are interested in first-order systems of the form:

$$\partial_t \mathbf{u} + \nabla \cdot \mathbb{f}(\mathbf{u}) = \mathbf{0}, \tag{1.2.1}$$



where  $\nabla \cdot \mathbf{f}(\mathbf{u})$  evaluated at  $(\mathbf{x}, t)$  is a column vector in  $\mathbb{R}^m$  with entries:

$$(\nabla \cdot \mathbf{f}(\mathbf{u}))_i := \sum_{l \in \{1:d\}} \partial_{x_l} \mathbb{f}_{il}(\mathbf{u}(\mathbf{x}, t)), \quad \forall i \in \{1 : m\}.$$

*Remark 1.2.1.* We say that a system of partial differential equations is a first-order system if it involves only first derivatives in space and time.

**Definition 1.2.1** (Hyperbolicity). We say that a system of the form (1.2.1) is hyperbolic if and only if the matrix  $\mathbb{A}(\mathbf{v}, \mathbf{n}) \in \mathbb{R}^{m \times m}$  with entries:

$$(\mathbb{A}(\mathbf{v}, \mathbf{n}))_{ij} := \sum_{l \in \{1:d\}} n_l \partial_{v_j} \mathbb{f}_{il}(\mathbf{v}), \quad \forall i, m \in \{1 : m\}, \quad (1.2.2)$$

is diagonalizable over  $\mathbb{R}$  for all  $\mathbf{v} \in \mathcal{A}$  and all unit vectors  $\mathbf{n} \in \mathbb{R}^d$ .

This definition of hyperbolicity is needed in Chapter 4 when we introduce a novel hyperbolic relaxation technique for solving the Serre–Green–Naghdi equations. In this work, we will consider the *Cauchy problem* for hyperbolic (and similar) systems such as the Saint-Venant equations and Serre–Green–Naghdi equations. The Cauchy problem is defined as follows:

**Definition 1.2.2** (Cauchy problem for hyperbolic system). For  $(\mathbf{x}, t) \in D \times \mathbb{R}_+$ , the Cauchy problem for a hyperbolic system of conservation laws is given by:

$$\begin{cases} \partial_t \mathbf{u} + \nabla \cdot \mathbf{f}(\mathbf{u}) = 0, \\ \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}) \in \mathcal{A}, \end{cases} \quad (1.2.3)$$

where  $\mathbb{R}_+ := [0, \infty)$ .

A particular Cauchy problem of interest in the context of hyperbolic systems is called the Riemann problem. The Riemann problem is a Cauchy problem with discontinuous initial data consisting of two constant states. More precisely, it is defined as follows:

**Definition 1.2.3** (Riemann Problem). Let  $(\mathbf{u}_L, \mathbf{u}_R) \in \mathcal{A} \times \mathcal{A}$  be a pair of constant admissible states. Let  $\mathbf{n} \in \mathbb{R}^d$  be a unit vector. Then, we define the one-dimensional Riemann problem as follows:

$$\partial_t \mathbf{u} + \partial_x (\mathbb{f}(\mathbf{u}) \cdot \mathbf{n}) = 0, \quad \mathbf{u}(\mathbf{x}, 0) = \begin{cases} \mathbf{u}_L, & \text{if } x \leq 0, \\ \mathbf{u}_R, & \text{if } x > 0, \end{cases} \quad (1.2.4)$$

Understanding the solution to (1.2.4) is essential for constructing robust numerical schemes. We refer the reader to Godlewski and Raviart [24, Chap. 1] and Ern and Guermond [18, Chap. 80] and for a more general discussion on the Riemann Problem for hyperbolic systems. We also refer the reader to Toro [61] for a detailed study on the Riemann Problem for the Saint-Venant shallow water equations.

### 1.3 Dispersive partial differential equations

In this work, we are also interested in dispersive partial differential equations. In particular, we are interested in dispersive PDEs because they admit special smooth solutions such as periodic waves and solitary waves which are of interest in the context of coastal hydrodynamics. The main mathematical model of interest is the Serre–Green–Naghdi equations which is a *dispersive* system. In this work, we follow Whitham [64, Chap. 1] and define a dispersive partial differential equation as follows:

**Definition 1.3.1** (Linear dispersive partial differential equation). Let  $(\mathbf{x}, t) \in D \times \mathbb{R}_+$ . Then, a linear partial differential equation of the form  $\partial_t \mathbf{u}(\mathbf{x}, t) = \mathcal{L}\mathbf{u}(\mathbf{x}, t)$  is dispersive if it admits a solution of the form:

$$\mathbf{u}(\mathbf{x}, t) = a \cos(\mathbf{k} \cdot \mathbf{x} - \omega t),$$

where the wave frequency  $\omega$  is a definite real function of the wave number  $k = |\mathbf{k}|$  and if  $\omega''(k) \neq 0$ .

*Remark 1.3.1* (Non-linear dispersive systems). We say that a system of non-linear partial differential equations is dispersive if its linearized counterpart satisfies Definition 1.3.1.

## 1.4 Shallow water flows

In this section, we introduce the notion of shallow water flows. In particular, we derive the two mathematical models of interest from the free-surface Euler equations for water waves: (i) the Saint-Venant shallow water equations and (ii) the Serre–Green–Naghdi equations for dispersive water waves. The derivation of the models is presented in an Eulerian reference frame and follows an asymptotic expansion approach based on two non-dimensional quantities  $\varepsilon$  and  $\mu$  (which are to be defined). The derivation presented in this section is not new and can be found throughout the literature, but is given for completeness. In particular, the derivation shown here loosely follows that of Lannes [44]. We refer the reader to [44] for a more thorough overview of the derivation of various shallow water models. We also refer the reader to Alvarez-Samaniego and Lannes [1] where a rigorous justification for such approach is given.

### 1.4.1 Free-surface Euler equations for water waves

Let  $\boldsymbol{x} \in \mathbb{R}^d$  ( $d = 1, 2$ ) be the horizontal Cartesian position vector and let  $z$  be the vertical Cartesian coordinate. Let  $t \geq 0$  denote the time variable. Let  $\nabla_{\boldsymbol{x}, z}(\cdot)$  denote the gradient operator on the full Cartesian position vector  $(\boldsymbol{x}, z)$  and let  $\nabla$  be the  $\mathbb{R}^d$ -dimensional gradient with respect to the horizontal position vector  $\boldsymbol{x}$ .

We consider the movement of a fluid with a free-surface under the action of gravity and assume that the fluid is incompressible (i.e., constant material density), inviscid (i.e., having negligible viscosity) and irrotational (i.e., no vorticity). Let  $\rho \in \mathbb{R}$  denote the fluid density. Assume that the fluid is bounded above by a free-surface  $z = \xi(\boldsymbol{x}, t)$  and bounded below by a bottom surface that is independent of time  $z = z_b(\boldsymbol{x})$ . We call this bottom surface the topography map (this nomenclature is often interchanged with bathymetry map). Assume that at rest, the free-surface is given by the rest state  $z = \xi_0 \in \mathbb{R}$  where  $\xi_0$  is a reference elevation. We define the deviation of the free-surface from its rest state by  $\eta(\boldsymbol{x}, t) := \xi(\boldsymbol{x}, t) - \xi_0$ . We introduce the water depth quantity  $h(\boldsymbol{x}, t) := \xi(\boldsymbol{x}, t) - z_b(\boldsymbol{x}) > 0$  so that  $\xi(\boldsymbol{x}, t) = h(\boldsymbol{x}, t) + z_b(\boldsymbol{x})$ . Equivalently, we can write  $h(\boldsymbol{x}, t) = \eta(\boldsymbol{x}, t) + \xi_0 - z_b(\boldsymbol{x})$ . Note that at rest, the reference water depth is given by

$h_0 := \xi_0 - z_b(\mathbf{x})$  (since at rest,  $\eta = 0$ ). The part of the domain that is occupied by the fluid is defined by:

$$D(t) := \{(\mathbf{x}, z) \in \mathbb{R}^d \times \mathbb{R} \mid z_b(\mathbf{x}) < z < \xi(\mathbf{x}, t)\}. \quad (1.4.1)$$

Note that the domain is time-dependent. We illustrate the fluid domain set-up in Figure 1.1. Let

$$\mathbb{R}^{d+1} \ni \mathbf{V}(\mathbf{x}, z, t) := (u(\mathbf{x}, z, t), v(\mathbf{x}, z, t), w(\mathbf{x}, z, t))^T,$$

denote the velocity of a fluid particle located at  $(\mathbf{x}, z)$  at time  $t$ . Let  $\mathbf{v} := (u, v)$  denote the horizontal component of the fluid velocity and  $w$  the vertical component. Then, the free-surface Euler equations for an inviscid, incompressible and irrotational fluid under the action of gravity are given as follows:

$$\partial_t \mathbf{V} + \mathbf{V} \cdot \nabla_{\mathbf{x}, z} \mathbf{V} = -\frac{1}{\rho} \nabla_{\mathbf{x}, z} p - g \mathbf{e}_z, \quad \text{in } D(t), \quad (1.4.2a)$$

$$\nabla_{\mathbf{x}, z} \cdot \mathbf{V} = 0, \quad \text{in } D(t), \quad (1.4.2b)$$

$$\nabla_{\mathbf{x}, z} \times \mathbf{V} = \mathbf{0}, \quad \text{in } D(t), \quad (1.4.2c)$$

where  $p = p(\mathbf{x}, z, t)$  denotes the pressure in the fluid domain  $D(t)$ ,  $g = 9.81 \text{ ms}^{-2}$  is the gravitational constant and  $\mathbf{e}_z$  denotes the unit normal vector in the  $z$ -direction. The equation (1.4.2a) describes the evolution of the fluid particle velocity under the action of gravity (and can be thought of as a balance of forces), (1.4.2b) represents the incompressibility condition of the fluid and (1.4.2c) represents the irrotational condition. Since the flow is incompressible (i.e.,  $\nabla_{\mathbf{x}, z} \cdot \mathbf{V} = 0$ ), the equation (1.4.2a) can be equivalently re-written as

$$\partial_t \mathbf{V} + \nabla_{\mathbf{x}, z} \cdot (\mathbf{V} \otimes \mathbf{V}) = -\frac{1}{\rho} \nabla_{\mathbf{x}, z} p - g \mathbf{e}_z \quad (1.4.3)$$

where  $\otimes$  denotes the outer product of two vectors. This can be shown with a simple expansion of the advection term.

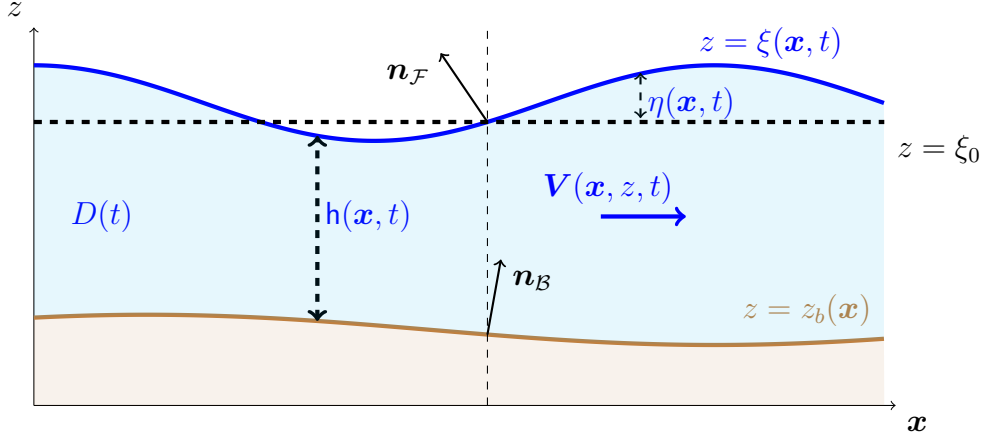


Figure 1.1: Two-dimensional schematic of the free-surface problem.

To complete the system (1.4.2), we need to consider boundary conditions on both the free-surface boundary and bottom boundary. The first boundary condition is the kinematic condition for the free-surface which says that the fluid particles can not cross the free-surface:

$$\mathbf{V} \cdot \mathbf{n}_F = \left( \frac{d}{dt} \mathbf{x}_F \right) \cdot \mathbf{n}_F, \quad z = \xi(\mathbf{x}, t). \quad (1.4.4)$$

Here,  $\mathbf{x}_F := (\mathbf{x}, z = \xi(\mathbf{x}, t))$  is the position coordinate evaluated on the free-surface. The quantity  $\mathbf{n}_F$  is the outward facing normal of the free-surface and is defined by:

$$\mathbf{n}_F := (-\nabla \xi, 1)^T. \quad (1.4.5)$$

Note that the kinematic condition (1.4.4) can be re-written as follows:

$$\partial_t \xi = \mathbf{V}_s \cdot \mathbf{n}_F, \quad (1.4.6)$$

where  $\mathbf{V}_s(\mathbf{x}, t) := \mathbf{V}(\mathbf{x}, z = \xi(\mathbf{x}, t), t)$  denotes the evaluation of the particle fluid velocity at the free-surface where the subscript “s” is meant to represent free-{s}urface. Note that we can expand

the kinematic condition even further:

$$\partial_t \mathbf{h} + \mathbf{v}_s \cdot \nabla (\mathbf{h} + z_b) = w_s, \quad (1.4.7)$$

where we used the definitions  $\xi(\mathbf{x}, t) := \mathbf{h}(\mathbf{x}, t) + z_b(\mathbf{x})$  and  $\mathbf{V} := (\mathbf{v}, w)^\top$ . The second boundary condition we consider is the impermeability of the bottom boundary:

$$\mathbf{V} \cdot \mathbf{n}_b = 0, \quad z = z_b(\mathbf{x}), \quad (1.4.8)$$

where  $\mathbf{n}_b$  is the outward facing normal for the bottom boundary defined by:

$$\mathbf{n}_b := (-\nabla z_b, 1)^\top. \quad (1.4.9)$$

This impermeable condition can be expanded as follows:

$$-\mathbf{v}_b \cdot \nabla z_b + w_b = 0, \quad (1.4.10)$$

where the subscript “b” denotes the evaluation at the bottom boundary  $z = z_b(\mathbf{x})$ . The final boundary condition that we consider is the so called dynamic boundary condition (i.e., we neglect surface tension effects):

$$p(\mathbf{x}, z = \xi(\mathbf{x}, t), t) = p_{\text{atm}} \in \mathbb{R}, \quad (1.4.11)$$

where the constant  $p_{\text{atm}}$  is the atmospheric pressure. For simplicity, we assume that  $p_{\text{atm}} = 0$ .

*Remark 1.4.1.* When discussing the evaluation of quantities at the free-surface or bottom boundaries, we are assuming *a priori* that these operations are justified. That is to say, we do not make any rigorous statements in this derivation. We refer the reader to Alvarez-Samaniego and Lannes [1] where such operations are justified.

The system of partial differential equations (1.4.2) coupled with the boundary conditions (1.4.7)–(1.4.10)–(1.4.11) form a free-surface problem. That is to say, the domain of interest  $D(t)$  is itself

an unknown since it is determined by the unknown free-surface quantity  $\xi(\mathbf{x}, t)$ . Thus, in order to solve this free-surface problem, it is necessary to recast it into an equivalent formulation in which the equations are solved in a fixed domain. In this work, we follow the derivation of [44] and recast the free-surface problem using a formulation based on the water depth and flow discharge.

### 1.4.2 The water depth and flow discharge formulation

In this section, we introduce the water depth and flow discharge formulation of the water waves problem. The goal of this formulation is to cast the free-surface problem in a fixed domain by removing the dependency of the vertical variable  $z$ .

We begin the derivation by introducing the horizontal flow discharge (also known as the momentum):

$$\mathbf{q}(\mathbf{x}, t) := \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \mathbf{v}(\mathbf{x}, z, t) dz. \quad (1.4.12)$$

Recalling that  $\xi(\mathbf{x}, t) := h(\mathbf{x}, t) + z_b(\mathbf{x})$ , we also define the average horizontal velocity:

$$\bar{\mathbf{v}}(\mathbf{x}, t) := \frac{1}{h(\mathbf{x}, t)} \int_{z_b(\mathbf{x})}^{h(\mathbf{x}, t) + z_b(\mathbf{x})} \mathbf{v}(\mathbf{x}, z, t) dz. \quad (1.4.13)$$

We now integrate the incompressibility condition (1.4.2b) over the water column (that is, from  $z = z_b(\mathbf{x})$  to  $z = \xi(\mathbf{x}, t)$ ):

$$\int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} (\nabla \cdot_{\mathbf{x}, z} \mathbf{V}) dz = 0, \quad (1.4.14)$$

$$\implies \quad (1.4.15)$$

$$\int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} (\nabla \cdot \mathbf{v}) dz + w_s - w_b = 0. \quad (1.4.16)$$

Note that an application of the Leibniz integral rule yields the following relation:

$$\nabla \cdot (h\bar{\mathbf{v}}) = \nabla \xi \cdot \mathbf{v}_s - \nabla z_b \cdot \mathbf{v}_b + \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} (\nabla \cdot \mathbf{v}) dz. \quad (1.4.17)$$

Then using the kinematic condition (1.4.7), impermeable condition (1.4.10) and (1.4.16), the above

can be re-written as follows:

$$\partial_t h + \nabla \cdot (h \bar{\mathbf{v}}) = 0.$$

This equation gives an evolution equation for the water depth  $h(\mathbf{x}, t)$ . This equation is also known as the mass conservation equation and can be equivalently written as:

$$\partial_t h + \nabla \cdot \mathbf{q} = 0. \tag{1.4.18}$$

Notice that this equation does not depend on the vertical coordinate  $z$ .

We now want to find an evolution equation for the flow discharge  $\mathbf{q} := h \bar{\mathbf{v}}$ . This can be done by integrating the horizontal part of the velocity evolution equations (1.4.2a) over the vertical water column. Recall that the horizontal components of (1.4.2a) can be expressed as follows:

$$\partial_t \mathbf{v} + \nabla \cdot (\mathbf{v} \otimes \mathbf{v}) + \partial_z (w \mathbf{v}) = -\frac{1}{\rho} \nabla p.$$

We will integrate each term individually and apply the Leibniz integral rule to simplify the time derivative term and advection term. We have that for the time derivative term:

$$\begin{aligned} \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \partial_t \mathbf{v} dz &= \partial_t \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \mathbf{v} dz \right) - \mathbf{v}_s \partial_t \xi \\ &= \partial_t (h \bar{\mathbf{v}}) - \mathbf{v}_s \partial_t h. \end{aligned}$$

The pressure term gives:

$$\begin{aligned} \frac{1}{\rho} \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \nabla p dz &= \frac{1}{\rho} \nabla \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} p dz \right) - \frac{1}{\rho} \underbrace{p_s \nabla (h + z_b)}_{=0} + \frac{1}{\rho} p_b \nabla z_b \\ &= \frac{1}{\rho} \nabla \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} p dz \right) + \frac{1}{\rho} p_b \nabla z_b. \end{aligned}$$

Note that when the topography map is flat, the last term on the right hand side disappears. We now want to integrate the advection term. For simplicity, we expand the advection term into its



respective components and then integrate. Recalling that  $\mathbf{v} := (u, v)$ , note that:

$$\mathbf{v} \otimes \mathbf{v} = \begin{pmatrix} u^2 & uv \\ uv & v^2 \end{pmatrix},$$

and

$$\nabla \cdot (\mathbf{v} \otimes \mathbf{v}) = \left( \partial_x(u^2) + \partial_y(uv), \partial_x(uv) + \partial_y(v^2) \right)^\top.$$

Then integrating each component gives:

$$\begin{aligned} \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} (\partial_x(u^2) + \partial_y(uv)) dz &= \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \partial_x(u^2) dz + \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \partial_y(uv) dz \\ &= \partial_x \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} u^2 dz \right) - u_s^2 \partial_x \xi + u_b^2 \partial_x z_b \\ &\quad + \partial_y \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} uv dz \right) - u_s v_s \partial_y \xi + u_b v_b \partial_y z_b, \end{aligned}$$

and

$$\begin{aligned} \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} (\partial_y(v^2) + \partial_x(uv)) dz &= \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \partial_y(v^2) dz + \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \partial_x(uv) dz \\ &= \partial_y \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} v^2 dz \right) - v_s^2 \partial_y \xi + v_b^2 \partial_y z_b \\ &\quad + \partial_x \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} uv dz \right) - u_s v_s \partial_x \xi + u_b v_b \partial_x z_b. \end{aligned}$$

Then, integrating the final term yields:

$$\int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \partial_z(w\mathbf{v}) dz = w_s \mathbf{v}_s - w_b \mathbf{v}_b.$$

Note that all the boundary terms can be combined as follows:

$$- \mathbf{v}_s \underbrace{(\partial_t \mathbf{h} + \mathbf{v}_s \cdot \nabla \xi - w_s)}_{=0} + \mathbf{v}_b \underbrace{(\mathbf{v}_b \cdot \nabla z_b - w_b)}_{=0} + \frac{1}{\rho} p_b \nabla z_b = \frac{1}{\rho} p_b \nabla z_b, \quad (1.4.19)$$

where the kinematic condition (1.4.7) and (1.4.10) were used to simplify the terms above. Then, the final form of the water depth and flow discharge formulation is given as follows:

$$\partial_t \mathbf{h} + \nabla \cdot \mathbf{q} = 0, \quad (1.4.20a)$$

$$\partial_t \mathbf{q} + \nabla \cdot \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} (\mathbf{v} \otimes \mathbf{v}) dz \right) + \frac{1}{\rho} \nabla \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} p dz \right) = -\frac{1}{\rho} p_b \nabla z_b(\mathbf{x}). \quad (1.4.20b)$$

This can be equivalently written as follows:

$$\partial_t \mathbf{h} + \nabla \cdot (\mathbf{h} \bar{\mathbf{v}}) = 0, \quad (1.4.21a)$$

$$\partial_t (\mathbf{h} \bar{\mathbf{v}}) + \nabla \cdot \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} (\mathbf{v} \otimes \mathbf{v}) dz \right) + \frac{1}{\rho} \nabla (\mathbf{h} \bar{p}) = -\frac{1}{\rho} p_b \nabla z_b(\mathbf{x}), \quad (1.4.21b)$$

where  $\bar{p}$  is the average pressure defined by:

$$\bar{p}(\mathbf{x}, t) = \frac{1}{\mathbf{h}} \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} p dz.$$

### 1.4.3 Decomposition of the pressure

The next step in the derivation of the shallow water models of interest is the decomposition of the pressure term in (1.4.21). We first note that the Euler Equations (1.4.2) admit a particular solution when the flow is at rest:  $\xi(\mathbf{x}, t) = \xi_0$  and  $\mathbf{V} = \mathbf{0}$ . At rest, the system (1.4.2) simplifies to the following ordinary differential equation for the pressure:

$$-\frac{1}{\rho} \partial_z p - g = 0, \quad p|_{\xi=\xi_0} = 0.$$

A simple integration gives the solution  $p(\mathbf{x}, z, t) = -\rho g z + \rho g \xi_0$ . This quantity is the so-called hydrostatic pressure. When the fluid is not at rest, the solution to following ordinary differential equation:

$$-\frac{1}{\rho} \partial_z p - g = 0, \quad p|_{\xi} = 0,$$

is also known as the hydrostatic pressure and is given by:  $p_H = -\rho g(z - \xi)$ . We will use this special solution to decompose the pressure in the depth/discharge formulation (1.4.21). That is to say, we decompose the pressure into its hydrostatic part ( $p_H$ ) and non-hydrostatic part ( $p_{NH}$ ):

$$p(\mathbf{x}, z, t) = p_H + p_{NH} = \rho g(\xi - z) + p_{NH}. \quad (1.4.22)$$

We can find an expression for the non-hydrostatic pressure by integrating the vertical component of (1.4.3) and using (1.4.22). Recall that the vertical component of (1.4.3) is given by:

$$\begin{aligned} \partial_t w + \nabla_{\mathbf{x}, z} \cdot (w \mathbf{V}) &= -\frac{1}{\rho} \partial_z p - g \\ &= -\frac{1}{\rho} \partial_z (\rho g(\xi - z) + p_{NH}) - g \\ &= -\frac{1}{\rho} \partial_z (p_{NH}). \end{aligned}$$

Then integrating with respect to  $z$  yields:

$$p_{NH}(\mathbf{x}, z, t) = \rho \int_z^{\xi(\mathbf{x}, t)} \left( \partial_t w + \nabla_{\mathbf{x}, z} \cdot (w \mathbf{V}) \right) dz. \quad (1.4.23)$$

Thus, we can write the full pressure as:

$$p(\mathbf{x}, z, t) = \rho g(\xi(\mathbf{x}, t) - z) + \rho \int_z^{\xi(\mathbf{x}, t)} \left( \partial_t w + \nabla_{\mathbf{x}, z} \cdot (w \mathbf{V}) \right) dz. \quad (1.4.24)$$

Consequently, we have the following relations:

$$\begin{aligned} p_b &= \rho g(\xi(\mathbf{x}, t) - z_b) + \rho \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \left( \partial_t w + \nabla_{\mathbf{x}, z} \cdot (w \mathbf{V}) \right) dz \\ &= \rho g h + \rho \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \left( \partial_t w + \nabla_{\mathbf{x}, z} \cdot (w \mathbf{V}) \right) dz \\ &= \rho g h + (p_{NH})_b, \end{aligned}$$

and

$$\begin{aligned}
h\bar{p} &= \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} p(\mathbf{x}, z, t) dz \\
&= \frac{1}{2}\rho g(\xi - z_b)^2 + \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} p_{NH} dz \\
&= \frac{1}{2}\rho g h^2 + \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} p_{NH} dz.
\end{aligned}$$

With this decomposition of the pressure, we can now re-write the water depth/flow discharge formulation (1.4.21) as follows:

$$\partial_t h + \nabla \cdot (h\bar{\mathbf{v}}) = 0, \quad (1.4.25a)$$

$$\begin{aligned}
\partial_t(h\bar{\mathbf{v}}) + \nabla \cdot \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} (\mathbf{v} \otimes \mathbf{v}) dz \right) + \nabla \left( \frac{1}{2}gh \right) + \frac{1}{\rho} \nabla \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} p_{NH} dz \right) \\
= -gh \nabla z_b(\mathbf{x}) - \frac{1}{\rho} (p_{NH})_b \nabla z_b(\mathbf{x}). \quad (1.4.25b)
\end{aligned}$$

By the Leibniz integral rule, the momentum equation (1.4.25b) is equivalent to:

$$\partial_t(h\bar{\mathbf{v}}) + \nabla \cdot \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} (\mathbf{v} \otimes \mathbf{v}) dz \right) + \nabla \left( \frac{1}{2}gh \right) + \frac{1}{\rho} \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} \nabla p_{NH} dz \right) = -gh \nabla z_b(\mathbf{x}).$$

#### 1.4.4 Decomposition of velocity vector

The equation (1.4.25b) is an evolution equation for the quantity  $h\bar{\mathbf{v}}$  where  $\bar{\mathbf{v}}(\mathbf{x}, t)$  is the average horizontal velocity that does not depend on vertical Cartesian coordinate  $z$ . Note that the advection term  $\nabla \cdot \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x},t)} (\mathbf{v} \otimes \mathbf{v}) dz \right)$  contributes to the system's dependence on the vertical coordinate  $z$ . To separate this dependence on the vertical coordinate, we introduce a decomposition of the horizontal velocity vector:

$$\mathbf{v}(\mathbf{x}, z, t) = \bar{\mathbf{v}}(\mathbf{x}, t) + \mathbf{v}^*(\mathbf{x}, z, t), \quad (1.4.26)$$

where  $\mathbf{v}^*(\mathbf{x}, z, t)$  measures the deviation of the average horizontal velocity  $\bar{\mathbf{v}}(\mathbf{x}, t)$  from the horizontal velocity  $\mathbf{v}$ . Recalling that  $\bar{\mathbf{v}} := \frac{1}{h} \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \mathbf{v} dz$ , integrating the decomposition (1.4.26) over the water column yields:

$$\begin{aligned} \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \mathbf{v} dz &= \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \left( \bar{\mathbf{v}}(\mathbf{x}, t) + \mathbf{v}^*(\mathbf{x}, z, t) \right) dz, \\ &\implies \\ h\bar{\mathbf{v}}(\mathbf{x}, t) &= h\bar{\mathbf{v}}(\mathbf{x}, t) + \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \mathbf{v}^*(\mathbf{x}, z, t) dz, \\ &\implies \\ \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \mathbf{v}^*(\mathbf{x}, z, t) dz &= 0. \end{aligned}$$

That is to say, the mean-value over the water column of the deviations  $\mathbf{v}^*$  is zero. Using this property and recalling the definition of the momentum quantity,  $\mathbf{q} := h\bar{\mathbf{v}}$ , we can express the advection term as follows:

$$\nabla \cdot \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} (\mathbf{v} \otimes \mathbf{v}) dz \right) = \nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) + \nabla \cdot \mathbf{R},$$

where

$$\mathbf{R} = \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \mathbf{v}^* \otimes \mathbf{v}^* dz.$$

Since the quantity  $\mathbf{v}^* = \mathbf{v} - \bar{\mathbf{v}}$  measures the deviation of the average horizontal velocity from the horizontal velocity,  $\mathbf{R}$  can be thought of as the contribution of these fluctuations to the momentum equations.

*Remark 1.4.2.* As stated in [44], the quantity  $\mathbf{R}$  can be thought of as analog to the Reynolds stress tensor in turbulence and can thus be used to measure the effects of “turbulence” in the model.

Proceeding, we can again re-write the depth/discharge formulation with this velocity decom-

position as follows:

$$\partial_t h + \nabla \cdot (h \bar{\mathbf{v}}) = 0, \quad (1.4.27a)$$

$$\partial_t (h \bar{\mathbf{v}}) + \nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) + \nabla \left( \frac{1}{2} g h \right) + \nabla \cdot \mathbf{R} + \frac{1}{\rho} \left( \int_{z_b(\mathbf{x})}^{\xi(\mathbf{x}, t)} \nabla p_{NH} dz \right) = -g h \nabla z_b(\mathbf{x}), \quad (1.4.27b)$$

$$p_{NH}(\mathbf{x}, z, t) = \rho \int_z^{\xi(\mathbf{x}, t)} \left( \partial_t w + \nabla_{\mathbf{x}, z} \cdot (w \mathbf{V}) \right) dz. \quad (1.4.27c)$$

*Remark 1.4.3.* The non-hydrostatic pressure term contributes towards the model's dispersive effects and will be important for describing realistic water waves.

### 1.4.5 Non-dimensionalization of the depth and flow discharge formulation

At its current state, the water depth and flow discharge formulation (1.4.27) is too complicated to solve due to the non-hydrostatic pressure term and the contribution of the deviations from the average velocity (i.e.,  $\nabla \cdot \mathbf{R}$ ). To proceed, we would like to derive *simpler* shallow water models based on asymptotic expansions in terms of the  $h$  and  $\mathbf{q} := h \bar{\mathbf{v}}$  variables.

We now introduce non-dimensionalized quantities to study the asymptotic expansions of the equations. These quantities are based on the typical scales of problems concerning water waves. The quantities are as follows:  $h_0$ , the typical water depth;  $a_{surf}$ , the order (i.e., size) of variation of the free-surface deviation (i.e.,  $\eta(\mathbf{x}, t)$ );  $a_{bott}$ , the order of the bottom boundary (i.e., topography map) variation; and typical horizontal length-scale  $L$ . We also define the typical reference speed  $V_0 := \sqrt{g h_0}$ . The following non-dimensional parameters are paramount for the asymptotic expansion:

$$\begin{aligned} \mu &= \frac{h_0^2}{L^2} \rightarrow \text{shallowness parameter,} \\ \varepsilon &= \frac{a_{surf}}{h_0} \rightarrow \text{amplitude parameter,} \\ \beta &= \frac{a_{bott}}{h_0} \rightarrow \text{topography parameter.} \end{aligned}$$

*Remark 1.4.4* (Shallow water flows). The term *shallow water flows* refers to when the shallowness

parameter  $\mu$  is assumed to be small. In particular, we are interested in scenarios when  $\mu \in (0, 1)$ .

*Remark 1.4.5 (Depth-averaged models).* The models of interest in this thesis are the Saint-Venant Shallow Water Equations and the Serre–Green–Naghdi Equations. These models are often referred to as “depth-averaged” models since they evolve (in space and time) the water depth,  $h$ , and momentum quantity,  $\mathbf{q} := h\bar{\mathbf{v}}$ , where  $\bar{\mathbf{v}}$  is the depth-averaged horizontal velocity.

Using the typical problem quantities defined above, we now define the dimensionless variables for the flow/discharge formulation:

$$\begin{aligned} \mathbf{x}' &= \frac{\mathbf{x}}{L}, & z' &= \frac{z}{h_0}, & t' &= \frac{t}{T}, \\ \eta' &= \frac{\eta}{a_{surf}}, & \mathbf{q}' &= \frac{\mathbf{q}}{a_{surf}V_0}, & z'_b &= \frac{z_b}{a_{bott}}, \end{aligned}$$

where  $T = L/V_0$ . We also set  $h' = h/h_0$  and  $\bar{\mathbf{v}}' = \bar{\mathbf{v}}/(\frac{a_{surf}}{h_0}V_0) = \bar{\mathbf{v}}/(\varepsilon V_0)$ . The vertical scale quantity  $W_0$  is defined to be  $W_0 = \varepsilon \frac{h_0}{L} V_0$  and is found by comparing the scales in the incompressibility condition (1.4.2b). Note that since  $h := \eta + \xi_0 - z_b$ , we have that  $h' = \varepsilon\eta' + \frac{\xi_0}{h_0} - \beta z'_b$ . For simplicity, we assume that reference elevation at rest  $\xi_0$  is of the scale  $h_0$ . That is to say, we replace  $\xi_0$  by  $h_0$  so that  $h' = \varepsilon\eta' + 1 - \beta z'_b$ .

We now want to substitute the non-dimensionalized quantities above into the depth/discharge formulation (1.4.27). For completeness, we illustrate the non-dimensionalization for a few terms

in the formulation (1.4.27):

$$\partial_t \mathbf{h} \rightarrow \frac{\mathbf{h}_0}{T} \partial_{t'} \mathbf{h}' = \mathbf{h}_0 \frac{V_0}{L} \partial_{t'} \mathbf{h}',$$

$$\nabla \cdot \mathbf{q} \rightarrow \frac{a_{surf} V_0}{L} \nabla_{\mathbf{x}'} \cdot \mathbf{q}' = \varepsilon \mathbf{h}_0 \frac{V_0}{L} \nabla_{\mathbf{x}'} \cdot \mathbf{q}',$$

$$\partial_t \mathbf{q} \rightarrow \frac{1}{T} a_{surf} V_0 \partial_{t'} \mathbf{q}' = \mathbf{h}_0 \varepsilon \frac{V_0^2}{L} \partial_{t'} \mathbf{q}',$$

$$\nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) \rightarrow \frac{1}{L} (\varepsilon V_0) a_{surf} V_0 \nabla_{\mathbf{x}'} \cdot (\bar{\mathbf{v}}' \otimes \mathbf{q}') = \mathbf{h}_0 \varepsilon^2 \frac{V_0^2}{L} \nabla_{\mathbf{x}'} \cdot (\bar{\mathbf{v}}' \otimes \mathbf{q}'),$$

$$\nabla \left( \frac{1}{2} g h^2 \right) \rightarrow \frac{1}{L} \frac{V_0^2}{\mathbf{h}_0} \mathbf{h}_0^2 \nabla_{\mathbf{x}'} \cdot \left( \frac{1}{2} (\mathbf{h}')^2 \right) = \mathbf{h}_0 \frac{V_0^2}{L} \nabla_{\mathbf{x}'} \cdot \left( \frac{1}{2} (\mathbf{h}')^2 \right),$$

$$\partial_t w \rightarrow \frac{1}{T} \varepsilon \frac{\mathbf{h}_0}{L} V_0 \partial_{t'} w' = \varepsilon \frac{\mathbf{h}_0}{L} \frac{V_0^2}{L} \partial_{t'} w'.$$

Note that we used the fact that  $g = V_0^2/\mathbf{h}_0$  in the second-to-last expansion. Now, substituting the non-dimensionalized quantities above into the full depth/discharge formulation (1.4.27) yields (dropping the prime symbol ' notation for simplicity):

$$\partial_t \mathbf{h} + \varepsilon \nabla \cdot \mathbf{q} = 0, \tag{1.4.28a}$$

$$\varepsilon \partial_t (\mathbf{h} \bar{\mathbf{v}}) + \varepsilon^2 \nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) + \nabla \cdot \left( \frac{1}{2} \mathbf{h}^2 \right) + \varepsilon^2 \nabla \cdot \mathbf{R} + \mu \varepsilon \int_{\beta z_b}^{\varepsilon \eta + 1} \nabla p_{NH} dz = -\mathbf{h} \nabla (\beta z_b), \tag{1.4.28b}$$

$$p_{NH} = \int_z^{\varepsilon \eta + 1} \left( \partial_t w + \varepsilon \nabla_{\mathbf{x}, z} \cdot (w \mathbf{V}) \right) dz, \tag{1.4.28c}$$

where  $\mathbf{R} = \int_{\beta z_b}^{\varepsilon \eta + 1} (\mathbf{v}^* \otimes \mathbf{v}^*) dz$ . All the computations up to this point have been exact. That is to say, no approximations of the equations have been made. The goal moving forward is to introduce approximations in terms of the  $\varepsilon$  and  $\mu$  quantities to derive the shallow water models of interest. In particular, we derive the Saint-Venant Shallow Water Equations which we will show



is the first approximation to the water waves problem and then we derive the more complicated Serre–Green–Naghdi Equations.

#### 1.4.6 Inner structure of the velocity and pressure fields

We begin the approximation by first studying the structure of the velocity field inside the fluid domain. In particular, we will study the incompressibility (1.4.2b) and irrotational conditions (1.4.2c) (along with the impermeable bottom condition (1.4.10)) to find an expression for  $\mathbf{v}^*$  in terms of  $\bar{v}$ . These conditions (in dimensionless form) can be re-written as follows:

$$\nabla \cdot \mathbf{v} + \partial_z w = 0, \quad (1.4.29a)$$

$$\partial_z \mathbf{v} - \mu \nabla w = 0, \quad (1.4.29b)$$

$$\nabla^\top \cdot \mathbf{v} = 0, \quad (1.4.29c)$$

$$w_b - \nabla(\beta z_b) \cdot \mathbf{v}_b = 0, \quad (1.4.29d)$$

where  $\nabla^\top \cdot \mathbf{v} = \partial_y v - \partial_x u$ .

*Remark 1.4.6* (Representation of the irrotational condition). Note that the irrotational condition is represented by (1.4.29b) and (1.4.29c). Respectively, the equation (1.4.29b) is the horizontal component of the irrotational condition and equation (1.4.29c) is the vertical component of the irrotational condition. For completeness, we show how the horizontal component (1.4.29b) is found from the irrotational condition. Expanding each component of the (1.4.2c) yields

$$\mu \partial_y w - \partial_z v = 0, \quad \partial_z u - \mu \partial_x w = 0, \quad \partial_x v - \partial_y u = 0.$$

Notice that the first and second equations can be “switched” and written as follows:

$$\partial_z u - \mu \partial_x w = 0, \quad \partial_z v - \mu \partial_y w = 0.$$

Careful observation shows that this is equivalent to  $\partial_z \mathbf{v} - \mu \nabla w = 0$ .

We continue the analysis on the inner structure of the velocity field by vertically integrating (1.4.29a) (from  $\beta z_b$  to  $z$ ) and using the condition (1.4.29d) to see that:

$$\begin{aligned}
w - w_b &= - \int_{\beta z_b}^z \nabla \cdot (\bar{\mathbf{v}} + \mathbf{v}^*) dz \\
&= - \int_{\beta z_b}^z \nabla \cdot \bar{\mathbf{v}} dz - \int_{\beta z_b}^z \nabla \cdot \mathbf{v}^* dz, \\
&\implies \\
w &= \beta \nabla_{z_b} \cdot \mathbf{v}_b - (z - \beta z_b) \nabla \cdot \bar{\mathbf{v}} - \int_{\beta z_b}^z \nabla \cdot \mathbf{v}^* dz \\
&= \beta \nabla_{z_b} \cdot (\bar{\mathbf{v}} + \mathbf{v}_b^*) - (z - \beta z_b) \nabla \cdot \bar{\mathbf{v}} - \int_{\beta z_b}^z \nabla \cdot \mathbf{v}^* dz \\
&= \nabla \cdot (\beta z_b \bar{\mathbf{v}}) - z \nabla \cdot \bar{\mathbf{v}} + \underbrace{\beta \nabla_{z_b} \mathbf{v}_b^* - \int_{\beta z_b}^z \nabla \cdot \mathbf{v}^* dz}_{\text{constant}} \\
&= -\nabla \cdot \left( (z - \beta z_b) \bar{\mathbf{v}} \right) - \nabla \cdot \left( \int_{\beta z_b}^z \mathbf{v}^* dz \right).
\end{aligned}$$

Then, notice that the equation (1.4.29b) gives the following:

$$\begin{aligned}
\partial_z \mathbf{v} - \mu \nabla w &= 0, \\
&\implies \\
\partial_z (\bar{\mathbf{v}}(\mathbf{x}, t) + \mathbf{v}^*(\mathbf{x}, z, t)) - \mu \nabla w &= 0, \\
&\implies \\
\partial_z \mathbf{v}^* &= \mu \nabla w, \\
&\implies \\
\text{constant} - \mathbf{v}^* &= \int_z^{\varepsilon \eta + 1} \mu \nabla w dz.
\end{aligned}$$

Averaging the last equation over the water column gives  $\text{constant} = \mu \int_{\beta z_b}^{\varepsilon \eta + 1} \int_z^{\varepsilon \eta + 1} \nabla w d\tilde{z} dz$  where we used the fact that the mean-value over the water column of  $\mathbf{v}^*$  is zero. Thus, we can

equivalently write the previous equation as

$$\mathbf{v}^* = -\mu \left( \int_z^{\varepsilon\eta+1} w dz \right) + \mu \int_{\beta z_b}^{\varepsilon\eta+1} \int_z^{\varepsilon\eta+1} \nabla w d\tilde{z} dz = -\mu \left( \int_z^{\varepsilon\eta+1} w dz \right)^*.$$

This equation can be expanded even further by using the expression for  $w$ :

$$\begin{aligned} \mathbf{v}^* &= -\mu \left( \int_z^{\varepsilon\eta+1} \nabla w dz \right)^* \\ &= \mu \left( \int_z^{\varepsilon\eta+1} \nabla \{ \nabla \cdot (z - \beta z_b) \bar{\mathbf{v}} \} dz \right)^* + \mu \left( \int_z^{\varepsilon\eta+1} \nabla \{ \nabla \cdot \left( \int_{\beta z_b}^z \mathbf{v}^* dz \right) dz \} \right)^* \end{aligned}$$

We now follow [44] and introduce the following two operators:

$$\mathbb{T}\mathbf{w} = \int_z^{\varepsilon\eta+1} \nabla \{ \nabla \cdot \left( \int_{\beta z_b}^z \mathbf{w} dz \right) dz \}, \quad (1.4.30a)$$

$$\mathbb{T}^*\mathbf{w} = (\mathbb{T}\mathbf{w})^*, \quad (1.4.30b)$$

where  $\mathbf{w}$  is a (sufficiently smooth)  $\mathbb{R}^d$ -valued function defined on the fluid domain. Note that since  $\int_{\beta z_b}^z \bar{\mathbf{v}} dz = (z - \beta z_b) \bar{\mathbf{v}}$ , we can use the above defined operators to rewrite the equation for  $\mathbf{v}^*$  as follows:

$$(1 - \mu \mathbb{T}^*) \mathbf{v}^* = \mu \mathbb{T}^* \bar{\mathbf{v}}.$$

Applying the operator  $(1 + \mu \mathbb{T}^*)$  to both sides of the above equation, gives

$$\mathbf{v}^* = \mu \mathbb{T}^* \bar{\mathbf{v}} + \mathcal{O}(\mu^2), \quad (1.4.31)$$

which gives that the total horizontal fluid velocity is given by:

$$\mathbf{v} = \bar{\mathbf{v}} + \mu \mathbb{T}^* \bar{\mathbf{v}} + \mathcal{O}(\mu^2). \quad (1.4.32)$$

We have now found an approximate expression for the horizontal velocity in the fluid domain in terms of only the average horizontal velocity  $\bar{\mathbf{v}}$  up to the order of  $\mathcal{O}(\mu^2)$ . Since  $\bar{\mathbf{v}}$  does not depend

on  $z$ , we can compute  $\mathbf{v}$  explicitly. Note that

$$\begin{aligned}\mathbb{T}^*\bar{\mathbf{v}} &= \left( \int_z^{\varepsilon\eta+1} \nabla\{\nabla\cdot((z - \beta z_b)\bar{\mathbf{v}})\} dz \right)^* \\ &= \int_z^{\varepsilon\eta+1} \nabla\{\nabla\cdot((z - \beta z_b)\bar{\mathbf{v}})\} dz - \frac{1}{h} \int_{\beta z_b}^{\beta z_b+h} \int_z^{\varepsilon\eta+1} \nabla\{\nabla\cdot((\tilde{z} - \beta z_b)\bar{\mathbf{v}})\} d\tilde{z} dz.\end{aligned}$$

For simplification purposes, we expand each term separately:

$$\begin{aligned}\nabla\cdot((z - \beta z_b(\mathbf{x}))\bar{\mathbf{v}}) &= (z - \beta z_b(\mathbf{x}))\nabla\cdot\bar{\mathbf{v}} - \nabla(\beta z_b(\mathbf{x}))\cdot\bar{\mathbf{v}}, \\ &\implies \\ \nabla\{\nabla\cdot((z - \beta z_b(\mathbf{x}))\bar{\mathbf{v}})\} &= (z - \beta z_b(\mathbf{x}))\nabla(\nabla\cdot\bar{\mathbf{v}}) - \nabla(\beta z_b(\mathbf{x}))\nabla\cdot\bar{\mathbf{v}} \\ &\quad - \nabla\{\nabla(\beta z_b(\mathbf{x}))\cdot\bar{\mathbf{v}}\}, \\ &\implies \\ \int_z^{\varepsilon\eta+1} \left( \nabla\{\nabla\cdot((z - \beta z_b(\mathbf{x}))\bar{\mathbf{v}})\} \right) dz &= \frac{1}{2}(\mathbf{h}^2 - (z - z_b(\mathbf{x}))^2)\nabla(\nabla\cdot\bar{\mathbf{v}}) \\ &\quad - (\mathbf{h} + z_b(\mathbf{x}) - z)\nabla(\beta z_b(\mathbf{x}))\nabla\cdot\bar{\mathbf{v}} \\ &\quad - (\mathbf{h} + z_b(\mathbf{x}) - z)\nabla\{\nabla(\beta z_b(\mathbf{x}))\cdot\bar{\mathbf{v}}\}.\end{aligned}$$

Similarly, we have that

$$\begin{aligned}\frac{1}{h} \int_{z_b(\mathbf{x})}^{\varepsilon\eta+1} \int_z^{\varepsilon\eta+1} \left( \nabla\{\nabla\cdot((z - \beta z_b(\mathbf{x}))\bar{\mathbf{v}})\} \right) dz &= \frac{1}{3}\mathbf{h}^2\nabla(\nabla\cdot\bar{\mathbf{v}}) - \frac{1}{2}\mathbf{h}\nabla(\beta z_b(\mathbf{x}))\nabla\cdot\bar{\mathbf{v}} \\ &\quad - \frac{1}{2}\mathbf{h}\nabla\{\nabla(\beta z_b(\mathbf{x}))\cdot\bar{\mathbf{v}}\}.\end{aligned}$$

Then, continuing the expansion for  $\mathbb{T}^*\bar{\mathbf{v}}$ :

$$\begin{aligned}\mathbb{T}^*\bar{\mathbf{v}} &= -\frac{1}{2} \left( (z - \beta z_b(\mathbf{x}))^2 - \frac{1}{3}\mathbf{h}^2 \right) \nabla(\nabla\cdot\bar{\mathbf{v}}) \\ &\quad - (\beta z_b(\mathbf{x}) - z + \frac{1}{2}\mathbf{h}) \left( \nabla(\beta z_b(\mathbf{x}))\nabla\cdot\bar{\mathbf{v}} + \nabla\{\nabla(\beta z_b(\mathbf{x}))\cdot\bar{\mathbf{v}}\} \right)\end{aligned}$$

Substituting the above expression into the definition for the full horizontal fluid velocity yields:

$$\begin{aligned} \mathbf{v} = \bar{\mathbf{v}} - \mu \frac{1}{2} \left( (z - \beta z_b(\mathbf{x}))^2 - \frac{1}{3} h^2 \right) \nabla(\nabla \cdot \bar{\mathbf{v}}) \\ - \mu (\beta z_b(\mathbf{x}) - z + \frac{1}{2} h) \left( \nabla(\beta z_b(\mathbf{x})) \nabla \cdot \bar{\mathbf{v}} + \nabla \{ \nabla(\beta z_b(\mathbf{x})) \cdot \bar{\mathbf{v}} \} \right) + \mathcal{O}(\mu^2). \end{aligned} \quad (1.4.33)$$

Consequently, we have that the definition of the vertical component of the velocity is given by:

$$w = -\nabla \cdot \left( (z - \beta z_b(\mathbf{x})) \bar{\mathbf{v}} \right) - \mu \nabla \cdot \int_{\beta z_b(\mathbf{x})}^{h + z_b(\mathbf{x})} (\mathbb{T}^* \bar{\mathbf{v}}) dz + \mathcal{O}(\mu^2). \quad (1.4.34)$$

Note that the first term in the above equation can be expanded as  $-(z - \beta z_b(\mathbf{x})) \nabla \cdot \bar{\mathbf{v}} + \nabla(\beta z_b(\mathbf{x})) \cdot \bar{\mathbf{v}}$ .

We can now use the approximate expressions for  $\mathbf{v}$  (1.4.33) and  $w$  (1.4.34) to study the structure of the pressure field in the fluid domain. Let us recall the full Euler pressure (1.4.24) in dimensionless variables:

$$\begin{aligned} p &= (\varepsilon \eta + 1 - \beta z_b(\mathbf{x})) + \mu \varepsilon \int_{\beta z_b}^{\varepsilon \eta + 1} \nabla p_{NH} dz, \\ p_{NH} &= \int_z^{\varepsilon \eta + 1} \left( \partial_t w + \varepsilon \nabla \cdot (w \mathbf{v}) + \varepsilon \partial_z (w^2) \right) dz. \end{aligned}$$

Since the hydrostatic part of the pressure is simpler, we only look at the non-hydrostatic part of the pressure. The goal is to substitute  $\mathbf{v}$  (1.4.33) and  $w$  (1.4.34) into the above non-hydrostatic pressure.

*Remark 1.4.7* (Flat topography – A simpler case). Note that when the topography is flat (i.e.,  $z_b(\mathbf{x}) \equiv 0$ ), a direct computation (with the help of the `Mathematica` software [65]) yields that the non-hydrostatic pressure is given by:

$$p_{NH} = \frac{1}{2} (z^2 - h^2) \left( \partial_t (\nabla \cdot \bar{\mathbf{v}}) - \varepsilon (\nabla \cdot \bar{\mathbf{v}})^2 + \varepsilon \bar{\mathbf{v}} \cdot \nabla (\nabla \cdot \bar{\mathbf{v}}) \right) + \mathcal{O}(\mu). \quad (1.4.36)$$

This gives

$$\int_{\beta z_b}^{\varepsilon \eta + 1} \nabla p_{NH} dz = -\nabla \left\{ \frac{1}{3} h^3 \left( \partial_t(\nabla \cdot \bar{\mathbf{v}}) - \varepsilon(\nabla \cdot \bar{\mathbf{v}})^2 + \varepsilon \bar{\mathbf{v}} \cdot \nabla(\nabla \cdot \bar{\mathbf{v}}) \right) \right\} + \mathcal{O}(\mu).$$

We now consider the case when the topography is not flat. Another direct computation (using the Mathematica software [65]) gives that:

$$\begin{aligned} p_{NH} = & \frac{1}{2} ((z - \beta z_b(\mathbf{x}))^2 - h^2) \left( \partial_t(\nabla \cdot \bar{\mathbf{v}}) - \varepsilon(\nabla \cdot \bar{\mathbf{v}})^2 + \varepsilon \bar{\mathbf{v}} \cdot \nabla(\nabla \cdot \bar{\mathbf{v}}) \right) \\ & + (h + \beta z_b(\mathbf{x}) - z) \left( \partial_t(\nabla(\beta z_b) \cdot \bar{\mathbf{v}}) + \varepsilon \bar{\mathbf{v}} \cdot \nabla(\nabla(\beta z_b) \cdot \bar{\mathbf{v}}) \right) + \mathcal{O}(\mu), \end{aligned} \quad (1.4.37)$$

consequently:

$$\begin{aligned} \int_{\beta z_b}^{\varepsilon \eta + 1} \nabla p_{NH} dz = & \nabla \left\{ -\frac{1}{3} h^3 \left( \partial_t(\nabla \cdot \bar{\mathbf{v}}) - \varepsilon(\nabla \cdot \bar{\mathbf{v}})^2 + \varepsilon \bar{\mathbf{v}} \cdot \nabla(\nabla \cdot \bar{\mathbf{v}}) \right) \right. \\ & \left. + \frac{1}{2} h^2 \left( \partial_t(\nabla(\beta z_b) \cdot \bar{\mathbf{v}}) + \varepsilon \bar{\mathbf{v}} \cdot \nabla(\nabla(\beta z_b) \cdot \bar{\mathbf{v}}) \right) \right\} + h \left( -\frac{1}{2} h \left( \partial_t(\nabla \cdot \bar{\mathbf{v}}) - \varepsilon(\nabla \cdot \bar{\mathbf{v}})^2 + \varepsilon \bar{\mathbf{v}} \cdot \nabla(\nabla \cdot \bar{\mathbf{v}}) \right) \right. \\ & \left. + \left( \partial_t(\nabla(\beta z_b) \cdot \bar{\mathbf{v}}) + \varepsilon \bar{\mathbf{v}} \cdot \nabla(\nabla(\beta z_b) \cdot \bar{\mathbf{v}}) \right) \right) \nabla(\beta z_b(\mathbf{x})) + \mathcal{O}(\mu). \end{aligned} \quad (1.4.38)$$

*Remark 1.4.8* (The structure of the non-hydrostatic pressure contribution). Although the structure of the non-hydrostatic pressure contribution (1.4.38) looks complicated, it can be re-written with a “simpler” mathematical structure. We show how to do this in Section 2.3.3.

Up to this point, we have been able to approximately describe the water waves problem in terms of the water depth  $h$ , averaged horizontal velocity  $\bar{\mathbf{v}}$  and topography map  $z_b(\mathbf{x})$ . That is to say, we have explicitly removed the dependencies of the vertical quantities  $z$  (the vertical Cartesian coordinate) and  $w$  (the vertical fluid velocity). Although these quantities do not show up in the following models directly, there is still an indirect influence of the vertical fluid acceleration (in particular in the Boussinesq and Serre–Green–Naghdi models). We will now introduce further approximations in terms of the amplitude parameter  $\varepsilon$ , shallowness parameter  $\mu$  and topography

parameter  $\beta$  to derive our models of interest.

#### 1.4.7 First order approximation – Saint-Venant Shallow Water Equations

The first mathematical model of interest is the Saint-Venant Shallow Water Equations. This model is an approximation of order  $\mathcal{O}(\mu)$  of the non-dimensional water waves problem described above. That is to say, all terms proportional to  $\mu$  are dropped from the equations. Note that we make no assumption on the amplitude parameter  $\varepsilon$ . In particular, we have that:

$$\begin{aligned} \mathbf{v} &= \bar{\mathbf{v}} + \mathcal{O}(\mu), \\ \nabla \cdot \mathbf{R} &= \mathcal{O}(\mu^2), \\ \mu \int_{\beta z_b(\mathbf{x})}^{\varepsilon \eta + 1} \nabla p_{NH} &= \mathcal{O}(\mu). \end{aligned}$$

Thus for  $\mathbf{x} \in \mathbb{R}^d$  and  $t \geq 0$ , the Saint-Venant model in non-dimensional variables with topography effects is given as follows:

$$\partial_t \mathbf{h} + \varepsilon \nabla \cdot \mathbf{q} = 0, \tag{1.4.40a}$$

$$\partial_t \mathbf{q} + \varepsilon \nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) + \frac{1}{\varepsilon} \nabla \left( \frac{1}{2} \mathbf{h}^2 \right) = -\frac{1}{\varepsilon} \mathbf{h} \nabla (\beta z_b(\mathbf{x})). \tag{1.4.40b}$$

where we used  $\mathbf{q} := \mathbf{h} \bar{\mathbf{v}}$ . Notice that by dropping all the terms proportional to  $\mu$ , we have neglected all effects induced by the vertical fluid acceleration. In particular, the pressure in the Saint-Venant model (1.4.40) is only the hydrostatic pressure. Consequently, the Saint-Venant model can not model dispersive effects induced by the non-hydrostatic pressure. Since  $\nabla \cdot \mathbf{R} = \mathcal{O}(\mu^2)$ , the Saint-Venant model does not model “turbulent” effects.

*Remark 1.4.9* (Dimensional form of the Saint-Venant model). Substituting the dimensional variables into (1.4.40) yields:

$$\partial_t \mathbf{h} + \nabla \cdot (\mathbf{q}) = 0, \tag{1.4.41a}$$

$$\partial_t \mathbf{q} + \nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) + \nabla \left( \frac{1}{2} g \mathbf{h}^2 \right) = -g \mathbf{h} \nabla (z_b(\mathbf{x})). \tag{1.4.41b}$$

### 1.4.8 Weakly non-linear second order approximation – The Boussinesq model

Although the work in this thesis is not concerned with the Boussinesq model, we will briefly discuss it here for completeness. The Boussinesq model is a weakly non-linear second order approximation of the full water waves problem. That is to say, it is an order  $\mathcal{O}(\mu^2)$  approximation with the additional assumption on the amplitude parameter  $\varepsilon = \mathcal{O}(\mu)$  (i.e., size of the elevation variation). Another common assumption that is made is on the topography variations:  $\beta = \mathcal{O}(\mu)$ . Thus, we can drop any terms in the full water waves problem proportional to  $\mu^2$ ,  $\varepsilon\mu$  and  $\beta\mu$ . Thus, we have the following approximate quantities:

$$\begin{aligned} \mathbf{v} &= \bar{\mathbf{v}} + \mu \mathbb{T}^* \bar{\mathbf{v}} + \mathcal{O}(\mu^2), \\ \nabla \cdot \mathbf{R} &= \mathcal{O}(\mu^2), \\ \mu \int_{\beta z_b(\mathbf{x})}^{\varepsilon \eta + 1} \nabla p_{NH} &= \mu \nabla \left\{ -\frac{1}{3} h^3 \partial_t (\nabla \cdot \bar{\mathbf{v}}) \right\} + \mathcal{O}(\mu^2). \end{aligned}$$

Then, for  $\mathbf{x} \in \mathbb{R}^d$  and  $t \geq 0$ , the Boussinesq model in non-dimensional variables with weak non-linear and weak topography effects is given as follows:

$$\partial_t \mathbf{h} + \varepsilon \nabla \cdot (\mathbf{q}) = 0, \quad (1.4.43a)$$

$$\partial_t \mathbf{q} + \varepsilon \nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) + \frac{1}{\varepsilon} \nabla \cdot \left( \frac{1}{2} h^2 \right) - \mu \nabla \cdot \left( -\frac{1}{3} h^3 \partial_t (\nabla \cdot \bar{\mathbf{v}}) \right) = -\frac{1}{\varepsilon} h \nabla \cdot (\beta z_b(\mathbf{x})). \quad (1.4.43b)$$

*Remark 1.4.10 (Weakly dispersive).* Mathematical models of order  $\mathcal{O}(\mu^2)$  are often labeled as weakly dispersive in the literature.

*Remark 1.4.11 (Dimensional form of the Boussinesq model).* Substituting the dimensional variables into (1.4.43) yields:

$$\partial_t \mathbf{h} + \nabla \cdot (\mathbf{q}) = 0, \quad (1.4.44a)$$

$$\partial_t \mathbf{q} + \nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) + \nabla \cdot \left( \frac{1}{2} g h^2 \right) + \nabla \cdot \left( -\frac{1}{3} h^3 \partial_t (\nabla \cdot \bar{\mathbf{v}}) \right) = -g h \nabla \cdot (z_b(\mathbf{x})). \quad (1.4.44b)$$



### 1.4.9 Strongly non-linear second order approximation – Serre–Green–Naghdi model

We now discuss how to derive the Serre–Green–Naghdi model which is the focus of this thesis. The SGN model is a strongly non-linear second order approximation of the water waves problem (1.4.28). In particular, it is an order of  $\mathcal{O}(\mu^2)$  approximation of the free-surface Euler Equations (similar to the Boussinesq model), but without assumptions on the amplitude parameter  $\varepsilon$  and topography variation parameter  $\beta$ . Due to the lack of additional assumptions on these two parameters, the complexity of the model grows significantly (in particular in the expression for the non-hydrostatic pressure). We have the following approximate quantities for the SGN model:

$$\mathbf{v} = \bar{\mathbf{v}} + \mu \mathbb{T}^* \bar{\mathbf{v}} + \mathcal{O}(\mu^2),$$

$$\nabla \cdot \mathbf{R} = \mathcal{O}(\mu^2),$$

$$\begin{aligned} \mu \int_{\beta z_b}^{\varepsilon \eta + 1} \nabla p_{NH} dz &= \mu \nabla \left\{ -\frac{1}{3} h^3 \left( \partial_t (\nabla \cdot \bar{\mathbf{v}}) - \varepsilon (\nabla \cdot \bar{\mathbf{v}})^2 + \varepsilon \bar{\mathbf{v}} \cdot \nabla (\nabla \cdot \bar{\mathbf{v}}) \right) \right. \\ &+ \frac{1}{2} h^2 \left( \partial_t (\nabla (\beta z_b) \cdot \bar{\mathbf{v}}) + \varepsilon \bar{\mathbf{v}} \cdot \nabla (\nabla (\beta z_b) \cdot \bar{\mathbf{v}}) \right) \left. \right\} + \mu h \left( -\frac{1}{2} h \left( \partial_t (\nabla \cdot \bar{\mathbf{v}}) - \varepsilon (\nabla \cdot \bar{\mathbf{v}})^2 + \varepsilon \bar{\mathbf{v}} \cdot \nabla (\nabla \cdot \bar{\mathbf{v}}) \right) \right. \\ &\left. + \left( \partial_t (\nabla (\beta z_b) \cdot \bar{\mathbf{v}}) + \varepsilon \bar{\mathbf{v}} \cdot \nabla (\nabla (\beta z_b) \cdot \bar{\mathbf{v}}) \right) \right) \nabla (\beta z_b(\mathbf{x})) + \mathcal{O}(\mu^2). \end{aligned}$$

Notice that we have  $\nabla \cdot \mathbf{R} = \mathcal{O}(\mu^2)$  which is also the case for the Saint-Venant and Boussinesq models. That is to say, none of these models exhibit “turbulent” effects. Before writing the full model in non-dimensional form, let us define the following two quantities:

$$\mathcal{P}_\varepsilon = -h \left( \partial_t (\nabla \cdot \bar{\mathbf{v}}) - \varepsilon (\nabla \cdot \bar{\mathbf{v}})^2 + \varepsilon \bar{\mathbf{v}} \cdot \nabla (\nabla \cdot \bar{\mathbf{v}}) \right), \quad (1.4.46a)$$

$$\mathcal{L}_\varepsilon = \partial_t (\nabla (\beta z_b) \cdot \bar{\mathbf{v}}) + \varepsilon \bar{\mathbf{v}} \cdot \nabla (\nabla (\beta z_b) \cdot \bar{\mathbf{v}}). \quad (1.4.46b)$$

The expression  $\mathcal{P}_\varepsilon + \mathcal{L}_\varepsilon$  can be thought of as the influence of vertical fluid acceleration at the free-surface and  $\mathcal{L}_\varepsilon$  can be thought of as the vertical fluid acceleration at the bottom boundary (i.e.,

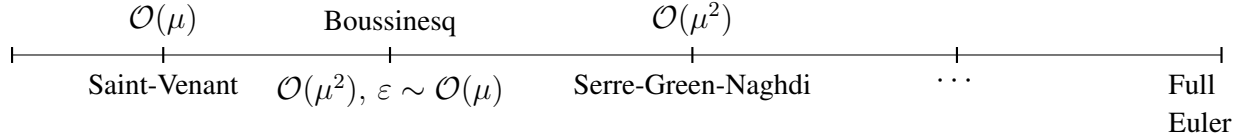


Figure 1.2: A representation of the physical accuracy of each model derived compared to the full free-surface Euler problem.

topography). Note that when the topography is flat, we have that  $\mathcal{L}_\varepsilon \equiv 0$ . For  $\mathbf{x} \in \mathbb{R}^d$  and  $t \geq 0$ , the full non-dimensional Serre–Green–Naghdi model can be written as follows:

$$\partial_t \mathbf{h} + \varepsilon \nabla \cdot (\mathbf{q}) = 0, \quad (1.4.47a)$$

$$\partial_t \mathbf{q} + \varepsilon \nabla \cdot (\bar{\mathbf{v}} \otimes \mathbf{q}) + \frac{1}{\varepsilon} \nabla \cdot \left( \frac{1}{2} \mathbf{h}^2 \right) - \mu \nabla \cdot \left( \mathbf{h}^2 \left( \frac{1}{3} \mathcal{P}_\varepsilon + \frac{1}{2} \mathcal{L}_\varepsilon \right) \right) = -\frac{1}{\varepsilon} \left( \mathbf{h} + \varepsilon \mu \frac{1}{2} \mathcal{P}_\varepsilon + \varepsilon \mu \mathcal{L}_\varepsilon \right) \nabla (\beta z_b(\mathbf{x})). \quad (1.4.47b)$$

In Figure 1.2, we give a representation of the physical accuracy of the models derived in this section. For an overview on open problems and discussion of boundary conditions regarding the models derived in this section, we refer the reader to [44].

## 2. SHALLOW WATER MODELS

### 2.1 Introduction

In this Chapter, we present the shallow water models of interest used for applications in inland flooding and coastal hydrodynamics: (i) the Saint-Venant shallow water equations and (ii) the Serre–Green–Naghdi model with topography effects for dispersive water waves. The goal of this chapter is to review the both mathematical models and discuss their mathematical and physical structures. In particular, we investigate steady-state solutions of the equations for developing robust numerical methods and verifying the accuracy of said methods. Although we review both mathematical models, the focus of this thesis is on dispersive water waves so we give particular attention to the Serre–Green–Naghdi equations. A fundamental question of this work is concerned with developing efficient methods for solving the SGN Equations and their use in applications in coastal hydrodynamics. Thus, it is important to understand the properties of the Serre–Green–Naghdi Equations and any potential drawbacks they exhibit. The discussion here lays the foundation for a robust and efficient hyperbolic relaxation technique shown in the following chapters. We then describe the mathematical formulation of the external physical sources that are considered in this work.

The chapter is organized as follows. In Section 2.2, we review the Saint-Venant model and give a brief discussion on its history and status in the literature. In Section 2.2.2, we describe some important properties of the Saint-Venant model and give a discussion on its drawbacks in Section 2.2.3. In Section 2.3, we give a general introduction and historical background to the Serre–Green–Naghdi equations. Then, in Section 2.3.1, we present the model problem for the Serre–Green–Naghdi equations along with a discussion of two important properties. Then in Section 2.3.3, we give two alternative representations of the equations. In Section 2.3.4, we discuss the challenges that are introduced when trying to solve the Serre Equations and its use for coastal hydrodynamics. The fundamental question of this thesis is presented in this section. In Section 2.3.5,

we derive important properties for the Serre model. In particular, a novel result concerning analytical solutions to the Serre–Green–Naghdi model is shown in Proposition 2.3.7. Finally, in Section 2.4 we discuss the external physical sources considered in this work.

## 2.2 The Saint-Venant model

The Saint-Venant equations are a hyperbolic system of conservation laws that were first derived in Saint-Venant [54] for one-dimensional flows in a channel. As shown in Section 1.4, the equations are an  $\mathcal{O}(\mu)$  approximation of the Free-Surface Euler Equations and neglect the effects of vertical acceleration of the fluid. The Saint-Venant equations have a wide range of applications from in-land flooding to atmospheric flows. As stated in Vreugdenhil [62], some of the early applications of digital computations were simulations done using the Saint-Venant equations for atmospheric flows in the 1940’s and for oceanographic flows in the 1950’s. Due to the general simplicity of the equations’ structure along with its mathematical properties (such as propagating shock waves), the Saint-Venant model has been studied and approximated extensively throughout the literature. For an overview of the Saint-Venant equations, we refer the reader to the books of Vreugdenhil [62] and Toro [61] and references therein.

### 2.2.1 The model problem

Let  $D$  be a polygonal domain in  $\mathbb{R}^d$ ,  $d \in \{1, 2\}$ , occupied by a body of water evolving in time under the action of gravity. Let  $\mathbf{u} = (h, \mathbf{q})^\top$  be the dependent variable, where  $h$  is the water height and  $\mathbf{q}$  the momentum vector (also known as the flow discharge). Let  $z(\mathbf{x})$  be the given topography (or bathymetry) map. The Saint-Venant model with topography effects can be written as follows:

$$\partial_t h + \nabla \cdot (h\mathbf{v}) = 0, \tag{2.2.1a}$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q} + p(\mathbf{u})\mathbb{I}_d) = -gh\nabla z \tag{2.2.1b}$$

where the hydrostatic pressure is defined by:

$$p(\mathbf{u}) = \frac{1}{2}gh^2. \tag{2.2.2}$$

Here,  $\mathbf{v}$  is the (depth-averaged) velocity vector and is defined such that  $\mathbf{q} := h\mathbf{v}$ .

### 2.2.2 Saint-Venant properties

In this section, we briefly discuss a few important properties of the Saint-Venant model that will be used later in this work.

**Proposition 2.2.1** (Saint-Venant Energy). *Let  $\mathbf{u}$  be a smooth solution to the Saint-Venant model (2.2.1).*

*Then, the Saint-Venant model admits a conservation equation for the total energy:  $\partial_t E(\mathbf{u}) + \nabla \cdot \mathbf{F}(\mathbf{u}) = 0$  with the energy functional and energy flux defined by:*

$$E(\mathbf{u}) := \frac{1}{2}gh^2 + gzh + \frac{1}{2}h\mathbf{v}^2, \quad (2.2.3a)$$

$$\mathbf{F}(\mathbf{u}) := \mathbf{v}(E(\mathbf{u}) + p(\mathbf{u})). \quad (2.2.3b)$$

We define the admissible set for the Saint-Venant model be:

$$\mathcal{A} = \{\mathbf{u} := (h, \mathbf{q})^\top \mid h > 0\}. \quad (2.2.4)$$

This (convex) set can be thought of as the set which contains physically relevant solutions to the Saint-Venant model.

**Proposition 2.2.2** (Lake-at-rest for the Saint-Venant Equations). *If the flow is at rest ( $\mathbf{q} \equiv \mathbf{0}$  for  $\mathbf{x} \in D$ ), the lake-at-rest steady-state problem for the Saint-Venant Equations (2.2.1) is given by:*

$$gh\nabla(h + z) = \mathbf{0}, \quad \mathbf{x} \in D. \quad (2.2.5a)$$

*Proof.* Assume the flow is at rest,  $\mathbf{q} := h\mathbf{v} \equiv \mathbf{0}$  (and  $\mathbf{v} \equiv \mathbf{0}$ ). Since we are interested in steady-state solutions (i.e., independent of time), the Saint-Venant system (2.2.1) reduces to the following

partial differential equation:

$$\begin{aligned}\nabla \cdot (\mathbf{h}\mathbf{v}) &= 0, \\ \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla \left( \frac{1}{2}gh^2 \right) &= -gh\nabla z.\end{aligned}$$

Using  $\mathbf{h}\mathbf{v} \equiv \mathbf{0}$ , the second equation above reduces to  $\nabla \left( \frac{1}{2}gh^2 \right) + gh\nabla z = \mathbf{0}$ . Expanding the first gradient term yields the lake-at-rest steady-state equation:

$$gh\nabla(\mathbf{h} + z) = \mathbf{0}.$$

□

### 2.2.3 Physical drawbacks of the Saint-Venant model

As discussed above, the Saint-Venant model has many applications. However, as shown in Chapter 1, the Saint-Venant model is only an  $\mathcal{O}(\mu)$  approximation of the free-surface Euler Equations for water waves and neglects the vertical acceleration effects. Consequently, the Saint-Venant model lacks dispersive properties. More precisely, its linear phase speed is given by  $c_p := \frac{\sigma}{k} = \sqrt{gh_0}$ . Here,  $\sigma$  is the wave frequency of a plane-wave solution and  $k$  is the wave number (or wave frequency in the spatial domain). Thus, the Saint-Venant model is a less accurate representation of the full water waves problem and can not properly propagate smooth solutions such as solitary waves or periodic waves. We illustrate this visually in Figure 2.1. It is clear from the figure that the Saint-Venant model produces a solution with shocks (i.e., breaking waves) instead of propagating the smooth waves. The lack of dispersive effects of the Saint-Venant model is a major drawback especially when studying wave dynamics in the near-shore region where phenomena such as wave shoaling and wave diffraction are important. The solution for correcting this drawback is by keeping the  $\mathcal{O}(\mu^2)$  terms (as shown in Chapter 1) such as is done with the Serre–Green–Naghdi equations.

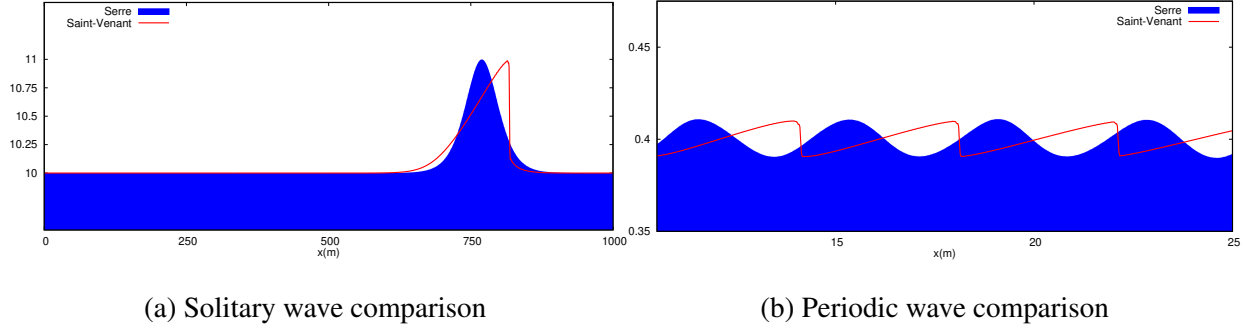


Figure 2.1: Comparison of smooth solutions with Hyperbolic Serre model and Saint-Venant model.

### 2.3 The Serre–Green–Naghdi Equations

The Serre–Green–Naghdi Equations are a non-linear system of partial differential equations used for modeling dispersive water waves and have a wide-range of applications in coastal hydrodynamics. The Serre Equations were first derived in one spatial-dimension in the work of Serre [57, Eq. (22), p. 860] as a generalization to the Saint-Venant model. The equations were then re-discovered in Su and Gardner [58, Eq. (A14)]. In both of these references, smooth exact solutions in the form of solitary waves were derived for the model (a key property for simulating tsunami waves). The equations were then derived in two spatial dimensions for a flat bottom in Green et al. [26, Eq. (4.12)–(4.15)] and then in Green and Naghdi [25, Eq. (4.27)–(4.30)] for variable topography. The equations were also derived with topography effects in Seabra-Santos et al. [56, Eq. (14)]. The authors of [56] supplemented their mathematical findings with physical experiments involving the propagation of solitary waves over variable topography. In the literature, the nomenclature of the Serre–Green–Naghdi Equations is often interchanged with the Serre equations, Green–Naghdi Equations, or the fully non-linear Boussinesq equations. The latter is a consequence of the Serre–Green–Naghdi model being a weakly dispersive model like the Boussinesq Equations (i.e.,  $\mathcal{O}(\mu^2)$  approximation to the free-surface Euler equations) but fully non-linear (as shown in §1.4). For the sake of brevity, we interchange the Serre–Green–Naghdi nomenclature with just the Serre model.

### 2.3.1 The model problem

Let  $D$  be a polygonal domain in  $\mathbb{R}^d$ ,  $d \in \{1, 2\}$ , occupied by a body of water evolving in time under the action of gravity. Let  $\mathbf{u} = (\mathbf{h}, \mathbf{q})^\top$  be the dependent variable of the system, where  $\mathbf{h}$  is the water height and  $\mathbf{q}$  the momentum vector (also known as the flow discharge). Then, the Serre model as first introduced in Serre [57] and extended in Green and Naghdi [25] and Seabra-Santos et al. [56] can be written as follows:

$$\partial_t \mathbf{h} + \nabla \cdot (\mathbf{h} \mathbf{v}) = 0, \quad (2.3.1a)$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q} + p(\mathbf{u}) \mathbb{I}_d) = -r(\mathbf{u}) \nabla z, \quad (2.3.1b)$$

with the pressure  $p(\mathbf{u})$  and the source  $r(\mathbf{u})$  defined by

$$p(\mathbf{u}) := \frac{1}{2} g h^2 + h^2 \left( \frac{1}{3} \ddot{\mathbf{h}} + \frac{1}{2} \dot{\mathbf{k}} \right), \quad (2.3.2a)$$

$$r(\mathbf{u}) = g h + h \left( \frac{1}{2} \ddot{\mathbf{h}} + \dot{\mathbf{k}} \right). \quad (2.3.2b)$$

Here, the “dot” objects are defined by:

$$\dot{\mathbf{h}} := \partial_t \mathbf{h} + \mathbf{v} \cdot \nabla \mathbf{h}, \quad (2.3.3a)$$

$$\ddot{\mathbf{h}} := \partial_t \dot{\mathbf{h}} + \mathbf{v} \cdot \nabla \dot{\mathbf{h}}, \quad (2.3.3b)$$

$$\dot{\mathbf{k}} := \partial_t (\mathbf{v} \cdot \nabla z) + \mathbf{v} \cdot \nabla (\mathbf{v} \cdot \nabla z). \quad (2.3.3c)$$

The quantity  $\mathbf{v}$  is the (depth-averaged) velocity vector field and is defined such that  $\mathbf{q} := \mathbf{h} \mathbf{v}$ . The quantity  $\dot{\mathbf{k}}$  can be interpreted as the vertical acceleration of the fluid particles induced by the topography while  $\ddot{\mathbf{h}} + \dot{\mathbf{k}}$  can be interpreted as the vertical acceleration of the fluid particles at the free surface ([56, 6]).



### 2.3.2 Admissible set and Lake-at-rest

The admissible set (or invariant domain) for the Serre Equations is defined by:

$$\mathcal{A} = \{\mathbf{u} := (\mathbf{h}, \mathbf{q})^\top \mid \mathbf{h} > 0.\} \quad (2.3.4)$$

The admissible set is physically meaningful since the water depth  $\mathbf{h}$  can be thought of as the distance between the free-surface elevation and topography map. A numerical method that preserves the invariant domain at the discrete level can be classified as invariant domain preserving. Since the admissible set contains only the positivity condition on the water depth, a numerical method that preserves this property can also be called *positivity-preserving*. These methods are robust and preserve the structure of their continuous counterparts without the requirement of tune-able parameters.

**Proposition 2.3.1** (Lake-at-rest for the Serre Equations). *If the flow is at rest ( $\mathbf{q} \equiv \mathbf{0}$  for  $\mathbf{x} \in D$ ), then the lake-at-rest steady-state problem for the Serre Equations (2.3.1) is given by:*

$$g\mathbf{h}\nabla(\mathbf{h} + z) = \mathbf{0}, \quad \mathbf{x} \in D. \quad (2.3.5)$$

*Proof.* Assume the flow is at rest  $\mathbf{q} \equiv \mathbf{0}$  (and  $\mathbf{v} \equiv \mathbf{0}$ ). Since we are interested in steady-state solutions (i.e., no dependence on the time variable  $t$ ), we have that  $\dot{\mathbf{h}} = \mathbf{v} \cdot \nabla \mathbf{h}$ . However, by assumption,  $\mathbf{v} = \mathbf{0}$  so  $\dot{\mathbf{h}} = 0$ . Consequently,  $\ddot{\mathbf{h}} = 0$ . The rest of the proof is exactly the same as Proof 2.2.2. □

*Remark 2.3.2* (Lake-at-rest solution). The lake-at-rest problem (2.3.5) has one of two solutions: (i)  $\mathbf{h} = 0$ ; or (ii)  $\mathbf{h} + z = \text{constant}$ . The physical interpretation of the solution  $\mathbf{h} + z = \text{constant}$  is as follows: If the flow is at rest, then the water elevation must remain constant.

We call numerical methods that preserve the lake-at-rest property at the discrete level *well-balanced*.

### 2.3.3 Mathematical reinterpretation of dispersive terms

In this section, we give different representation of the Serre model that are equivalent. Let  $D_t$  be the linear, total derivative operator defined by

$$D_t \phi := (\partial_t + \mathbf{v} \cdot \nabla) \phi, \quad (2.3.6)$$

where  $\phi$  is an arbitrary smooth function.

**Proposition 2.3.3** (First alternative representation). *A first alternative representation of the Serre Equations is given by:*

$$\partial_t h + \nabla \cdot (h \mathbf{v}) = 0, \quad (2.3.7a)$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla \left( \frac{1}{2} g h^2 + \frac{1}{3} h^2 D_t^2 \left( h + \frac{3}{2} z \right) \right) = - \left( g h + D_t^2 \left( \frac{1}{2} h + z \right) \right) \nabla z \quad (2.3.7b)$$

*Proof.* Assume  $\mathbf{u} := (h, \mathbf{q})^T$  is a smooth solution to the Serre model. Let  $z(\mathbf{x})$  be a given, smooth topography profile. Notice that only the non-hydrostatic pressure and source are represented differently from (2.3.1)–(2.3.2). Expanding the  $D_t^2$  term in the pressure gives:

$$\begin{aligned} \frac{1}{3} h^2 D_t^2 \left( h + \frac{3}{2} z \right) &= \frac{1}{3} h^2 D_t \left( D_t \left( h + \frac{3}{2} z \right) \right) \\ &= \frac{1}{3} h^2 D_t \left( \dot{h} + D_t \left( \frac{3}{2} z \right) \right) \\ &= \frac{1}{3} h^2 D_t \left( \dot{h} + \frac{3}{2} \mathbf{v} \cdot \nabla z \right) \\ &= \frac{1}{3} h^2 \left( \ddot{h} + \frac{3}{2} \dot{\mathbf{k}} \right) \\ &= h^2 \left( \frac{1}{3} \ddot{h} + \frac{1}{2} \dot{\mathbf{k}} \right) \end{aligned}$$

which is the expression in (2.3.2). A similar expansion can be performed for the topography source term. □

The representation shown in Proposition 2.3.3 was first shown in Green and Naghdi [25] (after

some algebraic manipulation of equations (4.27)–(4.31) therein) and can also be seen in Dellar and Salmon [14, Eq. (18)].

**Proposition 2.3.4** (Second alternative representation). *A second alternative representation of the Serre Equations is given by:*

$$\partial_t \mathbf{h} + \nabla \cdot (\mathbf{h} \mathbf{v}) = 0, \quad (2.3.8a)$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla \left( \frac{1}{2} g \mathbf{h}^2 + \bar{p}(\mathbf{u}) \right) = - (g \mathbf{h} + \bar{r}(\mathbf{u})) \nabla z, \quad (2.3.8b)$$

$$(2.3.8c)$$

where

$$\bar{p}(\mathbf{u}) = -\frac{1}{3} \mathbf{h}^2 \left( \mathbf{h} \left( \partial_t (\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2 + \mathbf{v} \cdot \nabla (\nabla \cdot \mathbf{v}) \right) - \frac{3}{2} \dot{\mathbf{k}} \right), \quad (2.3.9a)$$

$$\bar{r}(\mathbf{u}) = -\frac{1}{2} \mathbf{h} \left( \partial_t (\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2 + \mathbf{v} \cdot \nabla (\nabla \cdot \mathbf{v}) \right) + \dot{\mathbf{k}}. \quad (2.3.9b)$$

*Proof.* Assume  $\mathbf{u} := (\mathbf{h}, \mathbf{q})^\top$  is a smooth solution to the Serre model. Let  $z(\mathbf{x})$  be a given, smooth topography profile. In the above formulation, we have expanded the quantity  $\ddot{\mathbf{h}}$  using the mass

conservation equation (2.3.1a) as follows:

$$\begin{aligned}
\ddot{h} &= D_t(\dot{h}) \\
&= D_t \left( \underbrace{\partial_t h + \mathbf{v} \cdot \nabla h}_{\text{M.C (2.3.1a)}} \right) \\
&= D_t(-h \nabla \cdot \mathbf{v}) \\
&= -(\partial_t(h \nabla \cdot \mathbf{v}) + \mathbf{v} \cdot \nabla(h \nabla \cdot \mathbf{v})) \\
&= - \left( (\nabla \cdot \mathbf{v}) \underbrace{\partial_t h}_{\text{M.C (2.3.1a)}} + h \partial_t(\nabla \cdot \mathbf{v}) + \mathbf{v} \cdot ((\nabla \cdot \mathbf{v}) \nabla h + h \nabla(\nabla \cdot \mathbf{v})) \right) \\
&= - \left( (\nabla \cdot \mathbf{v})(-\mathbf{v} \cdot \nabla h - h(\nabla \cdot \mathbf{v})) + h \partial_t(\nabla \cdot \mathbf{v}) + \mathbf{v} \cdot ((\nabla \cdot \mathbf{v}) \nabla h + h \nabla(\nabla \cdot \mathbf{v})) \right) \\
&= -h \left( \partial_t(\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2 + \mathbf{v} \cdot \nabla(\nabla \cdot \mathbf{v}) \right)
\end{aligned}$$

which is the expression in (2.3.2). A similar expansion can be performed for the topography source term. □

The representation shown in Proposition 2.3.4 was first presented in Serre [57] and Su and Gardner [58] for one spatial dimension and no topography effects. This representation shows up naturally in the in the derivation of the Serre–Green–Naghdi model as shown in Section 1.4 (see equation (1.4.38)). Notice that we can re-write the momentum equation (2.3.8b) as an evolution equation for the velocity  $\mathbf{v}$  as follows:

$$\partial_t \mathbf{v} + (\mathbf{v} \cdot \nabla) \mathbf{v} + \frac{1}{h} \nabla \left( \frac{1}{2} g h^2 + \bar{p}(\mathbf{u}) \right) = -\frac{1}{h} (g h + \bar{r}(\mathbf{u})) \nabla z,$$

with

$$\begin{aligned}
\bar{p}(\mathbf{u}) &= -\frac{1}{3} h^2 \left( h \left( \partial_t(\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2 + \mathbf{v} \cdot \nabla(\nabla \cdot \mathbf{v}) \right) - \frac{3}{2} (\partial_t \mathbf{v} \cdot \nabla z + \mathbf{v} \cdot \nabla(\mathbf{v} \cdot \nabla z)) \right), \\
\bar{r}(\mathbf{u}) &= -\frac{1}{2} h \left( \partial_t(\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2 + \mathbf{v} \cdot \nabla(\nabla \cdot \mathbf{v}) \right) + \partial_t \mathbf{v} \cdot \nabla z + \mathbf{v} \cdot \nabla(\mathbf{v} \cdot \nabla z).
\end{aligned}$$

**Proposition 2.3.5** (Third alternative representation). *A third alternative representation of the Serre Equations is given by:*

$$\partial_t \mathbf{h} + \nabla \cdot (\mathbf{h} \mathbf{v}) = 0, \quad (2.3.11a)$$

$$\begin{aligned} \mathbf{h} \left( 1 + \frac{1}{\mathbf{h}} \mathbb{T} \right) [\partial_t \mathbf{v}] + \mathbf{h} (\mathbf{v} \cdot \nabla) \mathbf{v} + g \mathbf{h} \nabla (\mathbf{h} + z) \\ - \frac{1}{3} \nabla \{ \mathbf{h}^3 [(\mathbf{v} \cdot \nabla)(\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2] \} + \mathbf{h} \mathbb{Q}[\mathbf{v}] = 0. \end{aligned} \quad (2.3.11b)$$

Here, the operators  $\frac{1}{\mathbf{h}} \mathbb{T}[\cdot]$  and  $\mathbb{Q}[\cdot]$  are defined by:

$$\begin{aligned} \frac{1}{\mathbf{h}} \mathbb{T}[\mathbf{w}] &= -\frac{1}{3\mathbf{h}} \nabla (\mathbf{h}^3 (\nabla \cdot \mathbf{w})) - \frac{1}{2} \mathbf{h} (\nabla \cdot \mathbf{w}) \nabla z + \frac{1}{2\mathbf{h}} \nabla (\mathbf{h}^2 \nabla z \cdot \mathbf{w}) + (\nabla z \cdot \mathbf{w}) \nabla z, \\ \mathbb{Q}[\mathbf{w}] &= \frac{1}{2\mathbf{h}} \nabla (\mathbf{h}^2 (\mathbf{w} \cdot \nabla)^2 z) - \frac{1}{2} \mathbf{h} ((\mathbf{w} \cdot \nabla)(\nabla \cdot \mathbf{w}) - (\nabla \cdot \mathbf{w})^2) \nabla z + ((\mathbf{w} \cdot \nabla)^2 z) \nabla z. \end{aligned}$$

with  $\mathbf{w}$  is a sufficiently smooth vector-valued function.

*Proof.* Let  $\mathbf{u} := (\mathbf{h}, \mathbf{q})^\top$  be a smooth solution to the system (2.3.11). Since there is no change to the mass conservation equation, we want to show that (2.3.11b) is equivalent to (2.3.1b)–(2.3.2). To simplify the computations, we will give certain terms an integer marker to signify that they will be combined. Let us recall two identities that will be useful in the following analysis. When the solution  $\mathbf{u}$  and topographic map  $z(\mathbf{x})$  are smooth, we have that:

$$\begin{aligned} \dot{\mathbf{k}} &= \partial_t (\mathbf{v} \cdot \nabla z) + (\mathbf{v} \cdot \nabla)^2 z, \\ \mathbf{h}^2 \ddot{\mathbf{h}} &= -\mathbf{h}^3 (\partial_t (\nabla \cdot \mathbf{v}) + (\mathbf{v} \cdot \nabla)(\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2), \end{aligned}$$

where we used that  $(\mathbf{v} \cdot \nabla)^2 z = \mathbf{v} \cdot \nabla (\mathbf{v} \cdot \nabla z)$ . Using the second identity, we see that:

$$-\frac{1}{3} \nabla \{ \mathbf{h}^3 [(\mathbf{v} \cdot \nabla)(\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2] \} = \nabla \left( \frac{1}{3} \mathbf{h}^2 \ddot{\mathbf{h}} \right) + \underbrace{\nabla \left( \frac{1}{3} \mathbf{h}^3 \partial_t (\nabla \cdot \mathbf{v}) \right)}_{(1)}.$$

Note that  $gh\nabla(h+z) = \nabla(\frac{1}{2}gh^2) + gh\nabla z$ . Let us now focus on the first two terms of (2.3.11b).

We have that

$$\begin{aligned} h\left(1 + \frac{1}{h}\mathbb{T}\right) [\partial_t \mathbf{v}] + (\mathbf{v} \cdot \nabla) \mathbf{v} &= \underbrace{h(\partial_t \mathbf{v} + (\mathbf{v} \cdot \nabla) \mathbf{v})}_{\text{Prop. 3.2.1}} + h\frac{1}{h}\mathbb{T}[\partial_t \mathbf{v}] \\ &= \partial_t \mathbf{q} + \nabla \cdot (\mathbf{q} \otimes \mathbf{v}) + h\frac{1}{h}\mathbb{T}[\partial_t \mathbf{v}]. \end{aligned}$$

Expanding the last term separately yields:

$$h\frac{1}{h}\mathbb{T}[\partial_t \mathbf{v}] = \underbrace{-\frac{1}{3}\nabla(h^3(\nabla \cdot \partial_t \mathbf{v}))}_{(1)} - \underbrace{\frac{1}{2}h^2(\nabla \cdot \partial_t \mathbf{v})\nabla z}_{(2)} + \underbrace{\frac{1}{2}\nabla(h^2\nabla z \cdot \partial_t \mathbf{v})}_{(3)} + \underbrace{h(\nabla z \cdot \partial_t \mathbf{v})\nabla z}_{(4)},$$

We now expand each term the topography term  $h\mathbb{Q}[\mathbf{v}]$ :

$$\begin{aligned} \frac{1}{2}\nabla(h^2(\mathbf{v} \cdot \nabla)^2 z) &= \nabla(h^2\frac{1}{2}\dot{\mathbf{k}}) - \nabla(\frac{1}{2}h^2\partial_t(\mathbf{v} \cdot \nabla z)), \\ &= \nabla(h^2\frac{1}{2}\dot{\mathbf{k}}) - \underbrace{\nabla(\frac{1}{2}h^2\partial_t \mathbf{v} \cdot \nabla z)}_{(3)}, \\ -\frac{1}{2}h^2((\mathbf{v} \cdot \nabla)(\nabla \cdot \mathbf{v}) - (\nabla \cdot \mathbf{v})^2)\nabla z &= \frac{1}{2}h\ddot{h}\nabla z + \underbrace{\frac{1}{2}h^2(\nabla \cdot \partial_t \mathbf{v})\nabla z}_{(2)}, \\ h((\mathbf{v} \cdot \nabla)^2 z)\nabla z &= h\dot{\mathbf{k}}\nabla z - h\partial_t(\mathbf{v} \cdot \nabla z)\nabla z \\ &= h\dot{\mathbf{k}}\nabla z - \underbrace{h(\partial_t \mathbf{v} \cdot \nabla z)\nabla z}_{(4)}. \end{aligned}$$

We now combine all the expanded terms. Note that the terms with the integer markers will all cancel. Thus, we have that:

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{q} \otimes \mathbf{v}) + \nabla(\frac{1}{2}gh^2 + h^2(\frac{1}{3}\ddot{h} + \frac{1}{2}\dot{\mathbf{k}})) = -(gh + h(\frac{1}{2}\ddot{h} + \dot{\mathbf{k}}))\nabla z,$$

which is (2.3.1b)–(2.3.2). □

This representation was first introduced in Lannes and Bonneton [46, Eq. (26)] for the non-dimensional Serre–Green–Naghdi Equations and can be seen in Bonneton et al. [10, Eq. (13)] and Marche [50, Eq. (3)].

### 2.3.4 Challenges

In this section, we briefly discuss some of the challenges involved in the general scientific process for solving the Serre–Green–Naghdi equations. In particular, we discuss the time step restriction that arises in the numerical approximation of the equations, the process of verification and validation and the status of boundary conditions. The key scientific question of this thesis is presented in this section. For an overview on open problems regarding the Serre–Green–Naghdi equations, we refer the reader to Lannes [44].

#### 2.3.4.1 Dispersive time step restriction

As shown in Chapter 1, the Serre–Green–Naghdi equations are a better approximation of the free-surface Euler Equations than the Saint-Venant model and are able to accurately represent wave transformations in near-shore dynamics when dispersive and non-linear effects are strong. However, this physical accuracy comes at a price. A major drawback of the dispersive Serre model from a numerical perspective is that it involves third-order derivatives in space (this can be seen directly in (2.3.8b)–(2.3.9a)). More specifically, due to the presence of the  $\mathcal{O}(\mu^2)$  dispersive terms, the pressure mapping  $\mathbf{u} \mapsto p(\mathbf{u})$  defined in (2.3.2) is not a function, but a second-order differential operator in space and time. This second-order differential operator produces third-order spatial derivatives in the momentum equations (2.3.1b) and rule out any approximation technique that is explicit in time, since this would require that the time step  $\tau$  satisfy:

$$\tau \sim \mathcal{O}(h^3)V^{-1}L^{-2},$$

where  $h$  is the mesh-size,  $V$  is a characteristic wave speed scale, and  $L$  is a characteristic length scale. We call this condition the dispersive time step restriction. There are currently two popular classes of techniques for addressing this difficulty. The first one is based on Strang’s operator splitting and combines explicit and implicit time stepping, see for instance Bonneton et al. [10], Samii and Dawson [55], Duran and Marche [15]. Another approach consists of reinterpreting the dispersive system as a constrained first-order system and then relaxing the constraints to allow for using the typical hyperbolic CFL condition (see Chapters 3 and 4).

The goal of this work is concerned with answering the following scientific question:

*Is it possible to construct an explicit, well-balanced numerical method for solving the Serre–Green–Naghdi equations that is at least second-order in space and time and invariant-domain preserving under the hyperbolic CFL condition?*

#### 2.3.4.2 Verification and validation

Verification and validation are paramount processes in scientific computing. Analytical solutions are needed to verify the accuracy of numerical codes and can be useful for understanding the structure of mathematical models. At this time, the availability of analytical solutions for the Serre model with non-trivial topography is non-existent. The process of model validation is useful for investigating the physical limitations of the mathematical model of interest (such as the Serre Equations) and is often done by comparing numerical results with experimental data. However, the availability of reliable experimental data can be scarce.

#### 2.3.4.3 Boundary conditions

Most simulations involving water waves have large-time scales especially those modeling large domains such as the ocean. Since we can only perform computations on finite domains, these large-time scales can be an issue if we want to avoid wave reflections at the boundary. One can try to limit these reflections through wave absorption zones introduced in the domain, but these techniques are often limited and require more computational resources. Another approach is to carefully enforce the “correct” boundary conditions so that information flows out of the boundary.



However, this approach is highly non-trivial and is PDE dependent.

In the following chapters, we work towards answering the scientific equation 2.3.4.1 along with addressing the above the difficulties.

### 2.3.5 Properties

In this section, we derive important properties of the Serre–Green–Naghdi model (2.3.1)–(2.3.2). In particular, we derive a conservation equation for the total energy of the system. We then derive analytical solutions for the Serre model with topography effects which are used for verifying the accuracy of numerical methods. Finally, we derive the dispersion relation for the linearized model. The techniques and results shown in this section will be useful when studying the hyperbolic relaxation technique introduced in Chapter 4.

#### 2.3.5.1 Conservation of energy

We first derive the conservation of energy equation for the Serre system (2.3.1)–(2.3.2). This is done as follows. We first derive a conservation equation for the kinetic energy (see (2.3.13)). Then, we derive a conservation equation for the potential energy (see (2.3.14)). Finally, we combine the two equations and obtain a single conservation equation for the total energy of the system. This technique is a standard approach for deriving energy results for similar mathematical models such as the Saint-Venant model. The main result of this section is given in Proposition 2.3.6.

Let  $\mathbf{u}(\mathbf{x}, t) := (h(\mathbf{x}, t), \mathbf{q}(\mathbf{x}, t))^T$  be a smooth solution to the Serre model (2.3.1)–(2.3.2). We first derive the conservation equation for the kinetic energy by applying the operator  $\mathbf{v} \cdot (\cdot)$  to (2.3.1b) and combining the results. For completeness and clarity, we simplify each term individually and then combine the terms at the end (note that when convenient, we interchange the notation  $\dot{h} =$

$D_t \mathbf{h}$ ). Expanding  $\mathbf{v} \cdot (\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}))$  yields:

$$\begin{aligned}
\mathbf{v} \cdot \partial_t \mathbf{q} &= \mathbf{v} \cdot \partial_t (\mathbf{h} \mathbf{v}) \\
&= \|\mathbf{v}\|^2 \partial_t \mathbf{h} + \mathbf{h} \left( \frac{1}{2} \partial_t \|\mathbf{v}\|^2 \right) \\
&= \|\mathbf{v}\|^2 \partial_t \mathbf{h} + \frac{1}{2} \partial_t (\mathbf{h} \|\mathbf{v}\|^2) - \frac{1}{2} \|\mathbf{v}\|^2 \partial_t \mathbf{h} \\
&= \partial_t \left( \frac{1}{2} \mathbf{h} \|\mathbf{v}\|^2 \right) - \frac{1}{2} \|\mathbf{v}\|^2 \nabla \cdot (\mathbf{h} \mathbf{v}) \\
&= \partial_t \left( \frac{1}{2} \mathbf{h} \|\mathbf{v}\|^2 \right) - \nabla \cdot (\mathbf{v} \frac{1}{2} \mathbf{h} \|\mathbf{v}\|^2) + \mathbf{h} \mathbf{v} \cdot \nabla \left( \frac{1}{2} \|\mathbf{v}\|^2 \right),
\end{aligned}$$

$$\begin{aligned}
\mathbf{v} \cdot (\nabla \cdot (\mathbf{v} \otimes \mathbf{q})) &= \mathbf{v} \cdot (\nabla \cdot (\mathbf{v} \otimes \mathbf{h} \mathbf{v})), \\
&= \mathbf{v} \cdot \left( \mathbf{h} \mathbf{v} (\nabla \cdot \mathbf{v}) + (\mathbf{v} \cdot \nabla \mathbf{h}) \mathbf{v} + \mathbf{h} \nabla \left( \frac{1}{2} \|\mathbf{v}\|^2 \right) \right) \\
&= \|\mathbf{v}\|^2 (\mathbf{h} \nabla \cdot \mathbf{v}) + \|\mathbf{v}\|^2 (\mathbf{v} \cdot \nabla \mathbf{h}) + \mathbf{h} \mathbf{v} \cdot \left( \nabla \frac{1}{2} \|\mathbf{v}\|^2 \right) \\
&= \|\mathbf{v}\|^2 \nabla \cdot (\mathbf{h} \mathbf{v}) + \mathbf{h} \mathbf{v} \cdot \left( \nabla \frac{1}{2} \|\mathbf{v}\|^2 \right).
\end{aligned}$$

Note that:

$$\begin{aligned}
-\nabla \cdot \left( \mathbf{v} \frac{1}{2} \mathbf{h} \|\mathbf{v}\|^2 \right) + \mathbf{h} \mathbf{v} \cdot \nabla (\|\mathbf{v}\|^2) + \|\mathbf{v}\|^2 \nabla \cdot (\mathbf{h} \mathbf{v}) &= -\nabla \cdot \left( \mathbf{v} \frac{1}{2} \mathbf{h} \|\mathbf{v}\|^2 \right) + \nabla \cdot (\mathbf{h} \mathbf{v} \|\mathbf{v}\|^2) \\
&= \nabla \cdot \left( \mathbf{v} \frac{1}{2} \mathbf{h} \|\mathbf{v}\|^2 \right).
\end{aligned}$$

Expanding  $\mathbf{v} \cdot \nabla (p(\mathbf{u})) = \mathbf{v} \cdot \nabla \left( \frac{1}{2} g h^2 + h^2 \left( \frac{1}{3} \dot{\mathbf{h}} + \frac{1}{2} \dot{\mathbf{k}} \right) \right)$  yields:

$$\mathbf{v} \cdot \nabla \left( \frac{1}{2} g h^2 \right) = \mathbf{v} \cdot \nabla \left( \frac{1}{2} g h^2 \right),$$

$$\begin{aligned}
\mathbf{v} \cdot \nabla \left( \frac{1}{3} h^2 \ddot{h} \right) &= \nabla \cdot \left( \mathbf{v} \frac{1}{3} h^2 \ddot{h} \right) - \underbrace{\frac{1}{3} h^2 \ddot{h} (\nabla \cdot \mathbf{v})}_{= -\dot{h}/h} \\
&= \nabla \cdot \left( \mathbf{v} \frac{1}{3} h^2 \ddot{h} \right) + \frac{1}{3} h \ddot{h} \dot{h} \\
&= \nabla \cdot \left( \mathbf{v} \frac{1}{3} h^2 \ddot{h} \right) + \underbrace{\frac{1}{3} h D_t \left( \frac{1}{2} (\dot{h})^2 \right)}_{\text{Prop. 3.2.1}} \\
&= \nabla \cdot \left( \mathbf{v} \frac{1}{3} h^2 \ddot{h} \right) + \partial_t \left( \frac{1}{6} h (\dot{h})^2 \right) + \nabla \cdot \left( \mathbf{v} \frac{1}{6} h (\dot{h})^2 \right),
\end{aligned}$$

$$\begin{aligned}
\mathbf{v} \cdot \nabla \left( \frac{1}{2} h^2 \dot{\mathbf{k}} \right) &= \nabla \cdot \left( \mathbf{v} \frac{1}{2} h^2 \dot{\mathbf{k}} \right) - \frac{1}{2} h^2 \dot{\mathbf{k}} (\nabla \cdot \mathbf{v}) \\
&= \nabla \cdot \left( \mathbf{v} \frac{1}{2} h^2 \dot{\mathbf{k}} \right) + \underbrace{\frac{1}{2} h \dot{\mathbf{k}} \dot{h}}_{*}
\end{aligned}$$

Expanding  $\mathbf{v} \cdot (r(\mathbf{u}) \nabla z) = \mathbf{v} \cdot (gh + h(\frac{1}{2} \ddot{h} + \dot{\mathbf{k}}) \nabla z)$  yields:

$$\mathbf{v} \cdot (gh \nabla z) = gh \mathbf{v} \cdot \nabla z,$$

$$\begin{aligned}
\mathbf{v} \cdot \left( \frac{1}{2} h \ddot{h} \nabla z + h \dot{\mathbf{k}} \nabla z \right) &= \frac{1}{2} h \ddot{h} (\mathbf{v} \cdot \nabla z) + h (\mathbf{v} \cdot \nabla z) \dot{\mathbf{k}} \\
&= \underbrace{\frac{1}{2} h \ddot{h} (\mathbf{v} \cdot \nabla z)}_{*} + \underbrace{h D_t \left( \frac{1}{2} (\mathbf{v} \cdot \nabla z)^2 \right)}_{*}.
\end{aligned}$$

The terms with the ‘\*’ label can be simplified as follows:

$$\begin{aligned}
\frac{1}{2} h \dot{\mathbf{k}} \dot{h} + \frac{1}{2} h \ddot{h} (\mathbf{v} \cdot \nabla z) + h D_t \left( \frac{1}{2} (\mathbf{v} \cdot \nabla z)^2 \right) &= \underbrace{h D_t \left( \frac{1}{2} (\mathbf{v} \cdot \nabla z)^2 \right)}_{\text{Prop. 3.2.1}} + \underbrace{h D_t \left( \frac{1}{2} \dot{h} (\mathbf{v} \cdot \nabla z) \right)}_{\text{Prop. 3.2.1}} \\
&= \partial_t \left( \frac{1}{2} h (\mathbf{v} \cdot \nabla z)^2 \right) + \nabla \cdot \left( \mathbf{v} \frac{1}{2} h (\mathbf{v} \cdot \nabla z)^2 \right) \\
&\quad + \partial_t \left( \frac{1}{2} h \dot{h} (\mathbf{v} \cdot \nabla z) \right) + \nabla \cdot \left( \mathbf{v} \frac{1}{2} h \dot{h} (\mathbf{v} \cdot \nabla z) \right).
\end{aligned}$$

Then, the conservation equation for the kinetic energy can be written as follows:

$$\begin{aligned} & \partial_t \left( \frac{1}{2} h \|\mathbf{v}\|^2 + \frac{1}{6} h \left( \dot{h} + \frac{3}{2} (\mathbf{v} \cdot \nabla z) \right)^2 + \frac{3}{4} (\mathbf{v} \cdot \nabla z)^2 \right) \\ & + \nabla \cdot \left( \mathbf{v} \left( \frac{1}{2} h \|\mathbf{v}\|^2 + \frac{1}{6} h \left( \dot{h} + \frac{3}{2} (\mathbf{v} \cdot \nabla z) \right)^2 + \frac{3}{4} (\mathbf{v} \cdot \nabla z)^2 \right) \right) + \mathbf{v} \cdot \nabla \left( \frac{1}{2} g h^2 \right) + h \mathbf{v} \cdot \nabla (g z) = 0. \end{aligned} \quad (2.3.13)$$

To find a conservation equation for the potential energy, we multiply the mass conservation (2.3.1a) by  $g(h + z)$  and combine the results. Similarly as above, we expand each term individually. Expanding  $g(h + z)(\partial_t h + \nabla \cdot (h \mathbf{v}))$  yields

$$\begin{aligned} g(h + z) \partial_t h &= g h \partial_t h + g z \partial_t h \\ &= \partial_t \left( \frac{1}{2} g h^2 + g h z \right), \end{aligned}$$

$$\begin{aligned} g(h + z) \nabla \cdot (h \mathbf{v}) &= g h (\mathbf{v} \cdot \nabla h + h \nabla \cdot \mathbf{v}) + g z \nabla \cdot (h \mathbf{v}) \\ &= g \mathbf{v} \cdot h \nabla h + g h^2 \nabla \cdot \mathbf{v} + g z \nabla \cdot (h \mathbf{v}) \\ &= \mathbf{v} \cdot \nabla \left( \frac{1}{2} g h^2 \right) + \frac{1}{2} g h^2 \nabla \cdot \mathbf{v} + \frac{1}{2} g h^2 \nabla \cdot \mathbf{v} + g z \nabla \cdot (h \mathbf{v}) \\ &= \nabla \cdot \left( \frac{1}{2} g h^2 \mathbf{v} \right) + \frac{1}{2} g h^2 (\nabla \cdot \mathbf{v}) + g z \nabla \cdot (h \mathbf{v}). \end{aligned}$$

These terms are combined as follows to form the conservation equation for the potential energy:

$$\partial_t \left( \frac{1}{2} g h^2 + g h z \right) + \nabla \cdot \left( \mathbf{v} \frac{1}{2} g h^2 \right) + \frac{1}{2} g h^2 (\nabla \cdot \mathbf{v}) + g z \nabla \cdot (h \mathbf{v}) = 0. \quad (2.3.14)$$

Summing the equations for the kinetic energy (2.3.13) and potential energy (2.3.14) yields the following result.

**Proposition 2.3.6.** *Let  $u$  be a smooth solution to (2.3.1)–(2.3.2), then the following holds true:*

$\partial_t \mathcal{E}(\mathbf{u}) + \nabla \cdot (\mathcal{F}(\mathbf{u})) = 0$ , with

$$\mathcal{E}(\mathbf{u}) := \frac{1}{2}gh^2 + gzh + \frac{1}{2}h\mathbf{v}^2 + \frac{1}{6}h \left( \left( \dot{h} + \frac{3}{2}(\mathbf{v} \cdot \nabla z) \right)^2 + \frac{3}{4}(\mathbf{v} \cdot \nabla z)^2 \right), \quad (2.3.15a)$$

$$\mathcal{F}(\mathbf{u}) := \mathbf{v}(\mathcal{E}(\mathbf{u}) + p(\mathbf{u})). \quad (2.3.15b)$$

### 2.3.5.2 Analytical steady-state solution with topography

In this section, we derive a family of analytical steady-state solutions for the one-dimensional dispersive Serre equations. The solutions are used to validate the accuracy of the proposed relaxed model. The main interest of this exact solution is to help verify the accuracy numerical codes for the approximation of the dispersive Serre model with topography. These solutions are used in Chapter 6 to verify the accuracy of the proposed hyperbolic relaxation technique (introduced in Chapter 4) and associated numerical method (introduced in Chapter 5).

We restrict ourselves to one spatial dimension and assume that the solution to (2.3.1)–(2.3.2) is time-independent and smooth. Then, the steady state problem for the Serre equations with topography for  $\mathbf{u}(x) := (h(x), q(x))^T$  is given by:

$$\partial_x(hv) = 0, \quad (2.3.16a)$$

$$\partial_x \left( hv^2 + \frac{1}{2}gh^2 + h^2 \left( \frac{1}{3}\ddot{h} + \frac{1}{2}\dot{k} \right) \right) = - \left( gh + h \left( \frac{1}{2}\ddot{h} + \dot{k} \right) \right) \partial_x z, \quad (2.3.16b)$$

By (2.3.16a), we see that the flow discharge  $q := hv$  is constant so we let  $q \equiv q_0 \in \mathbb{R}$ . Since the solution variable  $\mathbf{u}$  is time-independent and  $v = \frac{q_0}{h}$ , we have that

$$\dot{h} = v\partial_x h = \frac{q_0}{h}\partial_x h,$$

$$\ddot{h} = v\partial_x(v\partial_x h) = \frac{q_0^2}{h}\partial_x \left( \frac{1}{h}\partial_x h \right),$$

$$\dot{k} = v\partial_x(v\partial_x z) = \frac{q_0^2}{h}\partial_x \left( \frac{1}{h}\partial_x z \right).$$

We now re-write (2.3.16b) as follows:

$$\partial_x \left( \frac{q_0^2}{h} + \frac{1}{2}gh^2 + h \left( \frac{1}{3}q_0^2 \partial_x \left( \frac{1}{h} \partial_x h \right) + \frac{1}{2}q_0^2 \partial_x \left( \frac{1}{h} \partial_x z \right) \right) \right) = - \left( gh + \frac{1}{2}q_0^2 \partial_x \left( \frac{1}{h} \partial_x h \right) + q_0^2 \partial_x \left( \frac{1}{h} \partial_x z \right) \right) \partial_x z.$$

Then, dividing the above equation by  $gh$  and re-arranging the terms using the Mathematica software [65] yields:

$$\partial_x \left( \frac{1}{gh} \left( gh^2 + ghz + \frac{1}{2} \frac{q_0^2}{h} + \frac{1}{3} q_0^2 \partial_{xx} h + \frac{1}{2} \partial_{xx} z - \frac{1}{6} \frac{q_0^2}{h} (\partial_x h)^2 + \frac{1}{2} \frac{q_0^2}{h} (\partial_x z)^2 \right) \right) = 0. \quad (2.3.17)$$

The above can be re-written even further by noticing that since  $q \equiv q_0$ , we have that:

$$h^2 \left( \frac{1}{3} \ddot{h} + \frac{1}{2} \dot{k} \right) + \frac{1}{6} h \left( (\dot{h} + \frac{3}{2} (v \partial_x z)^2 + \frac{3}{4} (v \partial_x z)^2) \right) = \frac{1}{3} q_0^2 \partial_{xx} h + \frac{1}{2} \partial_{xx} z - \frac{1}{6} \frac{q_0^2}{h} (\partial_x h)^2 + \frac{1}{2} \frac{q_0^2}{h} (\partial_x z)^2.$$

Then, multiplying (2.3.17) by  $gq_0$  gives:

$$\partial_x \left( \frac{q_0}{h} \left( gh^2 + ghz + \frac{1}{2} \frac{q_0^2}{h} + h^2 \left( \frac{1}{3} \ddot{h} + \frac{1}{2} \dot{k} \right) + \frac{1}{6} h \left( (\dot{h} + \frac{3}{2} (v \partial_x z)^2 + \frac{3}{4} (v \partial_x z)^2) \right) \right) \right) = 0. \quad (2.3.18)$$

Upon further inspection, we see this equation is equivalent to the steady-state problem for the energy equation introduced in Proposition 2.3.6:

$$\partial_x \left( v \left( \mathcal{E}(\mathbf{u}(x)) + p(\mathbf{u}(x)) \right) \right) = 0. \quad (2.3.19)$$

To find a steady-state solution for the Serre model, we need to solve (2.3.19). This is summarized in the following assertion which is the main result from this section.

**Proposition 2.3.7.** *Let  $q_0 \in \mathbb{R}$ ,  $a, r \in \mathbb{R}_+$  and let the bathymetry profile be defined by  $z(x) := -\frac{1}{2} \frac{a}{(\cosh(rx))^2}$ . Then  $h(x) = h_0 \left( 1 + \frac{a}{(\cosh(rx))^2} \right)$  with the constant discharge  $q_0$  is a steady state solution to (2.3.1)–(2.3.2) if*

$$q := \pm \sqrt{\frac{(1+a)gh_0^3}{2}}, \quad r := \frac{1}{h_0} \sqrt{\frac{3a}{1+a}}. \quad (2.3.20)$$

*Proof.* Let  $\mathbf{u}(x) := (\mathbf{h}(x), q(x))^\top$  be a smooth solution to the Serre model. Notice that the discharge  $q \equiv q_0$  is constant since the solution does not depend on time. As shown above, a steady-state solution can be found by solving the steady-state problem of the energy equation in Proposition 2.3.6. Solving this equation yields a Bernoulli-like relation for the dispersive Serre model (2.3.1)–(2.3.2). More precisely, from (2.3.19), we infer that  $\partial_x(\mathcal{F}(\mathbf{u})) = 0$ , which implies  $\mathcal{F}(\mathbf{u}(x)) = C_{\text{Ber}}q_0g$  where  $C_{\text{Ber}}$  is the Bernoulli constant. We look for a stationary wave with the following structure  $\mathbf{h}(x) = \mathbf{h}_0(1 + \frac{a}{(\cosh(rx))^2})$  and posit that the topography is of the form  $z(x) = \lambda(\mathbf{h}(x) - \mathbf{h}_0)$ . The problem now consists of finding relations between the parameters  $a$ ,  $r$ ,  $\mathbf{h}_0$ ,  $g$ , and  $\lambda$  so that the condition  $g^{-1}q_0^{-1}\mathcal{F}(\mathbf{u}(x)) = C_{\text{Ber}}$  is satisfied, i.e.,

$$\mathbf{h}(1 + \lambda) + \frac{q_0^2}{2g\mathbf{h}^2} - \frac{q_0^2}{6g\mathbf{h}^2}(1 - 3\lambda^2)(\partial_x\mathbf{h})^2 + \frac{q_0^2}{3g\mathbf{h}}(1 + \frac{3}{2}\lambda)\partial_{xx}\mathbf{h} = C_{\text{Ber}} + \lambda\mathbf{h}_0,$$

where we used  $v = \frac{q_0}{\mathbf{h}}$ . By taking the limit of this identity for  $|x| \rightarrow \infty$ , we find that  $C_{\text{Ber}} = \mathbf{h}_0 + \frac{q_0^2}{2g\mathbf{h}_0^2}$ . After inserting the ansatz  $\mathbf{h}(x) = \mathbf{h}_0(1 + \frac{a}{(\cosh(rx))^2})$  into the above identity, we find that the following must hold true for all  $x \in \mathbb{R}$ :

$$\begin{aligned} & \left( (1 + \lambda)g\mathbf{h}_0^3 + (\lambda + \frac{2}{3})2r^2q_0^2\mathbf{h}_0^2 - q_0^2 \right) \cosh(rx)^4 \\ & + \left( (1 + \lambda)2ag\mathbf{h}_0^3 + ((\lambda^2 + \lambda + \frac{1}{3})a - \frac{3}{2}\lambda - 1)2r^2q_0^2\mathbf{h}_0^2 - \frac{1}{2}aq_0^2 \right) \cosh(rx)^2 \\ & + \left( (1 + \lambda)ag\mathbf{h}_0 - 2(\lambda^2 + \frac{3}{2}\lambda + \frac{2}{3})r^2q_0^2 \right) a\mathbf{h}_0^2 = 0. \end{aligned}$$

This is equivalent to asserting that following nonlinear system of equations has a solution (i.e., setting the coefficients of the polynomial for  $\cosh(rx)$  equal to 0):

$$\begin{aligned} & (1 + \lambda)ag\mathbf{h}_0 - 2(\lambda^2 + \frac{3}{2}\lambda + \frac{2}{3})r^2q_0^2 = 0, \\ & (1 + \lambda)2ag\mathbf{h}_0^3 + ((\lambda^2 + \lambda + \frac{1}{3})a - \frac{3}{2}\lambda - 1)2r^2q_0^2\mathbf{h}_0^2 - \frac{1}{2}aq_0^2 = 0, \\ & (1 + \lambda)g\mathbf{h}_0^3 + (\lambda + \frac{2}{3})2r^2q_0^2\mathbf{h}_0^2 - q_0^2 = 0. \end{aligned}$$

The only nontrivial solution to the above system of equations is

$$\lambda = -\frac{1}{2}, \quad q_0 = \pm \sqrt{\frac{(1+a)gh_0^3}{2}}, \quad r = \frac{1}{h_0} \sqrt{\frac{3a}{1+a}}.$$

□

In Figure 2.2, we illustrate this steady-state solution with the values  $g = 9.81 \text{ ms}^{-2}$ ,  $h_0 = 1 \text{ m}$ ,  $a = 0.2 \text{ m}$  and  $q_0 \approx \sqrt{5.886} \text{ m}^2 \text{ s}^{-1}$ . The solid blue profile represents the wave elevation  $h(x) + z(x)$ .

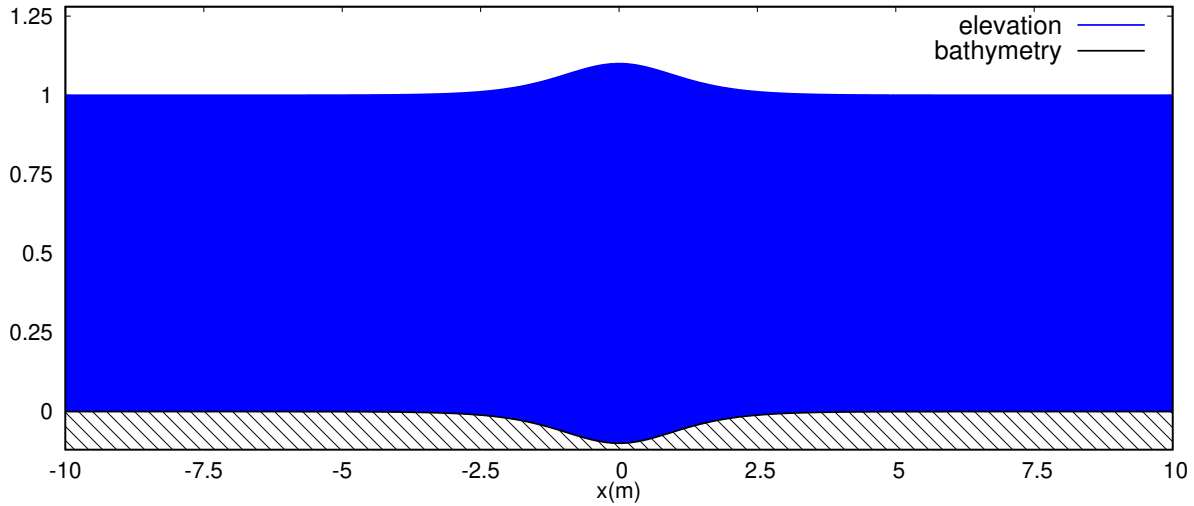


Figure 2.2: 1D Steady state solution for the Serre model

### 2.3.5.3 Dispersion relation

In this section, we derive the linear dispersion relation for the Serre model (2.3.1). The results shown here will be useful for analyzing the dispersive properties of a hyperbolic relaxed model of the Serre equations in Chapter 4. We note that the derivation shown here is not new and is found in the literature.



For simplicity, we restrict ourselves to one spatial dimension and assume the topography is flat. To derive the linear dispersion relation, we linearize the Serre equations about a rest state  $\mathbf{u}_0 := (H_0, 0)^\top$ . We consider solutions of the Serre equations of the form  $\mathbf{u} = \mathbf{u}_0 + \bar{\mathbf{u}}(x, t)$  where  $\bar{\mathbf{u}}(x, t)$  is a small perturbation of the rest state. Substituting this solution into the equations and dropping the product of the perturbation terms (since they are small) yields the following linear system of partial differential equations:

$$\partial_t \bar{h} + H_0 \partial_x \bar{u} = 0, \quad (2.3.21a)$$

$$\partial_t \bar{u} + g \partial_x \bar{h} + \frac{1}{3} H_0 \partial_{xxt} \bar{h} = 0. \quad (2.3.21b)$$

Assume now that  $\bar{\mathbf{u}}(x, t) = \mathbf{u}_A \exp(i(kx - \sigma t))$  is a plane-wave solution to (2.3.21) where  $\mathbf{u}_A$  is the wave amplitude,  $k$  is the wave number and  $\sigma$  is the wave frequency. Here  $i := \sqrt{-1}$  is the imaginary unit number. A direct substitution of the plane wave solution into (2.3.21) yields the following linear system for  $(h_A, u_A)^\top$ :

$$\begin{pmatrix} -\sigma & H_0 k \\ gk - \frac{1}{3} H_0 k \sigma^2 & -\sigma \end{pmatrix} \begin{pmatrix} h_A \\ u_A \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

Assuming that the solution to this linear system is non-trivial, we get the following equation for  $\sigma$ :

$$\sigma^2 \left( 1 + \frac{1}{3} H_0^2 k^2 \right) - g H_0 k^2 = 0. \quad (2.3.22)$$

Solving the above equation for the (squared) phase velocity squared  $(c_p^S)^2 := (\frac{\sigma}{k})^2$  yields:

$$(c_p^S)^2 = \frac{g H_0}{1 + \frac{1}{3} H_0^2 k^2}. \quad (2.3.23)$$

In Figure 2.3, we plot the squared phase velocity  $(c_p^S)^2$  as a function of the wave number  $k$ .

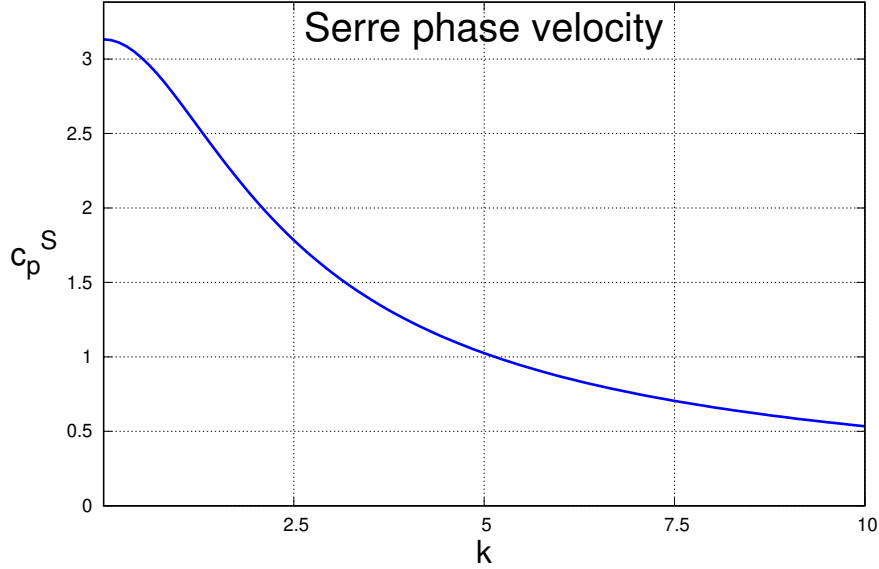


Figure 2.3: A plot of the Serre phase velocity (2.3.23) as a function of the wave number  $k$  with  $H_0 = 1$  m.

## 2.4 External physical sources

In this section we describe the mathematical formulation of the external physical sources used for applications in coastal hydrodynamics and inland flooding.

### 2.4.1 Preliminaries

For notation purposes, consider the condensed form of either the Saint-Venant model or the Serre–Green–Naghdi system with sources:

$$\partial_t \mathbf{u} + \nabla \cdot \mathbb{f}(\mathbf{u}) = \mathbf{R}(\mathbf{u}, \nabla z) + \mathbf{S}(\mathbf{u}).$$

Here, the vector  $\mathbf{R}(\mathbf{u}, \nabla z)$  is the natural topography source that arises in the formulation of the system of interest and is henceforth referred to as the PDE source. The quantity  $\mathbf{S}(\mathbf{u})$  represents the accumulation of the external physical sources described below. This term is henceforth referred to as the external source. The mathematical formulation of the source terms will be presented in model-independent fashion.

## 2.4.2 Gauckler-Manning friction

Let  $\mathbb{e}_q$  denote the characteristic vector for the momentum equations. We account for loss of discharge due to friction effects by adopting the Gauckler-Manning's friction law. The friction source is defined as follows:

$$\mathbf{S}_F(\mathbf{u}) := -gn^2\mathbf{h}^{-\gamma}\mathbf{q}\|\mathbf{v}\|_{\ell^2\mathbb{e}_q}. \quad (2.4.1)$$

The parameter  $n$  is the Gauckler-Manning's roughness coefficient and has units  $\text{m}^{\frac{\gamma-2}{2}}\text{s}$ . We take  $\gamma = \frac{4}{3}$  in the computations reported below in Chapter 6.

To illustrate the effects of the Gauckler-Manning friction term, we consider the propagation of a solitary wave over a flat-bottom with different values of the Manning's roughness coefficient. We set the computational domain to  $D = (0, 1000 \text{ m})$  and initialize the solitary wave with the profiles given by (6.3.1) at  $x_0 = 250 \text{ m}$ . The final time is set to  $T = 50 \text{ s}$ . In Figure 2.4, we give the final profiles for the water depth with roughness coefficients of  $n = \{0, 0.1, 0.2, 0.3, 0.4\}\text{m}^{-\frac{1}{3}}\text{s}$ . We see that as the roughness coefficient increases, the solitary wave slows down and loses amplitude as well. For an overview on selecting the Manning roughness coefficient  $n$  for natural channels and flood plains, we refer the reader to Arcement and Schneider [2] (and references therein) where this is discussed.

## 2.4.3 Wave generation and absorption

In applications that involve the propagation of periodic waves, a common technique in the literature is to introduce relaxation zones in a numerical wave tank to smoothly generate and absorb waves (see: Zhang et al. [68], Madsen et al. [48] and references therein). These generation and absorption zones are introduced as source terms in the equations.

For simplicity, let us assume we have a rectangular computational domain  $D$ . Assume that we want to generate uni-directional waves perpendicular to the inflow boundary so that wave profiles only depend on the  $x$ -direction (by convention  $x$  is the first Cartesian coordinate of the position vector  $\mathbf{x}$ ). Let  $\mathbf{u}_{\text{wave}}(\mathbf{x}, t)$  denote the theoretical wave profiles for each conserved variable. Denoting by  $h_{\text{wave}}$  the water depth component of  $\mathbf{u}_{\text{wave}}$ , we assume that  $h_{\text{wave}}(\mathbf{x}, t) \geq 0$  for all  $\mathbf{x} \in D$

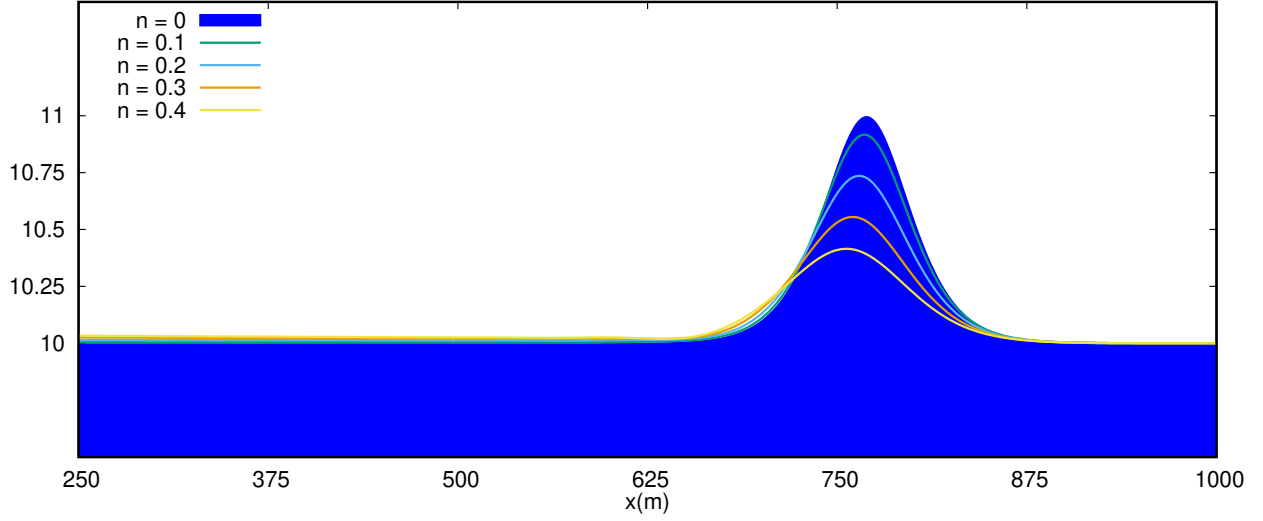


Figure 2.4: Propagation of solitary wave profile over flat-bottom with roughness coefficients of  $n = \{0, 0.1, 0.2, 0.3, 0.4\} \text{m}^{-\frac{1}{3}} \text{s}$ .

and all  $t > 0$ . To generate periodic waves through a relaxation zone, we introduce the following source:

$$\mathbf{S}_G(\mathbf{u}) := -\frac{\sqrt{gH_0}}{\epsilon} (\mathbf{u} - \mathbf{u}_{\text{wave}}(\mathbf{x}, t)) G\left(\frac{x-x_{\min}}{L_{\text{gen}}}\right), \quad (2.4.2)$$

where  $H_0$  is the still water depth and  $G(\xi)$  is a non-dimensional relaxation function defined as follows:

$$G(\xi) := \begin{cases} \frac{\exp(-|\log(\alpha)|\xi^2) - \alpha}{1 - \alpha} & \text{if } \xi < 1, \\ 0 & \text{otherwise.} \end{cases}$$

Here  $L_{\text{gen}}$  is the length of the generation zone. In this paper, we take  $\alpha := 0.005$ . We follow a similar methodology as above to absorb waves in a relaxation zone at the outflow boundary. The absorption zone is enforced via the following source term:

$$\mathbf{S}_A(\mathbf{u}) := -\frac{\sqrt{gH_0}}{\epsilon} G\left(\frac{x_{\max}-x}{L_{\text{abs}}}\right) \mathbf{u}_{\text{abs}}, \quad (2.4.3)$$

where  $\mathbf{u}_{\text{abs}}$  are the conserved variables to be ‘‘absorbed’’. That it to say, we enforce the zero value on

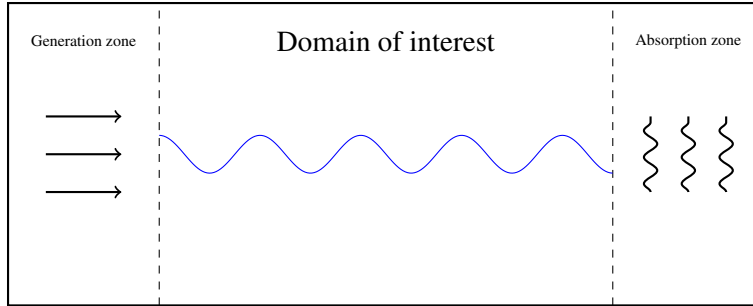


Figure 2.5: Generation and absorption zone schematic.

$\mathbf{u}_{\text{abs}}$  to dissipate the waves. In Figure 2.5, we show a simple schematic of the wave generation and absorption zone layout. We now illustrate numerically the wave generation and wave absorption source terms with an example. Let the computational domain be  $D = (0, 50 \text{ m})$ . We set the still water depth to  $h_0 = 0.4 \text{ m}$ . The wave amplitude is set to  $a = 0.01 \text{ m}$  and wave period is set to  $T_p = 2 \text{ s}^{-1}$ . The generation zone was set to 5 m and the absorption zone to 10 m. In Figure 2.6, we show a space-time plot of the periodic wave solution in the range  $t \in [0, 40 \text{ s}]$ . We see that the waves are smoothly introduced into the domain of interest as the time increases in the generation zone ( $x < 5 \text{ m}$ ). We then see that waves smoothly dissipate in the absorption zone ( $x > 40 \text{ m}$ ).

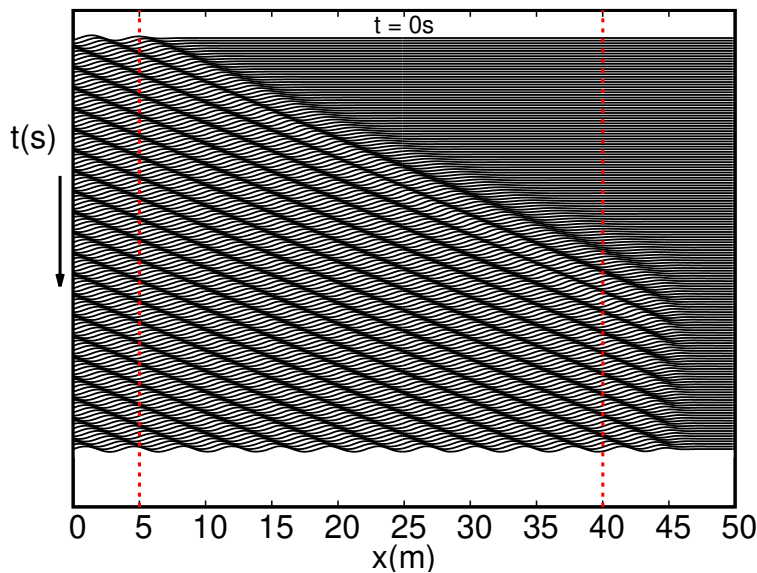


Figure 2.6: Numerical illustration of wave generation/absorption in the space-time domain.

### 3. REFORMULATION OF THE DISPERSIVE SERRE MODEL

#### 3.1 Introduction

In this chapter, we discuss the mathematical reformulation of the dispersive Serre Equations. The reformulation of the Serre model is of great interest to the coastal modeling community since it allows one to exploit different mathematical structures of the model which open the door for efficient implicit and explicit numerical methods. The strategy of reformulating higher-order partial differential equations is ubiquitous in the literature, so various reformulations of the Serre equations exist.

The chapter is organized as follows. In Section 3.2, we propose a novel reformulation of the Serre model that yields a first-order system under two algebraic constraints. This reformulation is the foundation for a hyperbolic relaxation technique introduced in Chapter 4. In Section 3.3, we give an overview of existing reformulations of the Serre model and give direct comparisons of the reformulation proposed in this work to others seen in the literature.

#### 3.2 Reformulation under two algebraic constraints

We now propose a novel reformulation of the Serre model. The motivation behind the reformulation is to reduce the higher-order dispersive terms in the pressure and topography source terms of the Serre model. This is achieved by introducing three first-order auxiliary equations coupled with two algebraic constraints that take the place of the dispersive terms.

We first recall the Serre model. Let  $\mathbf{u} := (h, \mathbf{q})^\top$  be the (smooth) solution variable where  $h$  is the positive water depth and  $\mathbf{q} \in \mathbb{R}^d$  is the momentum vector where  $d$  is the spatial dimension. Let  $\mathbf{v}$  be the velocity vector defined such that  $\mathbf{q} := h\mathbf{v}$ . Then, the Serre model is formulated as follows:

$$\begin{aligned} \partial_t h + \nabla \cdot (h\mathbf{v}) &= 0, \\ \partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla \cdot \left( \frac{1}{2}gh^2 + \underbrace{h^2 \left( \frac{1}{3}\ddot{h} + \frac{1}{2}\dot{\mathbf{k}} \right)}_{\tilde{p}(\mathbf{u})} \right) &= - \left( gh + \underbrace{h \left( \frac{1}{2}\ddot{h} + \dot{\mathbf{k}} \right)}_{\tilde{r}(\mathbf{u})} \right) \nabla z. \end{aligned}$$

The goal moving forward is to replace the non-hydrostatic pressure  $\tilde{p}(\mathbf{u})$  and non-hydrostatic part of the source  $\tilde{r}(\mathbf{u})$  by algebraic expressions. We begin the reformulation of the Serre model with the non-hydrostatic pressure. The reformulation of the topography source term will follow naturally. Rearranging the non-hydrostatic pressure yields:

$$\begin{aligned}\tilde{p}(\mathbf{u}) &= h^2 \left( \frac{1}{3} \ddot{h} + \frac{1}{2} \dot{k} \right) \\ &= \frac{1}{3} h^2 \left( \ddot{h} + \frac{3}{2} \dot{k} \right) \\ &= \frac{1}{3} h^2 D_t \left( \dot{h} + \frac{3}{2} k \right),\end{aligned}$$

where  $D_t$  denotes the total derivative operator (i.e., equivalent to the  $\dot{(\ )}$  notation). We use the following property to continue the analysis.

**Proposition 3.2.1.** *Let  $\mathbf{u} := (h, \mathbf{q})^\top$  be a smooth solution of (2.3.1). Assuming that  $\phi(\mathbf{x}, t)$  is an arbitrary smooth function, then*

$$h \dot{\phi} = \partial_t(h\phi) + \nabla \cdot (h\mathbf{v}\phi).$$

*Proof.* A direction computation of the left hand side yields:

$$\begin{aligned}h \dot{\phi} &= h (\partial_t \phi + \mathbf{v} \cdot \nabla \phi), \\ &= \partial_t(h\phi) - \phi \partial_t h + \nabla \cdot (h\mathbf{v}\phi) - \phi \nabla \cdot (h\mathbf{v}), \\ &= \partial_t(h\phi) + \phi \nabla \cdot (h\mathbf{v}) - \cancel{\phi (\partial_t h + \nabla \cdot (h\mathbf{v}))}, \\ &= \partial_t(h\phi) + \nabla \cdot (h\mathbf{v}\phi).\end{aligned}$$

□

Using Proposition 3.2.1 and continuing the expansion for the non-hydrostatic pressure, we see:

$$\begin{aligned}
\tilde{p}(\mathbf{u}) &= \frac{1}{3}h^2D_t\left(\dot{h} + \frac{3}{2}\mathbf{k}\right) \\
&= \frac{1}{3}h\left(hD_t(\dot{h}) + hD_t\left(\frac{3}{2}\mathbf{k}\right)\right) \\
&= \frac{1}{3}h\left(\partial_t(h\dot{h}) + \nabla\cdot(\mathbf{v}h\dot{h}) + \frac{3}{2}(\partial_t(\mathbf{q}\cdot\nabla z) + \nabla\cdot(\mathbf{v}(\mathbf{q}\cdot\nabla z)))\right).
\end{aligned}$$

Let us introduce the following constraint  $q_3 := \mathbf{q}\cdot\nabla z$  and variable  $\tilde{s}$  such that

$$\partial_t(q_3) + \nabla\cdot(\mathbf{v}q_3) = \tilde{s}.$$

Notice that this is a first-order conservation law with source  $\tilde{s}$ . Also note that the above equation is equivalent to  $h\dot{h} = \tilde{s}$ . We now introduce another variable  $s$  so that

$$-s = \partial_t(h\dot{h}) + \nabla\cdot(\mathbf{v}h\dot{h}) + \frac{3}{2}\tilde{s} \quad (3.2.1)$$

At this point, we can define the algebraic non-hydrostatic pressure of  $\tilde{p}(\mathbf{u}) = -\frac{1}{3}hs$  and non-hydrostatic source  $\tilde{r}(\mathbf{u}) = -\frac{1}{2}s + \frac{1}{4}\tilde{s}$ . However, the equation (3.2.1) still contains higher-order derivatives, so we continue the reduction. We introduce the auxiliary relation  $q_2 = h\dot{h} + \frac{3}{2}q_3$  so that substituting in the (3.2.1) gives:

$$\partial_t(q_2) + \nabla\cdot(\mathbf{v}q_2) = -s.$$

This is another first-order conservation law with a source term  $s$ . Since  $q_2$  is dependent on  $h\dot{h}$ , we can once again reduce the higher-order terms. By Proposition 3.2.1, we see that:

$$h\dot{h} = \partial_t(h^2) + \nabla\cdot(\mathbf{v}h^2).$$

Introducing the algebraic constraint  $q_1 := h^2$  and using the definitions for  $q_2$  and  $q_3$ , the previous



equation simplifies to

$$\partial_t(q_1) + \nabla \cdot (\mathbf{v}q_1) = q_2 - \frac{3}{2}q_3.$$

Again, we have arrived at a first-order conservation law with source  $q_2 - \frac{3}{2}q_3$ . Thus, if  $q_1 := h^2$  and  $q_3 := \mathbf{q} \cdot \nabla z$ , the full Serre pressure and topography source terms in (2.3.2) are equivalently written as:

$$\begin{aligned} p(\mathbf{u}) &= \frac{1}{2}gh^2 - \frac{1}{3}hs, \\ r(\mathbf{u}) &= gh - \frac{1}{2} + \frac{1}{4}\tilde{s}, \\ \partial_t q_1 + \nabla \cdot (\mathbf{v}q_1) &= q_2 - \frac{3}{2}q_3, \\ \partial_t q_2 + \nabla \cdot (\mathbf{v}q_2) &= -s, \\ \partial_t q_3 + \nabla \cdot (\mathbf{v}q_3) &= \tilde{s} \end{aligned}$$

Note that the variables  $s$  and  $\tilde{s}$  can be interpreted as Lagrange multipliers associated with the constraints  $q_1 = h^2$  and  $q_3 = \mathbf{q} \cdot \nabla z$ , respectively. The above ideas are summarized in the following lemma.

**Lemma 3.2.2.** *Let  $\mathbf{u} : D \times (0, T) \rightarrow \mathbb{R}_+ \times \mathbb{R}^d$  be a smooth function. Then  $\mathbf{u}$  solves the dispersive Serre model (2.3.1)–(2.3.2) iff  $(\mathbf{u}, q_1, q_2, q_3)$  solves*

$$\partial_t h + \nabla \cdot \mathbf{q} = 0, \tag{3.2.2a}$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla \cdot (\frac{1}{2}gh^2 - \frac{1}{3}hs) = -(gh - \frac{1}{2}s + \frac{1}{4}\tilde{s})\nabla z, \tag{3.2.2b}$$

$$\partial_t q_1 + \nabla \cdot (\mathbf{v}q_1) = q_2 - \frac{3}{2}q_3 \tag{3.2.2c}$$

$$\partial_t q_2 + \nabla \cdot (\mathbf{v}q_2) = -s \tag{3.2.2d}$$

$$\partial_t q_3 + \nabla \cdot (\mathbf{v}q_3) = \tilde{s} \tag{3.2.2e}$$

$$q_1 = h^2, \quad q_3 = \mathbf{q} \cdot \nabla z \tag{3.2.2f}$$

*Proof.* The forward direction of the proof is summarized in the steps above, so we now prove the

converse. Let us assume that  $(\mathbf{u}, q_1, q_2, q_3)$  solves (3.2.2a)–(3.2.2f).

We set  $\mathbf{u} := (h, \mathbf{q})$  and  $\mathbf{q} := h\mathbf{v}$  where  $\mathbf{v}$  is the velocity vector. Let us set  $\dot{\mathbf{k}} := \partial_t(\mathbf{v} \cdot \nabla z) + \mathbf{v} \cdot \nabla(\mathbf{v} \cdot \nabla z)$ . Then using (3.2.2e) and (3.2.2f) we obtain

$$\partial_t q_3 + \nabla \cdot (\mathbf{v} q_3) = \partial_t (h\mathbf{v} \cdot \nabla z) + \nabla \cdot (\mathbf{v} h(\mathbf{v} \cdot \nabla z)) = h D_t(\mathbf{v} \cdot \nabla z) = \tilde{s},$$

which gives  $h\dot{\mathbf{k}} = \tilde{s}$ . Similarly, Proposition 3.2.1 implies that  $h\dot{h} = \partial_t(h^2) + \nabla \cdot (\mathbf{v} h^2)$ . This identity, together with  $q_1 := h^2$  and (3.2.2c), gives  $h\dot{h} = q_2 - \frac{3}{2}q_3$ . Then

$$h\ddot{h} = \partial_t(h\dot{h}) + \nabla \cdot (\mathbf{v} h\dot{h}) = \partial_t(q_2 - \frac{3}{2}q_3) + \nabla \cdot (\mathbf{v}(q_2 - \frac{3}{2}q_3)) = -s - \frac{3}{2}\tilde{s}.$$

This implies that  $s = -h(\ddot{h} + \frac{3}{2}\dot{\mathbf{k}})$ . Let us set  $p := \frac{1}{2}gh^2 - \frac{1}{3}hs$ . Then  $p = \frac{1}{2}gh^2 + h^2(\frac{1}{3}\ddot{h} + \frac{1}{2}\dot{\mathbf{k}})$ . This is the expression of the pressure in (2.3.2a).

Finally let us set  $r := gh - \frac{1}{2}s + \frac{1}{4}\tilde{s}$ . Then the above computations imply that  $r = gh + \frac{1}{2}h(\ddot{h} + \frac{3}{2}\dot{\mathbf{k}}) + \frac{1}{4}h\dot{\mathbf{k}}$ , i.e.,  $r = gh + \frac{1}{2}h\ddot{h} + \dot{\mathbf{k}}$ . This is the expression of the source term  $r$  in (2.3.2b). Hence we have established that  $\mathbf{u}$  solves (2.3.1)–(2.3.2). This completes the proof.  $\square$

The following is another way to reformulate Lemma 3.2.2.

**Proposition 3.2.3** (Co-dimension 2). *Let  $\mathbf{u} : D \times (0, T) \rightarrow \mathbb{R}_+ \times \mathbb{R}^d$  be a smooth function. Then  $\mathbf{u}$  solves the dispersive Serre model (2.3.1)–(2.3.2) if and only if  $(\mathbf{u}, q_1, q_2, q_3)$  solves the quasi-linear first-order system (3.2.2a)–(3.2.2e) on the co-dimension 2 manifold  $\{(h, \mathbf{q}, q_1, q_2, q_3) \in \mathbb{R}_+ \times \mathbb{R}^d \times \mathbb{R}_+ \times \mathbb{R}^2 \mid q_1 = h^2, q_3 = \mathbf{q} \cdot \nabla z\}$ .*

*Remark 3.2.4* (Initial conditions). Let  $(h_0, \mathbf{q}_0)$  be the initial state for (2.3.1)–(2.3.2). Then the corresponding initial state for (3.2.2a)–(3.2.2f) is

$$h(\mathbf{x}, 0) = h_0(\mathbf{x}), \quad \mathbf{q}(\mathbf{x}, 0) = \mathbf{q}_0(\mathbf{x}), \quad q_1(\mathbf{x}, 0) = h_0(\mathbf{x})^2, \quad (3.2.3a)$$

$$q_3(\mathbf{x}, 0) = \mathbf{q}_0(\mathbf{x}) \cdot \nabla z, \quad q_2(\mathbf{x}, 0) = -h_0(\mathbf{x}) \nabla \cdot \mathbf{v}_0(\mathbf{x}) + \frac{3}{2}q_3(\mathbf{x}, 0). \quad (3.2.3b)$$

*Remark 3.2.5* (Reformulation with a flat-bottom). If no topography effects are considered (i.e.,  $z(\mathbf{x}) \equiv 0$ ), repeating the above reformulation yields the following (smaller) system of equations for the solution variable  $\mathbf{u} := (\mathbf{h}, \mathbf{q}, q_1, q_2)^\top$ :

$$\partial_t \mathbf{h} + \nabla \cdot \mathbf{q} = 0 \quad (3.2.4a)$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla \left( \frac{1}{2} g \mathbf{h}^2 - \frac{1}{3} \mathbf{h} s \right) = 0, \quad (3.2.4b)$$

$$\partial_t q_1 + \nabla \cdot (\mathbf{v} q_1) = q_2 \quad (3.2.4c)$$

$$\partial_t q_2 + \nabla \cdot (\mathbf{v} q_2) = -s \quad (3.2.4d)$$

$$q_1 = \mathbf{h}^2 \quad (3.2.4e)$$

### 3.3 Literature review

There are many ways to reformulate the dispersive Serre model. Two popular approaches in literature are: (i) reformulate the model as a system of first-order conservation equations with sources for the construction of explicit methods; (ii) reformulation of the model with linear differential operators to construct efficient implicit methods. For instance, in Gavriluk and Shugrin [22, Eqs. (5.12)–(5.15)] the authors reformulated the model with a flat bottom as a first-order system with a constraint on the divergence of the velocity. In Lannes and Bonneton [46, Eq. (26)], the authors reformulated the momentum evolution equation of the Serre Equations (called Green-Naghdi Equations therein) using two linear differential operators to induce regularizing effects for easing numerical computations. In Bristeau et al. [12, Eq. (50)], the authors recover the same first-order reformulation (with a flat-bottom) originally proposed in [22]. In Escalante et al. [19, Eq. (1)], the authors present a first-order reformulation with the assumption of mild bottom variation (i.e., dropping the terms containing  $\nabla \cdot (\nabla z)$  and  $\|\nabla z\|^2$ ). In Fernandez-Nieto et al. [21, Eq. (3.16)], the authors present a 1D first-order reformulation the Serre model with full topography effects with constraints on the divergence of the velocity and the quantity  $\mathbf{v} \cdot \nabla z$ .

To put the present work in perspective with respect to these techniques, we recall the reformulations proposed in [22, 12, 19], but contrary to [19] we keep all the effects induced by the

topography. Keeping the effects induced by the topography gives a reformulation equivalent to that proposed in [21] up to a change of constant. The starting point is again the Serre model. Assume that  $\mathbf{u} := (\mathbf{h}, \mathbf{q})^\top$  is a smooth solution to (2.3.1)–(2.3.2). Let

$$w := -\frac{1}{2}\mathbf{h}\nabla\cdot\mathbf{v} + \frac{3}{4}\mathbf{v}\cdot\nabla z. \quad (3.3.1)$$

Recall that by the mass conservation equation, we have  $\dot{\mathbf{h}} = -\mathbf{h}\nabla\cdot\mathbf{v}$ . Also recall the notation:  $\dot{\mathbf{k}} := \partial_t(\mathbf{v}\cdot\nabla z) + \mathbf{v}\cdot\nabla(\mathbf{v}\cdot\nabla z)$ . Then, applying the total derivative operator  $D_t$  to both sides of (3.3.1) yields:

$$\begin{aligned} D_t w &= \partial_t w + \mathbf{v}\cdot\nabla w \\ &= \frac{1}{2}\ddot{\mathbf{h}} + \frac{3}{4}\dot{\mathbf{k}} \\ &= \frac{3}{2}\left(\frac{1}{3}\ddot{\mathbf{h}} + \frac{1}{2}\dot{\mathbf{k}}\right) \\ &= \frac{3}{2}\frac{\bar{p}(\mathbf{u})}{\mathbf{h}}, \end{aligned}$$

where  $\bar{p}(\mathbf{u}) := \frac{1}{3}\mathbf{h}\ddot{\mathbf{h}} + \frac{1}{2}\mathbf{h}\dot{\mathbf{k}}$ . We multiply the above equation by  $\mathbf{h}$  and use Proposition 3.2.1 to obtain a first-order conservation law:

$$\partial_t(\mathbf{h}w) + \nabla\cdot(\mathbf{u}\mathbf{h}w) = \frac{3}{2}\bar{p}(\mathbf{u}).$$

Then, another first-order reformulation of the dispersive Serre model with topography is given as follows:

$$\partial_t \mathbf{h} + \nabla\cdot\mathbf{q} = 0, \quad (3.3.2a)$$

$$\partial_t \mathbf{q} + \nabla\cdot(\mathbf{v} \otimes \mathbf{q}) + \nabla\left(\frac{1}{2}g\mathbf{h}^2 + \mathbf{h}\bar{p}(\mathbf{u})\right) = -(g\mathbf{h} + \frac{3}{2}\bar{p}(\mathbf{u}) + \frac{1}{4}\mathbf{h}\dot{\mathbf{k}})\nabla z, \quad (3.3.2b)$$

$$\partial_t(\mathbf{h}w) + \nabla\cdot(\mathbf{u}\mathbf{h}w) = \frac{3}{2}\bar{p}(\mathbf{u}), \quad (3.3.2c)$$

$$\nabla\cdot\mathbf{v} + \frac{w - \frac{3}{4}\mathbf{v}\cdot\nabla z}{\frac{1}{2}\mathbf{h}} = 0. \quad (3.3.2d)$$

We recover Eq. (1) in [19], up to the term  $\frac{1}{4}\mathbf{h}\dot{\mathbf{k}}$ , which is neglected therein, and up to the coefficient  $\frac{3}{4}$  in (3.3.2d). The above system bears some resemblance to (3.2.2). In particular we observe that  $\mathbf{h}\bar{p}(\mathbf{u}) = -\frac{1}{3}s(\mathbf{u})$  and  $\mathbf{h}w = \frac{1}{2}q_2$ . Notice however that in our system (3.2.2) the two constraints (3.2.2f) are purely algebraic (i.e., these constraints are enforced in the phase space), whereas the constraint (3.3.2d) is differential. As a result, the technique proposed in [19] to relax the differential constraint (3.3.2d) is fundamentally different from (4.4.1h).

*Remark 3.3.1* (Reformulation in [21]). In [21] where the full topography effects of the Serre model are considered, two constraints are introduced for the reformulation:

$$w_s := -\mathbf{h}\nabla\cdot\mathbf{v} + \mathbf{v}\cdot\nabla z, \quad \tilde{w} := w_s + \frac{1}{2}\mathbf{h}\nabla\cdot\mathbf{v}$$

Combining these two constraints into one yields:

$$\tilde{w} := -\frac{1}{2}\mathbf{h}\nabla\cdot\mathbf{v} + \mathbf{v}\cdot\nabla z.$$

This constraint is of the form (3.3.1) up to the constant on the  $\mathbf{v}\cdot\nabla z$  term. Thus, we can repeat the above process and derive the first-order formulation introduced in [21]. Again, in [21], the constraint is differential and thus different from the proposed reformulation in §3.2.

## 4. A HYPERBOLIC RELAXATION TECHNIQUE FOR SOLVING THE SERRE EQUATIONS

### 4.1 Introduction

In this chapter, we introduce a hyperbolic relaxation technique for solving the Serre model (2.3.1) for dispersive water waves with topography effects. The work shown in this chapter addresses one of the main challenges for solving the Serre equations. In particular, the hyperbolic relaxation technique circumvents the dispersive time step restriction discussed in §2.3.4 since the technique reformulates the Serre model as a first-order hyperbolic system. This allows for explicit time-stepping when discretizing the equations in time under the usual hyperbolic CFL condition. The hyperbolic reformulation technique sets up the framework for introducing an explicit, continuous finite element approximation that is invariant-domain preserving and well-balanced in Chapter 5.

This chapter is organized as follows. In Section 4.2, we recall the reformulation of the Serre model introduced in Chapter 3 which is the starting point of the relaxation technique. We also introduce some notation preliminaries in this section. Then, in Section 4.3 we describe how to relax the constraints  $\{q_1 := h^2; q_3 := \mathbf{q} \cdot \nabla z\}$  discussed in §3.2 and introduce them in the model. This relaxation can be thought of as a penalty technique. Then, in Section 4.3.1, we discuss how to define the relaxed non-hydrostatic pressure so that the new model is an energy-consistent approximation of the original Serre Equations. In Section 4.4, we present a generic form of the relaxed model. Then, in Section 4.5, we derive important properties of the new relaxed model. In particular, we show the relaxed model is indeed hyperbolic. We also describe how to define the relaxation functional  $\Gamma$  in Section 4.5.2 so that the resulting hyperbolic system remains compatible with dry states. We then derive an energy inequality for the model and derive the linear dispersion properties. In Section 4.6, we review the literature on similar relaxation techniques for solving the Serre equations to put the work shown here in context. Finally, we summarize the chapter.

## 4.2 Preliminaries

For completeness, we recall here the first-order reformulation of the Serre Equations under the constraints  $\{q_1 := h^2; q_3 := \mathbf{q} \cdot \nabla z\}$ . Let  $\mathbf{u} := (h, \mathbf{q}, q_1, q_2, q_3)^\top$  be the solution variable for the following first-order system:

$$\partial_t h + \nabla \cdot \mathbf{q} = 0, \quad (4.2.1a)$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla \left( \frac{1}{2} g h^2 - \frac{1}{3} h s \right) = - \left( g h - \frac{1}{2} s + \frac{1}{4} \tilde{s} \right) \nabla z, \quad (4.2.1b)$$

$$\partial_t q_1 + \nabla \cdot (\mathbf{v} q_1) = q_2 - \frac{3}{2} q_3, \quad (4.2.1c)$$

$$\partial_t q_2 + \nabla \cdot (\mathbf{v} q_2) = -s, \quad (4.2.1d)$$

$$\partial_t q_3 + \nabla \cdot (\mathbf{v} q_3) = \tilde{s}, \quad (4.2.1e)$$

$$q_1 = h^2, \quad q_3 = \mathbf{q} \cdot \nabla z. \quad (4.2.1f)$$

To simplify the notation in what follows, we adopt a similar notation as in Favrie and Gavrilyuk [20] and introduce the following primitive variables for the auxiliary quantities:  $q_1 := h\eta$ ,  $q_2 := h\omega$ ,  $q_3 := h\beta$ . We also introduce the quantity  $\bar{\lambda}$  which is a non-dimensional number of order one. Finally, we introduce a relaxation parameter  $\epsilon$ , a small length scale, which we will later replace by the local mesh-size when the model is approximated in space.

## 4.3 Relaxing the $\{q_1 = h^2; q_3 = \mathbf{q} \cdot \nabla z\}$ constraints

In this section, we propose a technique to relax the constraints  $\{q_1 = h^2; q_3 = \mathbf{q} \cdot \nabla z\}$  so that the resulting system is hyperbolic and remains compatible with dry states. This technique is a loose adaptation of the Lagrangian penalty technique introduced in [20].

Let  $\Gamma \in C^2(\mathbb{R}; [0, \infty))$  be some smooth non-negative function with constraints  $\Gamma(1) = 0$  and  $\Gamma'(0) = 0$ . We replace the quantity  $s$  in (4.2.1) by

$$s \rightarrow -\frac{\bar{\lambda} g}{\epsilon} h^2 \Gamma' \left( \frac{\eta}{h} \right). \quad (4.3.1)$$

Notice that the quantity  $\frac{\bar{\lambda}g}{\epsilon}h^2$  has dimensions of  $m\frac{m}{s^2}$  which is what we expect since  $-\frac{\bar{\lambda}g}{\epsilon}h^2\Gamma'(\frac{\eta}{h})$  should be an ansatz for  $-h(\ddot{h} + \frac{3}{2}\dot{k})$ . The purpose of this term is to enforce the constraint  $q_1 = h^2$  through the  $\Gamma$  function; that is to say, we want to enforce the ratio  $\frac{\eta}{h}$  to be close to 1 as  $\epsilon \rightarrow 0$  since  $\Gamma'(1) = 0$ . We call (4.3.1) the penalty term for the  $q_2$  equation. We note that the definition of the  $\Gamma$  functional will be discussed in §4.5.2

Let  $\Phi \in C^0(\mathbb{R}; \mathbb{R})$  be a function such that  $\xi\Phi(\xi) \geq 0$  for all  $\xi \in \mathbb{R}$ . Let  $h_0$  be a reference water depth. We then replace  $\tilde{s}$  in (4.2.1) by:

$$\tilde{s} \rightarrow \frac{\bar{\lambda}g}{\epsilon}h_0h\Phi\left(\frac{\mathbf{v}\cdot\nabla z - \beta}{\sqrt{gh_0}}\right). \quad (4.3.2)$$

Notice that this quantity also has dimensions of  $m\frac{m}{s^2}$  since it is an ansatz for  $h\dot{k}$ . The purpose of this term is to enforce the constraint  $q_3 = \mathbf{q}\cdot\nabla z$  through the  $\Phi$  function; that is to say, we want to enforce the ratio  $\frac{\mathbf{v}\cdot\nabla z - \beta}{\sqrt{gh_0}}$  to be close to 0 as  $\epsilon \rightarrow 0$ . We call (4.3.2) the penalty term for the  $q_3$  equation. Note that the  $\Phi$  function is defined in Remark 4.4.1.

*Remark 4.3.1* ( $\epsilon$  vs.  $\frac{\bar{\lambda}}{\epsilon}$ ). Since  $\bar{\lambda}$  is a non-dimensional number, it is also possible to let the quantity  $\frac{\bar{\lambda}}{\epsilon}$  be the relaxation parameter instead of just  $\epsilon$ . Then,  $\epsilon \rightarrow 0$  is equivalent to  $\bar{\lambda} \rightarrow \infty$ . However, this condition on  $\bar{\lambda}$  is ad-hoc while the condition  $\epsilon \rightarrow 0$  will be directly correlated to decreasing the local-mesh size when we approximate the model in space.

### 4.3.1 The correlation between the energy functional and relaxed pressure and source terms

It is only natural to assume that the new relaxed pressure is given by  $p(\mathbf{u}) = \frac{1}{2}gh^2 - \frac{1}{3}hs \rightarrow \frac{1}{2}gh^2 + \frac{1}{3}\frac{\lambda g}{\epsilon}h^3\Gamma'(\frac{\eta}{h})$  (actually, this is similar to the augmented Lagrangian approach: see §4.5.2). However, this expression leads to a relaxed model that is not an energy-consistent relaxation of the original Serre Equations. That is to say, the relaxed energy functional with this pressure does not converge to the Serre energy functional as  $\epsilon \rightarrow 0$ . To find the correct expression for the relaxed pressure, we work through the derivation of an energy inequality (see: Proposition 4.5.3). The



derivation shown in §4.5.3 yields a relaxed pressure of the form:

$$p_\epsilon(\mathbf{u}) = \frac{1}{2}gh^2 + \tilde{p}_\epsilon(\mathbf{u}),$$

where

$$\begin{aligned} \tilde{p}_\epsilon(\mathbf{u}) &= -\frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}\eta^3\partial_x(x^{-2}\Gamma(x))|_{x=\eta h^{-1}} \\ &= -\frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}h^2\left(\eta\Gamma'(\frac{\eta}{h}) - 2h\Gamma(\frac{\eta}{h})\right) \\ &= -\frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}h^2\eta\Gamma'(\frac{\eta}{h}) + \frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}2h^3\Gamma(\frac{\eta}{h}). \end{aligned}$$

By working through the energy argument, we will see that we must also replace  $-\frac{3}{2}q_3$  on the right-hand side of (4.2.1c) by  $-\frac{3}{2}\mathbf{q}\cdot\nabla z$ . This is consistent since  $q_3$  should be equal to  $\mathbf{q}\cdot\nabla z$ .

*Remark 4.3.2* (Relaxed non-hydrostatic pressure). After replacing  $s$  by  $\frac{\bar{\lambda}g}{\epsilon}h^2\Gamma'(\frac{\eta}{h})$  in (3.2.2b), we observe that the definition of  $\tilde{p}_\epsilon(\mathbf{u})$  is compatible with the definition  $\tilde{p}(\mathbf{u}) = -\frac{1}{3}hs$  up to the remainder  $\frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}2h^3\Gamma(\frac{\eta}{h})$ . However, this remainder is small when the ratio  $\frac{\eta}{h}$  is close to 1. More precisely, using Taylor expansions at 1, we have

$$\Gamma(1) = 0 = \Gamma(\frac{\eta}{h}) + h^{-1}(h - \eta)\Gamma'(\frac{\eta}{h}) + h^{-2}\mathcal{O}(h - \eta)^2,$$

which shows that  $2h\Gamma(\frac{\eta}{h})/\eta|\Gamma'(\frac{\eta}{h})| = \mathcal{O}(\frac{|\eta-h|}{\eta})$ . Hence the ratio  $2h\Gamma(\frac{\eta}{h})/\eta|\Gamma'(\frac{\eta}{h})|$  is small as  $\eta \rightarrow h$ , which proves that  $\tilde{p}_\epsilon(\mathbf{u})$  is indeed a consistent approximation of  $\tilde{p}(\mathbf{u}) = -\frac{1}{3}hs$  as  $\eta \rightarrow h$ .

#### 4.4 The generic relaxed model

Finally, the relaxed system for the conserved variables  $\mathbf{u} := (\mathbf{h}, \mathbf{q}, q_1, q_2, q_3)^\top$  is formulated as follows:

$$\partial_t \mathbf{h} + \nabla \cdot \mathbf{q} = 0, \quad (4.4.1a)$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla p_\epsilon(\mathbf{u}) = -r_\epsilon(\mathbf{u}) \nabla z, \quad (4.4.1b)$$

$$\partial_t q_1 + \nabla \cdot (\mathbf{v} q_1) = q_2 - \frac{3}{2} \mathbf{q} \cdot \nabla z, \quad (4.4.1c)$$

$$\partial_t q_2 + \nabla \cdot (\mathbf{v} q_2) = -s_\epsilon(\mathbf{u}), \quad (4.4.1d)$$

$$\partial_t q_3 + \nabla \cdot (\mathbf{v} q_3) = \tilde{s}_\epsilon(\mathbf{u}), \quad (4.4.1e)$$

$$p_\epsilon(\mathbf{u}) := \frac{1}{2} g h^2 + \tilde{p}_\epsilon(\mathbf{u}), \quad \tilde{p}_\epsilon(\mathbf{u}) := -\frac{1}{3} \frac{\bar{\lambda} g}{\epsilon} h^2 \left( \eta \Gamma' \left( \frac{\eta}{h} \right) - 2h \Gamma \left( \frac{\eta}{h} \right) \right), \quad (4.4.1f)$$

$$r_\epsilon(\mathbf{u}) := g h - \frac{1}{2} s_\epsilon(\mathbf{u}) + \frac{1}{4} \tilde{s}_\epsilon(\mathbf{u}), \quad (4.4.1g)$$

$$s_\epsilon(\mathbf{u}) := \frac{\bar{\lambda} g}{\epsilon} h^2 \Gamma' \left( \frac{\eta}{h} \right), \quad \tilde{s}_\epsilon(\mathbf{u}) := \frac{\bar{\lambda} g}{\epsilon} h_0 h \Phi \left( \frac{\mathbf{v} \cdot \nabla z - \beta}{\sqrt{g h_0}} \right). \quad (4.4.1h)$$

The relaxed model (4.4.1) is incomplete since we have yet to explicitly define the  $\Gamma$  and  $\Phi$  functions. It is shown in Proposition 4.5.2 that the hyperbolicity of the relaxed model is dependent on the form of  $\Gamma$  so we defer defining it here.

*Remark 4.4.1* (Defining the  $\Phi$  function). The condition  $\xi \Phi(\xi) \geq 0$  for all  $\xi \in \mathbb{R}$  is important for deriving the energy inequality in Prop. 4.5.3. In particular, it is used to show that the energy equation is negative. This is the only information needed to construct  $\Phi$ . Thus, the simplest choice is to set  $\Phi(\xi) := \xi$  so that  $\xi \Phi(\xi) = \xi^2 \geq 0$ . This choice of  $\Phi$  gives:

$$\tilde{s}_\epsilon(\mathbf{u}) = \frac{\bar{\lambda}}{\epsilon} \sqrt{g h_0} (\mathbf{q} \cdot \nabla z - q_3).$$

*Remark 4.4.2* (Generality). Unless otherwise stated, we express the analysis in this work with generic functions  $\Gamma$  and  $\Phi$  to emphasize the generality of the relaxation procedure.

## 4.5 Properties

In this section, we derive important results for the relaxed model (4.4.1).

### 4.5.1 Hyperbolicity for the relaxed model

We now verify that the relaxed system (4.4.1) is hyperbolic as defined by Definition (1.2.1). Let  $\mathbf{n}$  be a unit vector in  $\mathbb{R}^d$ . Let  $\mathbb{f}(\mathbf{u})$  denote the conservative flux of (4.4.1). We introduce the notation  $\mathbf{q} \cdot \mathbf{n} := hv$  and  $\mathbf{q}^\perp := \mathbf{q} - (\mathbf{q} \cdot \mathbf{n})\mathbf{q}$ ,  $q^\perp := hv^\perp$  and make a change of coordinates so that the conserved variable can be re-written as  $\mathbf{u} = (h, (\mathbf{q} \cdot \mathbf{n}), \mathbf{q}^\perp, q_1, q_2, q_3)^\top$ . Recalling the notation  $q_1 := h\eta$ , we introduce the function  $\bar{p}(h, q_1)$  such that  $\bar{p}(h, q_1) := \tilde{p}_\epsilon(h, \eta)$ . The quasi-linear form of the equation in the new change of coordinates is given by

$$\partial_t \mathbf{u} + D(\mathbb{f}(\mathbf{u})\mathbf{n})\partial_x \mathbf{u} = \mathbf{S},$$

where  $D(\mathbb{f}(\mathbf{u})\mathbf{n})$  is the Jacobian matrix of the flux  $\mathbb{f}(\mathbf{u})\mathbf{n}$  defined as follows:

$$D(\mathbb{f}(\mathbf{u})\mathbf{n}) = \begin{pmatrix} 0 & 1 & \mathbf{0}^\top & 0 & 0 & 0 \\ gh - v^2 + \partial_h \bar{p} & 2v & \mathbf{0}^\top & \partial_{q_1} \bar{p} & 0 & 0 \\ -v\mathbf{v}^\perp & \mathbf{v}^\perp & v & 0 & 0 & 0 \\ -v\eta & \eta & \mathbf{0}^\top & v & 0 & 0 \\ -vw & w & 0 & 0 & v & 0 \\ -v\beta & \beta & 0 & 0 & 0 & v \end{pmatrix}.$$

Denoting by  $\mu$  a generic eigenvalue, the characteristic polynomial of the Jacobian  $D(\mathbb{f}(\mathbf{u})\mathbf{n})$  is

$$(\mu - v)^{2+d}((\mu - v)^2 - (gh + \eta\partial_{q_1} \bar{p} + \partial_h \bar{p})).$$

This characteristic polynomial yields eigenvalues:  $v$ , with multiplicity  $d+2$ ;  $v - \sqrt{gh + \eta\partial_{q_1} \bar{p} + \partial_h \bar{p}}$  and  $v + \sqrt{gh + \eta\partial_{q_1} \bar{p} + \partial_h \bar{p}}$ , each with multiplicity 1.

To verify hyperbolicity, we need to investigate the sign of  $gh + \eta\partial_{q_1} \bar{p} + \partial_h \bar{p}$ . In particular, we

would like this quantity to be positive. Using that  $\tilde{p}(\mathbf{h}, \eta) := \bar{p}(\mathbf{h}, q_1)$ , we observe by the chain rule that  $\partial_{\mathbf{h}}\tilde{p} = \eta\partial_{q_1}\bar{p} + \partial_{\mathbf{h}}\bar{p}$ . Let us introduce the following change of variable:  $\mathbf{h}(\tau) = \tau^{-1}$ , i.e.,  $\tau(\mathbf{h}) = \mathbf{h}^{-1}$ . Then,  $\hat{p}(\tau, \eta) := \tilde{p}(\tau^{-1}, \eta)$ , and recalling (4.4.1f), this gives  $\hat{p}(\tau, \eta) = -\frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}\partial_{\tau}(\tau^{-2}\Gamma(\eta\tau))$ . Using  $\partial_{\mathbf{h}}\tilde{p}(\mathbf{h}, \eta) = -\tau^2\partial_{\tau}\hat{p}(\tau, \eta)$ , this in turn implies that the following inequality

$$\begin{aligned} gh + \partial_{\mathbf{h}}\tilde{p}(\mathbf{h}, \eta) &= gh + \frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}\tau^2\partial_{\tau\tau}(\tau^{-2}\Gamma(\eta\tau)) \\ &= gh\left(1 + \frac{\bar{\lambda}}{3\epsilon}\eta \times (x^3\partial_{xx}(x^{-2}\Gamma(x)))\Big|_{x=\eta\mathbf{h}^{-1}}\right), \end{aligned}$$

Thus

$$gh\left(1 + \frac{\bar{\lambda}}{3\epsilon}\eta \times (x^3\partial_{xx}(x^{-2}\Gamma(x)))\Big|_{x=\eta\mathbf{h}^{-1}}\right) \geq 0,$$

is a necessary and sufficient condition for hyperbolicity. The results are summarized in the following proposition.

**Proposition 4.5.1** (Hyperbolicity). *Let  $\mathbb{f}(\mathbf{u})$  be the conservative flux of the system (4.4.1). For any unit vector  $\mathbf{n} \in \mathbb{R}^d$ , the  $d + 4$  eigenvalues of the Jacobian matrix of the flux  $\mathbb{f}(\mathbf{u})\mathbf{n}$  are  $\mu_k = \mathbf{v} \cdot \mathbf{n}$ ,  $k \in \{2:d+3\}$  and*

$$\mu_1 = \mathbf{v} \cdot \mathbf{n} - \sqrt{gh + \partial_{\mathbf{h}}\tilde{p}(\mathbf{h}, \eta)}, \quad \mu_{d+4} = \mathbf{v} \cdot \mathbf{n} + \sqrt{gh + \partial_{\mathbf{h}}\tilde{p}(\mathbf{h}, \eta)}. \quad (4.5.1)$$

The system (4.4.1) is hyperbolic iff the following holds for all  $\eta \in \mathbb{R}$  and all  $\mathbf{h} \geq 0$ :

$$gh \left(1 + \frac{1}{3}\frac{\bar{\lambda}}{\epsilon}\eta \left(x^3\partial_{xx}(x^{-2}\Gamma(x))\Big|_{x=\eta\mathbf{h}^{-1}}\right)\right) \geq 0. \quad (4.5.2)$$

## 4.5.2 Defining the $\Gamma$ function

In this section, we construct the smooth functional  $\Gamma$  that: (i) satisfies the constraints  $\Gamma \in C^2(\mathbb{R}; [0, \infty))$ ,  $\Gamma \geq 0$ ,  $\Gamma(1) = \Gamma'(1) = 0$ ; (ii) satisfies the hyperbolicity condition (4.5.2); (iii) yields a relaxed model that is compatible with dry states.

The hyperbolicity condition (4.5.2) suggests that the explicit form for  $\Gamma$  should satisfy:

$$\frac{1}{3} \frac{\bar{\lambda}}{\epsilon} \eta (x^3 \partial_{xx}(x^{-2}\Gamma(x)))|_{x=\eta h^{-1}} \geq 0.$$

Note that a priori we don't know that  $\eta$  is non-negative (it is assumed, but not proven). Thus, it is important to construct  $\Gamma(\eta h^{-1})$  for all  $\eta \in \mathbb{R}$ . Since  $h > 0$ , we set  $\eta = h \frac{\eta}{h} = hx$  where  $x = \frac{\eta}{h}$  and rewrite the above inequality:

$$\frac{1}{3} \frac{\bar{\lambda}}{\epsilon} h (x^4 \partial_{xx}(x^{-2}\Gamma(x)))|_{x=\eta h^{-1}} \geq 0.$$

The quantity  $\frac{1}{3} \frac{\bar{\lambda}}{\epsilon} h$  is positive, so finding  $\Gamma$  is reduced to solving the following ordinary differential equation:

$$\begin{cases} 6\Gamma - 4x\Gamma' + x^2\Gamma'' = \alpha, \\ \Gamma(1) = 0, \quad \Gamma'(1) = 0, \end{cases}$$

where  $\alpha$  is a positive constant. Solving the ordinary differential equation yields  $\Gamma(x) = \frac{\alpha}{6}(1 + 2x)(x - 1)^2$ . Note that the condition  $\alpha > 0$  is important since  $\alpha = 0$  would yield  $\Gamma(x) = 0$ . For simplicity, we take  $\alpha = 6$ . This yields a relaxed pressure:

$$\begin{aligned} \tilde{p}_\epsilon(h, \eta) &= \frac{\bar{\lambda}g}{\epsilon} \frac{2}{3} (h^3 - \eta^3), \\ &\implies \\ \bar{p}(h, q_1) &= \frac{\bar{\lambda}g}{\epsilon} \frac{2}{3} \left( h^3 - \frac{q_1^3}{h^3} \right), \end{aligned}$$

and relaxed source:

$$\begin{aligned} s_\epsilon(h, \eta) &= \frac{\bar{\lambda}g}{\epsilon} 2\eta(\eta - h), \\ &\implies \\ \bar{s}(h, q_1) &= \frac{\bar{\lambda}g}{\epsilon} 2 \left( \frac{q_1^2}{h^2} - q_1 \right), \end{aligned}$$

The above relaxed pressure yields a relaxed model that is uniformly hyperbolic. However, it has been observed numerically that this model does not behave properly with dry states ( $h \rightarrow 0$ ) when  $\eta \leq h$ . This can be seen as follows. Since  $\lim_{h \rightarrow 0} \frac{1}{2}gh^2 = 0$ , we consider the ratio of the non-hydrostatic relaxed pressure to the hydrostatic pressure:  $\tilde{p}_\epsilon(h, \eta)/(\frac{1}{2}gh^2)$ :

$$\frac{\tilde{p}_\epsilon(h, \eta)}{(\frac{1}{2}gh^2)} = \frac{\bar{\lambda} 4 h^3 - \eta^3}{\epsilon 3 h^2}$$

We see that as  $h \rightarrow 0$ , this quantity is unbounded. Thus, this particular  $\Gamma$  functional gives a relaxed model that is uniformly hyperbolic, but not compatible with dry states.

To find a suitable workaround, we adapt the augmented Lagrangian approach of Favrie and Gavrilyuk [20]. That is to say, we set the relaxed pressure to  $\tilde{p}(h, \eta) = -\frac{1}{3}hs = -\frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}hh^2\Gamma'(\frac{\eta}{h})$  where  $s(h, \eta) = \frac{\bar{\lambda}g}{\epsilon}h^2\Gamma'(\frac{\eta}{h})$ . However, recall that by the energy argument in Proposition 4.5.3, this new relaxed non-hydrostatic pressure also has to satisfy:  $\tilde{p}(h, \eta) = -\frac{1}{3}\frac{\bar{\lambda}g}{\epsilon}\eta^3\partial_x(x^{-2}\Gamma(x))|_{x=\eta h^{-1}}$ . This is only possible if  $h^3\Gamma'(\frac{\eta}{h}) = \eta^3\partial_x(x^{-2}\Gamma(x))|_{x=\eta h^{-1}}$ . This condition yields the following ordinary differential equation for  $\Gamma$ :

$$\begin{cases} \Gamma'(x) = -2\Gamma(x) + x\Gamma'(x), \\ \Gamma(1) = 0. \end{cases}$$

The solution to this ODE is given by  $\Gamma(x) = c(x-1)^2$  where  $c$  is an arbitrary, positive constant. For simplicity, we take  $c = 3$ . Substituting this definition of  $\Gamma$  into the hyperbolicity condition (4.5.2) yields:

$$gh + 2\frac{\bar{\lambda}g}{\epsilon}(3h - 2\eta) \geq 0.$$

We see that hyperbolicity may be lost when  $2\eta > 3h$  i.e.,  $x > \frac{3}{2}$ .

The solution we propose moving forward is to combine both functionals in a piecewise manner.

That is, we define  $\Gamma$  as follows:

$$\Gamma(x) = \begin{cases} 3(1-x)^2 & \text{if } x \leq 1, \\ (1+2x)(1-x)^2 & \text{if } 1 \leq x. \end{cases} \quad (4.5.3)$$

Notice that  $\Gamma(x) \geq 0$  for any  $x \in \mathbb{R}$ ,  $\Gamma(1) = 0$ ,  $\Gamma'(1) = 0$ , and  $\Gamma \in C^2(\mathbb{R}; [0, \infty))$ . A plot of this function can be seen in Figure 4.1. The final source term  $h^2\Gamma'(\eta h^{-1})$  in the balance equation for

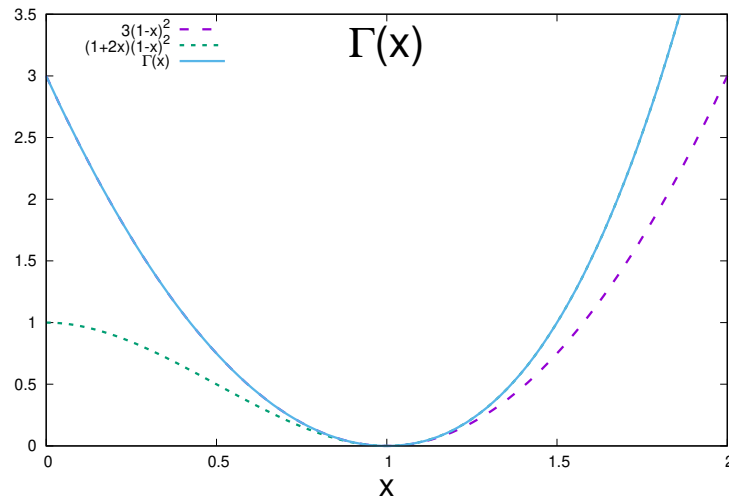


Figure 4.1: A plot of  $\Gamma(x)$ .

$h\omega$  and pressure  $\tilde{p}(h, \eta)$  are given by

$$h^2\Gamma'(\eta h^{-1}) = \begin{cases} 6(\eta h - h^2), & \text{if } \eta \leq h \\ 6(\eta^2 - \eta h), & \text{if } h \leq \eta, \end{cases} \quad (4.5.4)$$

$$\tilde{p}(h, \eta) = -\frac{\bar{\lambda}g}{3\epsilon} \times \begin{cases} 6h(\eta h - h^2), & \text{If } \eta \leq h \\ 2(\eta^3 - h^3), & \text{If } h \leq \eta. \end{cases} \quad (4.5.5)$$

The above definition of  $\Gamma$  implies that

$$\partial_h \tilde{p}(h, \eta) = gh \frac{\bar{\lambda}}{3\epsilon} \times \begin{cases} 6h + 12(h - \eta), & \text{If } \eta \leq h, \\ 6h, & \text{If } h \leq \eta. \end{cases} \quad (4.5.6)$$

Hence the hyperbolicity condition (4.5.2) is satisfied for any pair  $(\eta, h) \in \mathbb{R} \times \mathbb{R}_+$  and the model is compatible with dry states.

*Remark 4.5.2* (Infinitely many models). The  $\Gamma$  and  $\Phi$  functionals chosen above are seemingly the most efficient, but they are not unique. Technically, there exists a family of infinitely many relaxed models assuming  $\Gamma$  and  $\Phi$  satisfy their respective constraints.

### 4.5.3 Derivation of energy inequality

In this section, we derive an energy inequality for the hyperbolic relaxed model (4.4.1). The key results of this section are (i) Proposition 4.5.3 which a relaxed counterpart to Proposition 2.3.6; (ii) Corollary 4.5.5 which shows that the extra energy induced by the relaxation is small and positive.

Let  $\mathbf{u} := (h, \mathbf{q}, q_1, q_2, q_3)^\top$  be a smooth solution to (4.4.1a)–(4.4.1h) and assume  $\epsilon$  is constant. To simplify the presentation, we use the primitive representation of the auxiliary variables  $q_1 := h\eta$ ,  $q_2 := h\omega$ ,  $q_3 := h\beta$ . We first derive an equation for the potential energy. Recall this is done by multiplying the mass conservation equation (4.4.1a) with  $g(h + z)$ . Since the mass equation of the relaxed model is equivalent to that of the Serre model (2.3.1), the potential energy equation is the same as (2.3.14):

$$\partial_t \left( \frac{1}{2} gh^2 + ghz \right) + \nabla \cdot \left( \mathbf{v} \frac{1}{2} gh^2 \right) + \frac{1}{2} gh^2 (\nabla \cdot \mathbf{v}) + gz \nabla \cdot (h\mathbf{v}) = 0. \quad (4.5.7)$$

To derive the kinetic energy equation, we apply the operator  $(\mathbf{v} \cdot)$  to (4.4.1b):

$$\partial_t \left( \frac{1}{2} h \|\mathbf{v}\|^2 \right) + \nabla \cdot \left( \mathbf{v} \left( \frac{1}{2} h \|\mathbf{v}\|^2 \right) \right) + \mathbf{v} \cdot \nabla p_\epsilon(\mathbf{u}) = \frac{1}{2} s_\epsilon(\mathbf{u}) \mathbf{v} \cdot \nabla z - \frac{1}{4} \tilde{s}_\epsilon(\mathbf{u}) \mathbf{v} \cdot \nabla z.$$

We now find the contributions to the total energy induced by the auxiliary variables. We first



multiply (4.4.1c) by  $\frac{\bar{\lambda}g}{3\epsilon}h\Gamma'(\frac{\eta}{h}) := \frac{1}{3h}s_\epsilon(\mathbf{u})$ :

$$\frac{\bar{\lambda}g}{3\epsilon}h\Gamma'(\frac{\eta}{h})(\partial_t(h\eta) + \nabla \cdot (h\eta\mathbf{v})) = \frac{\bar{\lambda}g}{3\epsilon}h^2\Gamma'(\frac{\eta}{h})(\partial_t\eta + \mathbf{v} \cdot \nabla\eta).$$

We want to simplify the expression on the right hand side. By the chain rule, we have

$$\partial_t\Gamma(\frac{\eta}{h}) = \Gamma'(\frac{\eta}{h})(\frac{1}{h}\partial_t\eta - \frac{\eta}{h^2}\partial_t h), \quad \mathbf{v} \cdot \nabla\Gamma(\frac{\eta}{h}) = \Gamma'(\frac{\eta}{h})(\frac{1}{h}\mathbf{v} \cdot \nabla\eta - \frac{\eta}{h^2}\mathbf{v} \cdot \nabla h).$$

Using the above identities, we see

$$\begin{aligned} \frac{\bar{\lambda}g}{3\epsilon}h^2\Gamma'(\frac{\eta}{h})(\partial_t\eta + \mathbf{v} \cdot \nabla z) &= \frac{\bar{\lambda}g}{3\epsilon}h^3(\partial_t\Gamma(\frac{\eta}{h}) + \mathbf{v} \cdot \nabla\Gamma(\frac{\eta}{h})) + \frac{\bar{\lambda}g}{3\epsilon}\Gamma'(\frac{\eta}{h})h^2\frac{\eta}{h}(\partial_t h + \mathbf{v} \cdot \nabla h) \\ &= \frac{\bar{\lambda}g}{3\epsilon}h^3(\partial_t\Gamma(\frac{\eta}{h}) + \mathbf{v} \cdot \nabla\Gamma(\frac{\eta}{h})) + \frac{\bar{\lambda}g}{3\epsilon}\Gamma'(\frac{\eta}{h})\eta h(-h\nabla \cdot \mathbf{v}) \\ &= \frac{\bar{\lambda}g}{3\epsilon}h^3D_t(\Gamma(\frac{\eta}{h})) - \frac{\bar{\lambda}g}{3\epsilon}\Gamma'(\frac{\eta}{h})\eta h^2\nabla \cdot \mathbf{v} \\ &= \frac{\bar{\lambda}g}{3\epsilon}h^2(D_t(h\Gamma(\frac{\eta}{h})) + h\Gamma(\frac{\eta}{h})\nabla \cdot \mathbf{v}) - \frac{\bar{\lambda}g}{3\epsilon}\Gamma'(\frac{\eta}{h})\eta h^2\nabla \cdot \mathbf{v} \\ &= \frac{\bar{\lambda}g}{3\epsilon}h(D_t(h^2\Gamma(\frac{\eta}{h})) + 2h^2\Gamma(\frac{\eta}{h})\nabla \cdot \mathbf{v}) - \frac{\bar{\lambda}g}{3\epsilon}\Gamma'(\frac{\eta}{h})\eta h^2\nabla \cdot \mathbf{v} \\ &= \frac{\bar{\lambda}g}{3\epsilon}(D_t(h^3\Gamma(\frac{\eta}{h})) + \nabla \cdot (\mathbf{v}h^3\Gamma(\frac{\eta}{h}))) + \underbrace{\frac{\bar{\lambda}g}{3\epsilon}(2h^3\Gamma(\frac{\eta}{h}) - h^2\eta\Gamma(\frac{\eta}{h}))}_{\tilde{p}_\epsilon(\mathbf{u})}\nabla \cdot \mathbf{v} \\ &= \frac{\bar{\lambda}g}{3\epsilon}(D_t(h^3\Gamma(\frac{\eta}{h})) + \nabla \cdot (\mathbf{v}h^3\Gamma(\frac{\eta}{h}))) + \tilde{p}_\epsilon(\mathbf{u})\nabla \cdot \mathbf{v}. \end{aligned}$$

Notice that as mentioned in §4.3.1 the definition of the relaxed non-hydrostatic pressure  $\tilde{p}_\epsilon$  is important for the above simplification.

Multiplying the source terms of (4.4.1c) yields:

$$\frac{\bar{\lambda}g}{3\epsilon}h\Gamma'(\frac{\eta}{h})(h\omega - \frac{3}{2}\mathbf{q} \cdot \nabla z) = \frac{1}{3}s_\epsilon(\mathbf{u})\omega - \frac{1}{2}s_\epsilon(\mathbf{u})\mathbf{v} \cdot \nabla z.$$

Notice again, as mentioned in §4.3.1, replacing  $-\frac{3}{2}q_3$  by  $-\frac{3}{2}h\mathbf{v} \cdot \nabla z$  in (4.4.1c) is important here.

Without this substitution we would have  $-\frac{1}{2}s_\epsilon(\mathbf{u})\beta$  on the right-hand side of the above identity

instead of  $-\frac{1}{2}s_\epsilon(\mathbf{u})\mathbf{v}\cdot\nabla z$ . The above terms combine to the following equation:

$$\partial_t\left(\frac{\bar{\lambda}g}{3\epsilon}h^3\Gamma\left(\frac{\eta}{h}\right)\right) + \nabla\cdot\left(\frac{\bar{\lambda}g}{3\epsilon}h^3\Gamma\left(\frac{\eta}{h}\right)\mathbf{v}\right) + \tilde{p}_\epsilon(h,\eta)\nabla\cdot\mathbf{v} = \frac{1}{3}s_\epsilon(\mathbf{u})\omega - \frac{1}{2}s_\epsilon(\mathbf{u})\mathbf{v}\cdot\nabla z.$$

We continue by multiplying (4.4.1d) by  $\frac{1}{3}\omega$  and we obtain

$$\partial_t\left(\frac{1}{6}h\omega^2\right) + \nabla\cdot\left(\frac{1}{6}h\omega^2\mathbf{v}\right) = -\frac{1}{3}s_\epsilon(\mathbf{u})\omega.$$

Finally we multiply (4.4.1e) by  $\frac{1}{4}\beta$  and obtain

$$\partial_t\left(\frac{1}{8}h\beta^2\right) + \nabla\cdot\left(\frac{1}{8}h\beta^2\mathbf{v}\right) = \frac{1}{4}\tilde{s}_\epsilon(\mathbf{u})\beta.$$

Notice that the definition of  $\tilde{s}_\epsilon(\mathbf{u})$  and  $\Phi(\xi)$  gives

$$\frac{1}{4}\tilde{s}_\epsilon(\mathbf{u})(\beta - \mathbf{v}\cdot\nabla z) = -\frac{gh_0}{4}\frac{h}{\epsilon}\Phi\left(\frac{\mathbf{v}\cdot\nabla z - \beta}{\sqrt{gh_0}}\right)(\mathbf{v}\cdot\nabla z - \beta) \leq 0.$$

Combining the above details, we have the following energy inequality.

**Proposition 4.5.3** (Energy inequality for relaxed system). *Let  $\mathbf{u}$  be a smooth solution to (4.4.1a)–(4.4.1h). Then the following holds true:  $\partial_t\mathcal{E}_\epsilon(\mathbf{u}) + \nabla\cdot(\mathcal{F}_\epsilon(\mathbf{u})) = \frac{1}{4}\tilde{s}_\epsilon(\mathbf{u})(\beta - \mathbf{v}\cdot\nabla z) \leq 0$ , with*

$$\mathcal{E}_\epsilon(\mathbf{u}) := \frac{1}{2}gh^2 + gzh + \frac{1}{2}h\|\mathbf{v}\|^2 + \frac{1}{6}h\omega^2 + \frac{1}{8}h\beta^2 + \frac{\bar{\lambda}g}{3\epsilon}h^3\Gamma\left(\frac{\eta}{h}\right), \quad (4.5.8a)$$

$$\mathcal{F}_\epsilon(\mathbf{u}) := \mathbf{v}(\mathcal{E}_\epsilon(\mathbf{u}) + p_\epsilon(\mathbf{u})). \quad (4.5.8b)$$

*Remark 4.5.4* ((4.5.8a) vs. (2.3.15a)). By comparing the expression (4.5.8a) to (2.3.15a), and recalling that  $\omega$  is meant to be an approximation for  $\dot{h} + \frac{3}{2}\mathbf{v}\cdot\nabla z$  and  $\beta$  an approximation for  $\mathbf{v}\cdot\nabla z$ , we see that  $\frac{1}{6}h\omega^2 + \frac{1}{8}h\beta^2 = \frac{1}{6}h(\omega^2 + \frac{3}{4}\beta^2)$  in (4.5.8a) is the approximation of  $\frac{1}{6}h\left(\left(\dot{h} + \frac{3}{2}(\mathbf{v}\cdot\nabla z)\right)^2 + \frac{3}{4}(\mathbf{v}\cdot\nabla z)^2\right)$  in (2.3.15a). Notice that we observe an extra term in the energy  $h^3\Gamma\left(\frac{\eta}{h}\right)$ , but this was

shown in Remark 4.3.2 to be a small, positive quantity.

**Corollary 4.5.5** (Bound on relaxation induced energy term). *Let  $\mathbf{u} := (\mathbf{h}, \mathbf{q}, q_1, q_2, q_3)^\top$  be a smooth solution to (4.4.1a)–(4.4.1h) and assume  $\epsilon$  is constant (i.e., does not depend on  $\mathbf{x}$  and  $t$ ). Let  $T > 0$  be some final time. Assume that the boundary conditions for  $\mathbf{u}$  are such that  $\mathcal{F}(\mathbf{u}) \cdot \mathbf{n}_{\partial D} = 0$  for all  $t \in (0, T)$ . Also assume that the frame of reference for the model problem is such that the topography is positive (i.e.,  $z(\mathbf{x}) > 0$ ). Then, there is a  $c(\mathbf{u}_0)$  such that:*

$$\int_D (\mathbf{h}^3 \Gamma(\frac{\eta}{\mathbf{h}}))|_{t=T} \leq c(\mathbf{u}_0)\epsilon. \quad (4.5.9)$$

*Proof.* Integrating  $\partial_t \mathcal{E}_\epsilon(\mathbf{u}) + \nabla \cdot \mathcal{F}_\epsilon(\mathbf{u}) \leq 0$  over the spatial domain and applying the boundary condition  $\mathcal{F}(\mathbf{u}) \cdot \mathbf{n}_{\partial D} = 0$  gives:

$$\int_D \partial_t \mathcal{E}_\epsilon(\mathbf{u}) \leq 0.$$

Then, integrating again over the time interval  $(0, T)$  yields:

$$\begin{aligned} \int_D (\mathcal{E}_\epsilon(\mathbf{u}(T)) - \mathcal{E}_\epsilon(\mathbf{u}(0))) &\leq 0, \\ &\implies \\ \int_D \mathcal{E}_\epsilon(\mathbf{u}(T)) &\leq c(\mathbf{u}_0). \end{aligned}$$

Then, by dropping the positive quantities  $(\frac{1}{2}gh^2 + gzh + \frac{1}{2}\mathbf{h}\|\mathbf{v}\|^2 + \frac{1}{6}\mathbf{h}\omega^2 + \frac{1}{8}\mathbf{h}\beta^2)_{t=T}$ , we have

$$\int_D \frac{\bar{\lambda}g}{3\epsilon} \mathbf{h}^3 \Gamma(\frac{\eta}{\mathbf{h}}) \leq c(\mathbf{u}_0).$$

The result follow immediately. □

*Remark 4.5.6* (Energy consistent approximation). The previous result shows that the extra energy term induced by the relaxation vanishes as  $\epsilon \rightarrow 0$ . That is to say, the energy functional (4.5.8a) converges to the Serre energy functional (2.3.15a). Thus, our hyperbolic relaxed model is a consistent approximation of the original Serre equations. In Figure 4.2 we show an updated version of



topography effects (3.2.4). Then,

$$\begin{aligned}
\partial_t \mathcal{E}_{\text{flat}}(\mathbf{u}) + \nabla \cdot (\mathcal{F}_{\text{flat}}(\mathbf{u})) &= (\nabla_{\mathbf{u}} \mathcal{E}_{\text{flat}})^\top \partial_t \mathbf{u} + \nabla \cdot (\mathcal{F}_{\text{flat}}(\mathbf{u})) \\
&= (\nabla_{\mathbf{u}} \mathcal{E}_{\text{flat}})^\top (-\nabla \cdot \mathbb{f}_{\text{flat}}(\mathbf{u}) + \mathbf{R}(\mathbf{u})) + \nabla \cdot (\mathcal{F}_{\text{flat}}(\mathbf{u})) \\
&= 0, \\
&\implies \\
\nabla \cdot (\mathcal{F}_{\text{flat}}(\mathbf{u})) &= (\nabla \mathcal{E}_{\text{flat}}(\mathbf{u}))^\top (\nabla \cdot \mathbb{f}_{\text{flat}}(\mathbf{u}) - \mathbf{R}(\mathbf{u})).
\end{aligned}$$

□

#### 4.5.4 Dispersion relation

In this section, we derive the dispersion relation for the hyperbolic relaxed model (4.4.1). We then compare the relaxed phase velocity to that of the original Serre model (2.3.1). For simplicity, we restrict ourselves to one space-dimension and assume the topography is flat. Notice that since the topography is flat, we can drop the dependency of the  $q_3$  variable (see: Remark 3.2.5).

We proceed similarly as in §2.3.5.3 by considering the linearized problem of (4.4.1) about a rest state  $\mathbf{u}_0 := (H_0, 0, H_0^2, 0)^\top$ :

$$\partial_t h + H_0 \partial_x u = 0, \tag{4.5.12a}$$

$$\partial_t u + g \partial_x h + 2 \frac{\bar{\lambda} g}{\epsilon} H_0 \partial_x (h - \eta) = 0, \tag{4.5.12b}$$

$$\partial_t \eta = \omega, \tag{4.5.12c}$$

$$\partial_t \omega = 6 \frac{\bar{\lambda} g}{\epsilon} (h - \eta). \tag{4.5.12d}$$

We look for solutions in the form of  $\mathbf{u}(x, t) = \mathbf{u}_A \exp(i(kx - \sigma t))$  where  $\mathbf{u}_A$  are the wave amplitudes,  $k$  is the wave number and  $\sigma$  is the wave frequency. Here  $i := \sqrt{-1}$  is the imaginary unit number. A direct substitution of the plane wave solution into (4.5.12) yields the following linear

system:

$$\begin{pmatrix} -\sigma & H_0 k & 0 & 0 \\ gk + 2\frac{\bar{\lambda}g}{\epsilon} H_0 k & -\sigma & -2\frac{\bar{\lambda}g}{\epsilon} k H_0 & 0 \\ 0 & 0 & i\sigma & 1 \\ -6\frac{\bar{\lambda}g}{\epsilon} & 0 & 6\frac{\bar{\lambda}g}{\epsilon} & -i\sigma \end{pmatrix} \begin{pmatrix} h_A \\ u_A \\ \eta_A \\ \omega_A \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

Assuming that the solution to the linear system is non-trivial, we get the following equation for  $\sigma$ :

$$\sigma^4 - \left( gH_0 + 6\frac{\bar{\lambda}g}{\epsilon} \left( \frac{1 + \frac{1}{3}H_0^2 k^2}{k^2} \right) \right) k^2 \sigma^2 + 6\frac{\bar{\lambda}g}{\epsilon} gH_0 k^2 = 0. \quad (4.5.13)$$

We want to solve this equation for  $\sigma$  to find the linear dispersion relation for the relaxed model (4.4.1).

However, we first show that this equation can be simplified to find the Serre dispersion relation.

Dividing the equation by  $6\frac{\bar{\lambda}g}{\epsilon} gH_0$  yields:

$$\frac{\epsilon}{6\bar{\lambda}g} \frac{1}{gH_0} \sigma^4 - \left( \frac{\epsilon}{6\bar{\lambda}g} + \left( \frac{1 + \frac{1}{3}H_0^2 k^2}{gH_0 k^2} \right) \right) k^2 \sigma^2 + k^2 = 0.$$

Taking the limit  $\epsilon \rightarrow 0$  the above equation yields:

$$- \left( \frac{1 + \frac{1}{3}H_0^2 k^2}{gH_0 k^2} \right) \sigma^2 k^2 + k^2 = 0.$$

Solving for  $\sigma$  in the above equation yields the Serre dispersion relation:

$$\sigma^2 = \frac{gH_0 k^2}{1 + \frac{1}{3}H_0^2 k^2},$$

which was previously derived in §2.3.5.3. Thus, the dispersion relation for the relaxed system converges to that of the original Serre model as  $\epsilon \rightarrow 0$ .

We continue by giving the expression for the square of the phase velocity  $c_p^2 := \left(\frac{\sigma}{k}\right)^2$ . This is

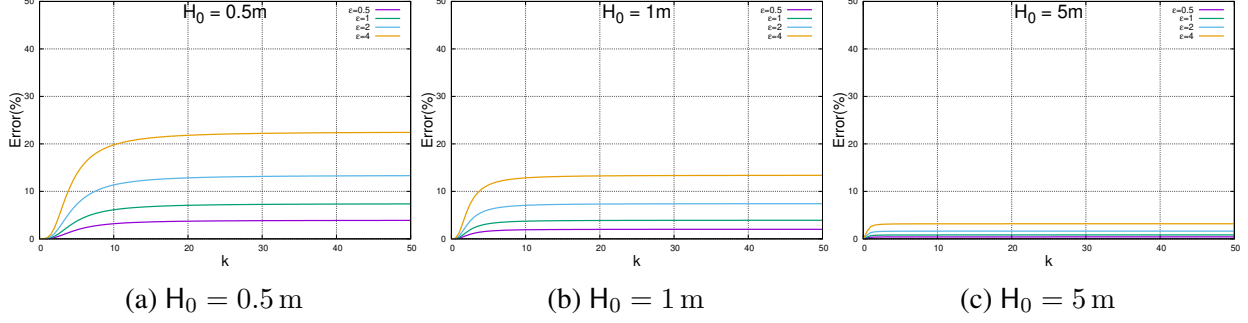


Figure 4.3: The error  $\frac{c_p^S - c_p^-}{c_p^S}$  as a function of the wave number  $k$  for different water depths.

found by solving the polynomial equation (4.5.13) and dividing by  $k^2$ :

$$(c_p^\pm)^2 = \frac{1}{2}gH_0 + 3\frac{\bar{\lambda}g}{\epsilon} \left( \frac{1 + \frac{1}{3}H_0^2k^2}{gH_0k^2} \right) \pm \frac{1}{2k} \sqrt{k^2 \left( gH_0 + 6\frac{\bar{\lambda}g}{\epsilon} \left( \frac{1 + \frac{1}{3}H_0^2k^2}{gH_0k^2} \right) \right)^2 - 24\frac{\bar{\lambda}g}{\epsilon}gH_0}. \quad (4.5.14)$$

The phase velocity above bears resemblance to the phase velocity derived in Favrie and Gavriluk [20, Eq. 19] for a different hyperbolic relaxed model. Just as in [20], we define a slow phase velocity (corresponding to  $(c_p^-)^2$ ) and a rapid phase velocity (corresponding to  $(c_p^+)^2$ ). In Figure 4.4, we plot the slow phase velocity  $c_p^-$  as a function of the wave number  $k$  for the values  $\bar{\lambda} = 1$ ,  $H_0 = 1$  m and  $\epsilon = \{1, 2, 4, 8\}$ . We see that as  $\epsilon$  decreases the relaxed slow phase velocity approaches that of the Serre model. In Figure 4.3, we show the error (in percent) between the relaxed slow phase velocity and the Serre phase velocity for several reference water depths:  $H_0 = \{0.5 \text{ m}, 1 \text{ m}, 5 \text{ m}\}$ . It can be seen that as the water depth decreases, we must decrease the relaxation parameter to capture the dispersive effects. In Figure 4.5, we plot the rapid phase velocity  $c_p^+$  for the same values. We see that  $c_p^+$  tends away from the Serre phase velocity as  $\epsilon$  decreases. As described in [20], this can be interpreted as the evolution of parasitic high-frequency waves induced by the relaxation. The above results are summarized in the following proposition.

**Proposition 4.5.9** (Linear dispersion relation for the relaxed model (4.4.1)). *The hyperbolic re-*

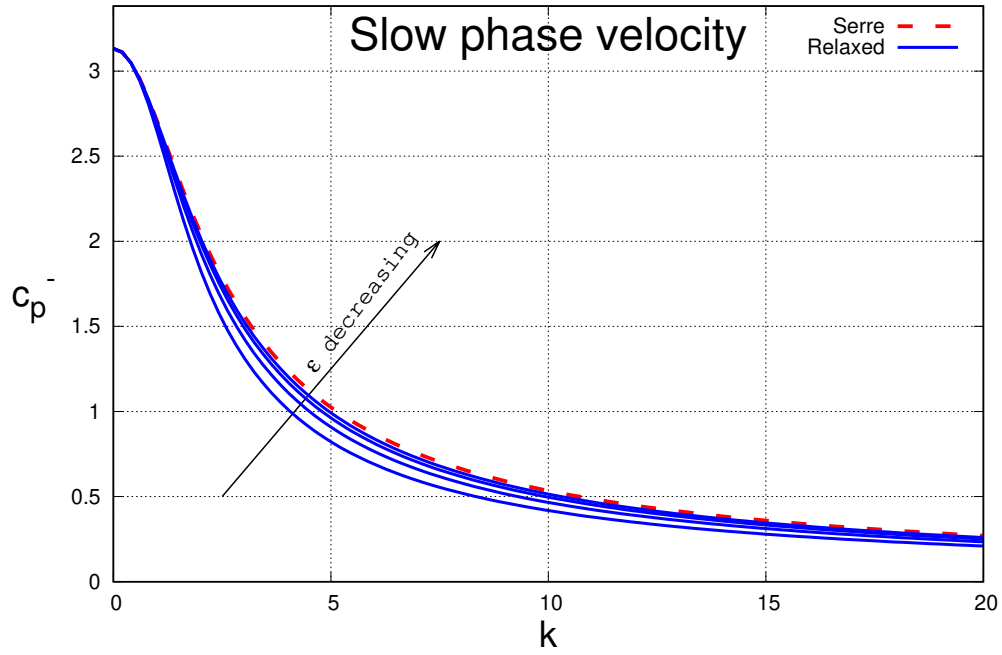


Figure 4.4: Comparison of slow phase velocity of hyperbolic relaxed model (4.4.1) with  $\epsilon = \{1, 2, 4, 8\}$  and Serre phase velocity (2.3.23).

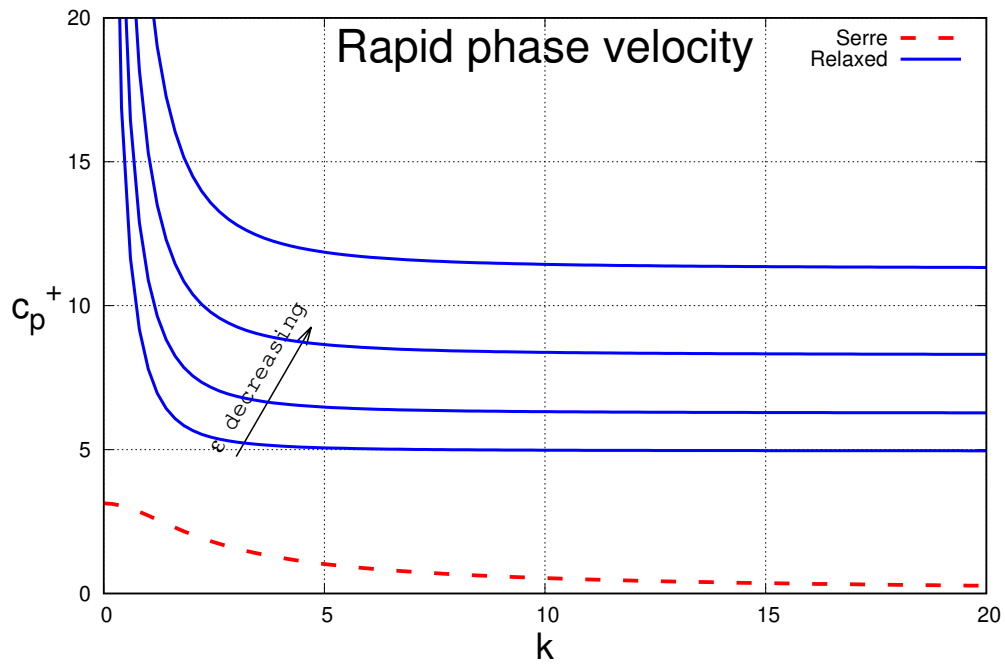


Figure 4.5: Comparison of rapid phase velocity of hyperbolic relaxed model (4.4.1) with  $\epsilon = \{1, 2, 4, 8\}$  and Serre phase velocity (2.3.23).



laxed system (4.4.1) admits the linear (squared) phase velocity  $(c_p^\pm)^2 := (\frac{\sigma}{k})^2$ :

$$(c_p^\pm)^2 = \frac{1}{2}gH_0 + 3\frac{\bar{\lambda}g}{\epsilon} \left( \frac{1 + \frac{1}{3}H_0^2k^2}{gH_0k^2} \right) \pm \frac{1}{2k} \sqrt{k^2 \left( gH_0 + 6\frac{\bar{\lambda}g}{\epsilon} \left( \frac{1 + \frac{1}{3}H_0^2k^2}{gH_0k^2} \right) \right)^2 - 24\frac{\bar{\lambda}g}{\epsilon}gH_0}.$$

Additionally, the phase velocities  $(c_p^\pm)$  satisfy the following condition:

$$0 < c_p^-(\epsilon) < c_p^S < c_p^+(\epsilon),$$

where  $c_p^S$  is the Serre phase velocity (2.3.23).

## 4.6 Literature review

In this section, we discuss similar hyperbolic relaxation techniques for solving the Serre Equations seen in the literature. In particular, we give special attention to the work of Favrie and Gavrilyuk [20] which was the motivation for the hyperbolic relaxation technique described above (and many other relaxation techniques). We note that the technique described in [20] was done in one spatial dimension  $\mathbb{R}^1$  and extended in Tkachenko [60] for  $\mathbb{R}^2$ . The other relaxation techniques discussed are those that are extensions to the reformulations discussed in §3.3.

### 4.6.1 Favrie and Gavrilyuk model

The relaxation approach introduced in [20] is based on a Hamiltonian mechanics point of view. The authors relax the constraint  $q_1 := h^2$  (or equivalently  $\eta := h$ ) by modifying the Lagrangian of the Serre model:

$$\mathcal{L}(h, \mathbf{q}) = \int_D \left( \left( \frac{1}{2}h^{-1}\|\mathbf{q}\|^2 + \frac{1}{6}hh^2 \right) - \frac{1}{2}gh^2 \right) dx,$$

with an augmented Lagrangian:

$$\widehat{\mathcal{L}}(h, \mathbf{q}) = \int_D \left( \frac{1}{2}h^{-1}\|\mathbf{q}\|^2 - \frac{1}{2}gh^2 + \frac{1}{6}hh^2 - \frac{1}{6}\bar{\lambda}\left(\frac{\eta}{h} - 1\right)^2 \right) dx,$$

and deriving a new system from this augmented Lagrangian. Recall that the Lagrangian of a mechanical system can be thought of as a functional measuring the difference of the system's

kinetic energy and potential energy.

The relaxed system (with a flat-bottom) derived from the augmented Lagrangian (see [20, 60] for details) is of the form:

$$\partial_t \mathbf{h} + \nabla \cdot (\mathbf{h} \mathbf{v}) = 0 \quad (4.6.1a)$$

$$\partial_t \mathbf{q} + \nabla \cdot (\mathbf{v} \otimes \mathbf{q}) + \nabla \cdot \left( \frac{1}{2} g \mathbf{h}^2 - \frac{1}{3} \eta \bar{\lambda} \left( \frac{\eta}{\mathbf{h}} - 1 \right) \right) = 0 \quad (4.6.1b)$$

$$\partial_t q_1 + \nabla \cdot (q_1 \mathbf{v}) = q_2 \quad (4.6.1c)$$

$$\partial_t q_2 + \nabla \cdot (q_2 \mathbf{v}) = -\bar{\lambda} \left( \frac{\eta}{\mathbf{h}} - 1 \right) \quad (4.6.1d)$$

By directly comparing this model to the one defined by (4.4.1), we observe that the Hamiltonian strategy proposed in [20] can be re-interpreted of as:

- (i) Define the pressure as  $p := \frac{1}{2} g \mathbf{h}^2 - \frac{1}{3} \frac{q_1}{\mathbf{h}} s = \frac{1}{2} g \mathbf{h}^2 - \frac{1}{3} \eta s$ ;
- (ii) Replacing  $s$  by  $-\bar{\lambda} \left( \frac{\eta}{\mathbf{h}} - 1 \right)$  to relax the constraint  $q_1 := \mathbf{h}^2$  where  $\bar{\lambda} > 0$  is the relaxation parameter which scales like the square of a velocity.

This relaxed model is shown in [20] to be hyperbolic, but not compatible with dry states.

## 5. APPROXIMATION

### 5.1 Introduction

In this chapter, we introduce a space/time approximation of the Saint-Venant model (described in Chapter 2) and the hyperbolic Serre model (described in Chapter 4). The approximation in space is done using continuous finite elements with explicit time stepping. The goal of this chapter is to introduce a high-order (accurate in space) invariant-domain-preserving method that is well-balanced (with respect to rest states) for solving the two models with external physical sources. The approximation is carried out as follows. We first introduce a low-order (i.e., first-order accurate in space) method that is robust with respect to dry states and well-balanced with respect to rest states. We then introduce a provisional high-order method (i.e., second-order accurate in space) using the entropy viscosity methodology that may violate the invariant-domain-preserving property. We then describe a convex limiting technique used for correcting (i.e., limiting) the provisional higher-order method. This is done by enforcing local physical bounds that are naturally satisfied by the low-order method. We give an emphasis on how to apply the convex limiting methodology to hyperbolic systems with sources since it is not well-documented in the literature. The key results of this chapter are: Proposition 5.4.2 and Proposition 5.4.3 which prove the low-order method is invariant-domain-preserving and well-balanced for each model; Theorem 5.7.8 which proves the limited high-order method is invariant-domain-preserving for the Hyperbolic Serre model; and Proposition 5.7.9 which proves that the higher-order method is well-balanced. We note that the algorithms introduced in this chapter are presented in a generic fashion (i.e., independent of the model) unless otherwise specified.

The chapter is organized as follows. In Section 5.2, we introduce the continuous finite element setting for approximating the models in space. In Section 5.2.1, we introduce the finite element representation of the approximate solution. Then, in Section 5.3, we describe the low-order approximation that is invariant-domain-preserving. The details on how to construct the graph-viscosity

for both models is given in this section along with a discussion on the discretization of the external physical sources. In Section 6.4, we give a discussion on the invariant domain preserving and well-balancing properties of the low-order method. Then, in Section 5.6, we introduce a provisional higher-order method that is more accurate, but may violate the invariant domain preserving property. Finally, in Section 5.7, we describe how to apply the convex limiting technique to the provisional higher-order method.

## 5.2 Finite element setting

In this section, we introduce the continuous finite element setting used for the approximation of the Saint-Venant model (2.2.1) and the hyperbolic Serre model (4.4.1). We note that the techniques shown here can be also be adapted using discontinuous finite elements and finite volumes as discussed in Guermond et al. [32].

Let  $(\mathcal{T}_h)_{h>0}$  be a shape-regular family of matching meshes where  $h$  can be thought of as the typical mesh-size. Let  $(\widehat{K}, \widehat{P}, \widehat{\Sigma})$  be a reference finite element. For every cell  $K \in \mathcal{T}_h$  we denote by  $T_K : \widehat{K} \rightarrow K$  the geometric bijective transformation that maps the reference element  $\widehat{K}$  to the current element  $K$ . Let  $\mathcal{V}$  denote the set of integers that represent the degrees of freedom. Given some mesh  $\mathcal{T}_h$ , we consider a scalar-valued finite element space  $P(\mathcal{T}_h)$  with positive global shape functions  $\{\varphi_i\}_{i \in \mathcal{V}}$  associated with the Lagrange nodes  $\{\mathbf{a}_i\}_{i \in \mathcal{V}}$ :

$$P(\mathcal{T}_h) := \{v \in C^0(D; \mathbb{R}) \mid v|_K \circ T_K \in \widehat{P}, \forall K \in \mathcal{T}_h\}. \quad (5.2.1)$$

Note that  $\dim(P(\mathcal{T}_h)) := \text{card}(\mathcal{V})$ . The approximation in space of the conserved variable  $\mathbf{u}$  is done in the space of  $\mathbb{R}^m$ -valued finite elements  $\mathbf{P}(\mathcal{T}_h) := [P(\mathcal{T}_h)]^m$  where  $m$  is the PDE system size. Letting  $d$  represent the spatial dimension, we have that  $m := d + 1$  for the Saint-Venant model or  $m := d + 4$  for the hyperbolic Serre model. For both models, the bottom topography  $z(\mathbf{x})$  is approximated in  $P(\mathcal{T}_h)$ . For every  $i \in \mathcal{V}$ , we call the stencil of the shape function,  $\varphi_i$ , the index set

$$\mathcal{I}(i) := \{j \in \mathcal{V} \mid \text{supp}(\varphi_i) \cap \text{supp}(\varphi_j) \neq \emptyset\}. \quad (5.2.2)$$

We also define  $\mathcal{I}^*(i) := \mathcal{I}(i) \setminus \{i\}$ . The following mesh-dependent quantities play an important role for the space and time approximation:

$$m_{ij} := \int_D \varphi_i(\mathbf{x}) \varphi_j(\mathbf{x}) \, dx, \quad m_i := \int_D \varphi_i(\mathbf{x}) \, dx, \quad (5.2.3a)$$

$$\mathbf{c}_{ij} := \int_D \varphi_i \nabla \varphi_j \, dx, \quad \mathbf{n}_{ij} := \frac{\mathbf{c}_{ij}}{\|\mathbf{c}_{ij}\|_{\ell^2}}, \quad (5.2.3b)$$

where  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$ . Here  $m_{ij}$  represents the entries of the consistent mass matrix and  $m_i$  the entries of the lumped mass matrix. By the partition of unity property:

$$\sum_{i \in \mathcal{V}} \varphi_i = 1,$$

we have that  $m_i = \sum_{j \in \mathcal{I}(i)} m_{ij}$ . The following three properties are essential to establish conservation of the numerical scheme:

- (i)  $\sum_{i \in \mathcal{V}} \mathbf{c}_{ij} = \mathbf{0}$  (partition of unity property);
- (ii)  $\mathbf{c}_{ij} = -\mathbf{c}_{ji}$  if either  $\varphi_i$  or  $\varphi_j$  is zero on  $\partial D$  (integration by parts);
- (iii)  $\sum_{i \in \mathcal{V}} \mathbf{c}_{ij} = \mathbf{0}$  if  $\varphi_j$  is zero on  $\partial D$  (partition of unity property).

The quantity  $\mathbf{c}_{ij}$  is used to approximate the gradient (or divergence) operator. Let us expand on this. Let  $f(\mathbf{x})$  be a sufficiently smooth scalar function and let  $\sum_{i \in \mathcal{V}} F_i \varphi_i(\mathbf{x})$  be its finite element approximation. Then, we define its discrete gradient as follows:  $(\nabla f)_i := \sum_{j \in \mathcal{I}(i)} F_j \mathbf{c}_{ij}$ . Now, let  $\mathbf{f}(\mathbf{x})$  be a sufficiently smooth vector-valued function and let  $\sum_{i \in \mathcal{V}} \mathbf{F}_i \varphi_i(\mathbf{x})$  be its finite element approximation. Then we define the approximate divergence of  $\mathbf{f}$  as:  $(\nabla \cdot \mathbf{f})_i := \sum_{j \in \mathcal{I}(i)} \mathbf{F}_j \cdot \mathbf{c}_{ij}$

*Remark 5.2.1* (Finite element spaces). In the numerical simulations reported in Chapter 6, we use the following finite element spaces:

- (i) The linear, continuous  $\mathbb{P}_1$  Lagrange finite elements on unstructured, Delaunay meshes.
- (ii) The bi-linear, continuous  $\mathbb{Q}_1$  finite elements on quadrangular meshes.

### 5.2.1 Finite element representations

Since we would like to present the space/time approximation in a model-independent fashion, we first represent quantities and definitions in a generic way. Then, when appropriate, we give explicit definitions for each respective model. The finite element approximation of the conserved variable  $\mathbf{u}$  at time  $t$  is denoted by  $\mathbf{u}_h(t)$  where

$$\mathbf{u}_h(t) = \sum_{i \in \mathcal{V}} \mathbf{U}_i(t) \varphi_i \in \mathbf{P}(\mathcal{T}_h).$$

We denote by  $z_h := \sum_{i \in \mathcal{V}} Z_i \varphi_i \in P(\mathcal{T}_h)$  the approximation of the topography map. Let  $(\nabla Z)_i := \sum_{j \in \mathcal{I}(i)} Z_j \mathbf{c}_{ij}$  denote the approximate gradient of the topography.

Let  $H_{0,\max}$  be some reference scale for the water depth. For instance we can take  $H_{0,\max} := \text{ess sup}_{\mathbf{x} \in D} h_0(\mathbf{x})$ , where  $h_0(\mathbf{x})$  is the initial water depth profile. We introduce the regularized water depth quantity  $H^\delta$ :

$$H_i^\delta := \left( \frac{2H_i}{H_i^2 + \max(H_i, \delta H_{0,\max})^2} \right)^{-1} \quad (5.2.4)$$

where  $\delta$  is a small dimensionless parameter. The quantity  $H^\delta$  is used for defining the approximate primitive variables and is used to make sense of  $1/H^\delta$ . For example, the approximate velocity is defined by  $\mathbf{v}_h := \sum_{i \in \mathcal{V}} \mathbf{V}_i \varphi_i$  with

$$\mathbf{V}_i := \frac{\mathbf{Q}_i}{H_i^\delta} \quad (5.2.5)$$

Notice that this regularization is active only when dry state occurs. For example,  $\mathbf{V}_i := \frac{1}{H_i} \mathbf{Q}_i$  if  $H_i \geq \delta H_{0,\max}$ . When  $0 \leftarrow H_i < \delta H_{0,\max}$ , we see that  $1/H^\delta$  approaches 0; thus avoiding the division by 0. Though this may seem un-physical, when  $H_i$  is small the momentum  $\mathbf{Q}_i$  will be close to 0. The reader is referred to Kurganov and Petrova [43, Eq. (2.17)], Chertock et al. [13, Eq. (3.10)], and [5, 28, §5.1], where this technique is also adopted.

*Remark 5.2.2 (Setting the  $\delta$  quantity).* We set  $\delta = 10^{-5}$  in the numerical illustrations described in Chapter 6 for both models. We note that it is possible to take  $\delta$  to be smaller for each model. However, for the hyperbolic Serre model, a smaller  $\delta$  requires a smaller CFL number.

*Remark 5.2.3* (Saint-Venant representation). The finite element representation for the Saint-Venant model is given by  $\mathbf{u}_h(t) := (\mathbf{h}_h, \mathbf{q}_h)^\top$  with

$$\mathbf{U}_i(t) := \left( H_i(t), \mathbf{Q}_i(t) \right)^\top. \quad (5.2.6)$$

*Remark 5.2.4* (Hyperbolic Serre representation). The finite element representation for the hyperbolic Serre model is given by:  $\mathbf{u}_h(t) := (\mathbf{h}_h, \mathbf{q}_h, q_{1,h}, q_{2,h}, q_{3,h})^\top$  with

$$\mathbf{U}_i(t) := (H_i(t), \mathbf{Q}_i(t), Q_{1,i}(t), Q_{2,i}(t), Q_{3,i}(t))^\top. \quad (5.2.7)$$

and approximate auxiliary quantities  $\eta_h, \omega_h, \beta_h$  defined with regularization for all  $i \in \mathcal{V}$ :

$$\mathbf{N}_i := \frac{Q_{1,i}}{H_i^\delta}, \quad \mathbf{W}_i := \frac{Q_{2,i}}{H_i^\delta}, \quad \mathbf{B}_i := \frac{Q_{3,i}}{H_i^\delta}, \quad (5.2.8)$$

The relaxation parameter  $\epsilon$  is chosen to be proportional to the local mesh-size. Recalling that  $m_i := \int_D \varphi_i dx$  is proportional to the volume of the support of the shape function  $\varphi_i$ , we set  $\epsilon_h := \sum_{i \in \mathcal{V}} \mathcal{E}_i \varphi_i$  with  $\mathcal{E}_i := m_i^{\frac{1}{d}}$  (recall that  $d \in \{1, 2\}$  is the space dimension).

### 5.3 The low-order method

In this section, we describe the low-order approximation of the shallow water flow models (2.2.1) and (2.3.1) using the finite element setting shown above. The method is formally first-order accurate in space and is presented with forward Euler time stepping. Higher-order accuracy in time is achieved by using any explicit strong stability preserving Runge-Kutta method (more on this is discussed in §5.3.8).

#### 5.3.1 Numerical flux and hydrostatic pressure/source

Let  $\mathbf{u}_h^0 := \sum_{i \in \mathcal{V}} \mathbf{U}_i^0 \varphi_i \in \mathbf{P}(\mathcal{T}_h)$  be some reasonable approximation of the initial data  $\mathbf{u}_0$  at time  $t^0$ . Let  $n \in \mathbb{N}$  and  $t^n$  be the current time. Let  $\mathbf{u}_h^n := \sum_{i \in \mathcal{V}} \mathbf{U}_i^n \varphi_i \in \mathbf{P}(\mathcal{T}_h)$  be the current admissible approximation of  $\mathbf{u}$ . Recall that the admissible set for both models is:  $\mathcal{A} = \{\mathbf{u} \mid \mathbf{h} > 0\}$ .

We define the numerical gas dynamics flux  $\nabla \cdot (\mathbf{v} \otimes \mathbf{q})$  for all  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$ :

$$\mathbf{G}_{ij}^n := \mathbf{U}_j^n (\mathbf{V}_j^n \cdot \mathbf{c}_{ij}) \quad (5.3.1)$$

We now want to discretize the hydrostatic part of the pressure  $\nabla \frac{1}{2}gh^2$  and topography source term  $gh\nabla z$ . To ensure well-balancing, we combine the two terms as follows:

$$\begin{aligned} \nabla \frac{1}{2}gh^2 + gh\nabla z &= gh\nabla h + gh\nabla z \\ &= gh\nabla h + gh\nabla z \\ &= gh\nabla(h + z). \end{aligned}$$

Notice that the quantity  $gh\nabla(h + z)$  shows up directly in the lake-at-rest problem for both models (see: Prop. 2.2.6). We discretize  $gh\nabla(h + z)$  as follows:

$$\mathbf{P}_{ij}^{\text{SV}} = g\mathbf{H}_i^n (\mathbf{H}_j^n + Z_j) \mathbf{c}_{ij} \quad (5.3.2)$$

*Remark 5.3.1* (Saint-Venant numerical flux). For all  $i \in \mathcal{V}$  and all  $j \in \mathcal{I}(i)$ , we define the Saint-Venant numerical flux as follows:

$$\mathbf{F}_{ij}^n := \mathbf{G}_{ij}^n + (0, \mathbf{P}_{ij}^{\text{SV}})^\top. \quad (5.3.3)$$

*Remark 5.3.2* (Hyperbolic Serre numerical flux). For all  $i \in \mathcal{V}$  and all  $j \in \mathcal{I}(i)$ , we define the Hyperbolic Serre numerical flux as follows:

$$\mathbf{F}_{ij}^n := \mathbf{G}_{ij}^n + \left(0, \mathbf{P}_{ij}^{\text{SW}} + \tilde{\mathbf{P}}(\mathbf{U}_j^n) \mathbf{c}_{ij}, 0, 0, 0\right)^\top, \quad (5.3.4a)$$

$$\tilde{\mathbf{P}}(\mathbf{U}) := -\frac{\bar{\lambda}g}{3\mathcal{E}} \times \begin{cases} 6H(Q_1 - H^2), & \text{if } Q_1 \leq H^2 \\ 2\frac{(Q_1 - H^2)}{H^3} (N^2 + Q_1 + H^2), & \text{if } H^2 < Q_1. \end{cases} \quad (5.3.4b)$$



where  $\tilde{P}(\mathbf{U})$  represents the numerical relaxed non-hydrostatic pressure.

### 5.3.2 Well-balancing star states

We now define the hydrostatic reconstructed star states  $\mathbf{U}_i^{*,j,n}$  which ensure that the numerical method is well-balanced with respect to rest states. The concept of well-balancing (originally introduced for the Saint-Venant model) has roots in the work of Bermúdez and Vázquez [9, Def. 1] and Greenberg and Le Roux [27]. In this work, we follow the ideas presented in Audusse et al. [4] (and used in [5, 28]) and define the hydrostatic reconstructed water depth for all  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$  as follows:

$$H_i^{*,j,n} := \max(0, H_i^n + Z_i - \max(Z_i, Z_j)). \quad (5.3.5)$$

The hydrostatic reconstructed water depth  $H_i^{*,j,n}$  is used to define the full reconstructed state  $\mathbf{U}_i^{*,j,n}$  for each model below.

*Remark 5.3.3* (Saint-Venant star states). We define the hydrostatic reconstructed star states for the Saint-Venant model for all  $i \in \mathcal{V}$  and all  $j \in \mathcal{I}(i)$ :

$$\mathbf{U}_i^{*,j,n} := \frac{H_i^{*,j,n}}{H_i^\delta} (H_i^n, \mathbf{Q}_i^n)^\top, \quad (5.3.6a)$$

$$\mathbf{U}_j^{*,i,n} := \frac{H_j^{*,i,n}}{H_j^\delta} (H_j^n, \mathbf{Q}_j^n)^\top, \quad (5.3.6b)$$

*Remark 5.3.4* (Hyperbolic Serre star states). We define the Hyperbolic Serre hydrostatic reconstructed star states for all  $i \in \mathcal{V}$  and all  $j \in \mathcal{I}(i)$ :

$$\mathbf{U}_i^{*,j,n} := \frac{H_i^{*,j,n}}{H_i^\delta} \left( H_i^n, \mathbf{Q}_i^n, \frac{H_i^{*,j,n}}{H_i^\delta} Q_{1,i}^n, Q_{2,i}^n, Q_{3,i}^n \right)^\top, \quad (5.3.7a)$$

$$\mathbf{U}_j^{*,i,n} := \frac{H_j^{*,i,n}}{H_j^\delta} \left( H_j^n, \mathbf{Q}_j^n, \frac{H_j^{*,i,n}}{H_j^\delta} Q_{1,j}^n, Q_{2,j}^n, Q_{3,j}^n \right)^\top, \quad (5.3.7b)$$

Note that the star state for the quantity  $Q_1$  in (5.3.7) is defined with the square of the ratio  $\left(\frac{H_i^{*,j,n}}{H_i^\delta}\right)^2$ ; that is  $Q_{1,i}^{*,j,n} := \left(\frac{H_i^{*,j,n}}{H_i^\delta}\right)^2$ . This is crucial for well-balanced and is seen in the proof of Proposi-

tion 5.4.3. We expand on this in Remark 5.4.1.

### 5.3.3 PDE source

We now discretize the PDE source for both models.

*Remark 5.3.5* (Saint-Venant PDE source). Recall  $\mathbf{R}(\mathbf{u}, \nabla z) = gh\nabla z$  for the Saint-Venant model.

However, since we discretized the hydrostatic pressure and source together in (5.3.2), we set

$$\mathbf{R}_i^n = \mathbf{0}. \quad (5.3.8)$$

*Remark 5.3.6* (Hyperbolic Serre PDE source). To approximate the PDE source term  $\mathbf{R}(\mathbf{u}, \nabla z)$  in (4.4.1), we introduce the following quantities:

$$\mathbf{R}_1(\mathbf{U}, \nabla Z) := \mathbf{Q}_2 - \frac{3}{2}\mathbf{Q} \cdot \nabla Z, \quad (5.3.9a)$$

$$\mathbf{R}_2(\mathbf{U}) := \frac{\bar{\lambda}g}{\mathcal{E}} \times \begin{cases} 6(\mathbf{Q}_1 - \mathbf{H}^2), & \text{if } \mathbf{Q}_1 \leq \mathbf{H}^2, \\ 6\mathbf{N} \frac{(\mathbf{Q}_1 - \mathbf{H}^2)}{\mathbf{H}^{\delta}}, & \text{if } \mathbf{H}^2 < \mathbf{Q}_1, \end{cases} \quad (5.3.9b)$$

$$\mathbf{R}_3(\mathbf{U}, \nabla Z) := \frac{\bar{\lambda}}{\mathcal{E}} \sqrt{g\mathbf{H}_{0,\max}} (\mathbf{Q} \cdot \nabla Z - \mathbf{Q}_3). \quad (5.3.9c)$$

Then, the Hyperbolic Serre discrete PDE source is defined by

$$\mathbf{R}_i^n := \left( 0, \left( \frac{1}{2}\mathbf{R}_2(\mathbf{U}_i^n) - \frac{1}{4}\mathbf{R}_3(\mathbf{U}_i^n, \frac{1}{m_i}(\nabla Z)_i) \right) \frac{1}{m_i}(\nabla Z)_i, \right. \\ \left. \mathbf{R}_1(\mathbf{U}_i^n, \frac{1}{m_i}(\nabla Z)_i), -\mathbf{R}_2(\mathbf{U}_i^n), \mathbf{R}_3(\mathbf{U}_i^n, \frac{1}{m_i}(\nabla Z)_i) \right)^{\top}. \quad (5.3.10)$$

### 5.3.4 External physical sources

We discretize the Gauckler-Manning friction source term (2.4.1) by setting

$$\mathbf{S}_F(\mathbf{U}_i^n) := \frac{-2gn^2\mathbf{Q}_i^n \|\mathbf{V}_i\|_{\ell^2}}{(\mathbf{H}_i^n)^{\gamma} + \max((\mathbf{H}_i^n)^{\gamma}, 2gn^2\tau_n \|\mathbf{V}_i\|_{\ell^2})} \mathbb{e}_{\mathbf{q}}. \quad (5.3.11)$$

Note that a regularization for the quantity  $h^{-\gamma}$  was introduced to avoid division by 0 when  $h \rightarrow 0$ . This expression was introduced in Guermond et al. [28] and is shown to be stable under the usual hyperbolic CFL time step restriction, i.e., no iterations or semi-implicit time stepping is needed to advance in time with this definition.

We discretize the wave generation source (2.4.2) as follows:

$$\mathbf{S}_G(\mathbf{U}_i^n) := -\frac{\sqrt{gH_0}}{\mathcal{E}_i}(\mathbf{U}_i^n - \mathbf{u}_{\text{wave}}(\mathbf{a}_i, t^n))G\left(\frac{\mathbf{a}_i - x_{\min}}{L_{\text{gen}}}\right), \quad (5.3.12)$$

The absorption zone source term (2.4.3) is approximated as follows:

$$\mathbf{S}_A(\mathbf{U}_i^n) := -\frac{\sqrt{gH_0}}{\mathcal{E}_i}\mathbf{U}_{\text{abs},i}^n G\left(\frac{x_{\max} - \mathbf{a}_i}{L_{\text{abs}}}\right), \quad (5.3.13)$$

*Remark 5.3.7* (Saint-Venant absorption). The quantities enforced in the absorption zone for the Saint-Venant model are:

$$\mathbf{U}_{\text{abs},i}^n := (0, \mathbf{Q}_i^n)^\top.$$

*Remark 5.3.8* (Hyperbolic Serre absorption). The quantities enforced in the absorption zone for the hyperbolic Serre model are:

$$\mathbf{U}_{\text{abs},i}^n := (0, \mathbf{Q}_i^n, 0, \mathbf{Q}_{2,i}, 0)^\top.$$

### 5.3.5 Low-order graph-viscosity coefficients

We now define the low-order graph-viscosity coefficients  $d_{ij}^{\text{L},n}$  and  $\mu_{ij}^{\text{L},n}$  that make the numerical method invariant domain preserving and well-balanced. Here, the superindex <sup>L</sup> denotes that the quantities are associated with a (L)ow-order solution (i.e., first-order accurate in space).

For all  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$ , we set

$$\mu_{ij}^{L,n} := \max(|\mathbf{V}_i^n \cdot \mathbf{n}_{ij}| \|\mathbf{c}_{ij}\|_{\ell^2}, |\mathbf{V}_j^n \cdot \mathbf{n}_{ji}| \|\mathbf{c}_{ji}\|_{\ell^2}), \quad (5.3.14a)$$

$$d_{ij}^{L,n} := \max(\mu_{ij}^{L,n}, \max(\lambda_{ij}^n \|\mathbf{c}_{ij}\|_{\ell^2}, \lambda_{ji}^n \|\mathbf{c}_{ji}\|_{\ell^2})), \quad (5.3.14b)$$

such that

$$d_{ij}^{L,n} = d_{ji}^{L,n}, \mu_{ij}^{L,n} = \mu_{ji}^{L,n}, d_{ij}^{L,n} \geq \mu_{ij}^{L,n} \geq 0, \quad i \neq j.$$

The quantity  $\lambda_{ij}^n$  is some wave speed yet to be defined for each model.

*Remark 5.3.9* (Finding  $\lambda_{ij}^n$  for the Saint-Venant model). The wave speed  $\lambda_{ij}^n$  for the Saint-Venant model can be found by deriving a guaranteed maximum wave speed from the associated Riemann problem (see Azerad et al. [5, Sec. 3.3] and Guermond et al. [28, Sec. 4]). For the sake of brevity, we only give the exact formula here and not the derivation.

Let  $\partial_t \mathbf{w} + \partial_x(\mathbb{f}_{1D}(\mathbf{w})) = 0$  be the associated directional 1D Riemann problem with left states

$$(\mathbf{h}_L, \mathbf{q}_L \cdot \mathbf{n})^\top := (\mathbf{H}_i^n, \mathbf{Q}_i^n \cdot \mathbf{n}_{ij})^\top,$$

and right states

$$(\mathbf{h}_R, \mathbf{q}_R \cdot \mathbf{n})^\top := (\mathbf{H}_j^n, \mathbf{Q}_j^n \cdot \mathbf{n}_{ji})^\top,$$

where  $\mathbf{n}_{ij} := \frac{\mathbf{c}_{ij}}{\|\mathbf{c}_{ij}\|}$ . Let  $\mathbf{h}_{\min} = \min(\mathbf{h}_L, \mathbf{h}_R)$  and  $\mathbf{h}_{\max} = \max(\mathbf{h}_L, \mathbf{h}_R)$ . We set the left and right sound speeds as  $a_L = \sqrt{g\mathbf{h}_L}$  and  $a_R = \sqrt{g\mathbf{h}_R}$ , respectively. Then,  $\lambda_{ij}^n$  for the Saint-Venant model is defined as follows:

$$\lambda_{ij}^n = \max(\lambda_1^-(\bar{\mathbf{h}}_*), \lambda_1^+(\bar{\mathbf{h}}_*)) \quad (5.3.15)$$

with

$$\lambda_1^-(\mathbf{h}_*) = u_L - a_L \sqrt{\left(1 + \left(\frac{\mathbf{h}_* - \mathbf{h}_L}{2\mathbf{h}_L}\right)_+\right) \left(1 + \left(\frac{\mathbf{h}_* - \mathbf{h}_L}{\mathbf{h}_L}\right)_+\right)}, \quad (5.3.16a)$$

$$\lambda_2^-(\mathbf{h}_*) = u_R + a_R \sqrt{\left(1 + \left(\frac{\mathbf{h}_* - \mathbf{h}_R}{2\mathbf{h}_R}\right)_+\right) \left(1 + \left(\frac{\mathbf{h}_* - \mathbf{h}_R}{\mathbf{h}_R}\right)_+\right)}, \quad (5.3.16b)$$

$$\bar{h}_* = \begin{cases} \frac{1}{16g} (\max(0, u_L - u_R + 2a_L + 2a_R))^2, & \text{if } 0 \leq f(x_0 h_{\min}), \\ \sqrt{h_{\min} h_{\max}} \left(1 + \frac{\sqrt{2}(u_L - u_R)}{\sqrt{g h_{\min}} + \sqrt{g h_{\max}}}\right), & \text{if } f(x_0 h_{\max}) < 0, \\ \left(-\sqrt{2h_{\min}} + \sqrt{3h_{\min}} + 2\sqrt{2h_{\min} h_{\max}} + \sqrt{\frac{2}{g} h_{\min}}(u_L - u_R)\right)^2, & \text{otherwise.} \end{cases} \quad (5.3.17)$$

Here,  $x_0 = (2\sqrt{2} - 1)^2$  and  $f(h) := f_L(h) + f_R(h_r) - u_R - u_L$  is the depth function where for  $Z = \{L, R\}$

$$f_Z(h) = \begin{cases} 2(\sqrt{gh} - \sqrt{gh_Z}) & \text{if } h \leq h_Z, \\ (h - h_Z) \sqrt{\frac{g(h+h_Z)}{2hh_Z}} & \text{if } h > h_Z. \end{cases} \quad (5.3.18)$$

*Remark 5.3.10* (Finding  $\lambda_{ij}^n$  for the Hyperbolic Serre model). The Riemann problem associated with the system (4.4.1) is quite complicated, so in this work, we avoid solving it and define  $\lambda_{ij}^n$  using the eigenvalues established in Proposition 4.5.1. That is, for all  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$ :

$$\lambda_{ij}^n = \max(|\mathbf{V}_i^n \cdot \mathbf{n}_{ij} - (gH_i^n + \theta_i^n)^{\frac{1}{2}}|, |\mathbf{V}_j^n \cdot \mathbf{n}_{ij} + (gH_j^n + \theta_j^n)^{\frac{1}{2}}|) \quad (5.3.19)$$

with  $\theta_i^n := \partial_h \tilde{p}(H_i^n, \mathbf{N}_i^n) \left(\frac{\mathcal{E}_i}{\max(\mathcal{E}_i, H_i^n)}\right)^2$ .

Notice that when  $H_i^n \leq \mathcal{E}_i$ , the expressions  $\mathbf{V}_i^n \cdot \mathbf{n}_{ij} - (gH_i^n + \theta_i^n)^{\frac{1}{2}}$  and  $\mathbf{V}_i^n \cdot \mathbf{n}_{ij} + (gH_i^n + \theta_i^n)^{\frac{1}{2}}$  are the two extreme eigenvalues of the relaxed system (4.4.1) as established in Proposition 4.5.1. This means that the viscosity accounts for the large wave speeds induced by dispersion only when  $H_i^n \lesssim \mathcal{E}_i$ .

### 5.3.6 Defining the time-step

We define the time-step in the following way:

$$\tau_n := \text{CFL} \times \max_{i \in \mathcal{V}} \left( \frac{m_i}{\sum_{j \in \mathcal{I}^*(i)} d_{ij}^{L,n}} \right), \quad (5.3.20)$$

where CFL is a user-defined positive constant.

### 5.3.7 Generic low-order update for both models

Let us set  $t^{n+1} := t^n + \tau_n$ . Let  $\mathbf{u}_h^{n+1} := \sum_{i \in \mathcal{V}} \mathbf{U}_i^{n+1} \varphi_i$  be the update of  $\mathbf{u}$  at  $t^{n+1}$ . Let  $\mathbf{S}_i^n$  denote the total contribution of the external sources at time  $t^n$ , i.e.,  $\mathbf{S}_i^n := \mathbf{S}_F(\mathbf{U}_i^n) + \chi \mathbf{S}_G(\mathbf{U}_i^n) + \mathbf{S}_A(\mathbf{U}_i^n)$ . Here  $\chi := 1$  if the wave generation source term is active and  $\chi := 0$  otherwise. Then, the generic low-order update  $\mathbf{U}_i^{n+1}$  for all  $i \in \mathcal{V}$  is computed as follows:

$$\begin{aligned} \frac{m_i}{\tau_n} (\mathbf{U}_i^{L,n+1} - \mathbf{U}_i^n) &= m_i (\mathbf{R}_i^n + \mathbf{S}_i^n) + \sum_{j \in \mathcal{I}(i)} -\mathbf{F}_{ij}^n \\ &+ \sum_{j \in \mathcal{I}^*(i)} \left( (d_{ij}^{L,n} - \mu_{ij}^{L,n}) (\mathbf{U}_j^{*,i,n} - \mathbf{U}_i^{*,j,n}) + \mu_{ij}^{L,n} (\mathbf{U}_j^n - \mathbf{U}_i^n) \right). \end{aligned} \quad (5.3.21)$$

**Definition 5.3.1** (Saint-Venant low-order scheme). Let  $\mathbf{T}_{SV} : \mathbf{u}_h^h \mapsto \mathbf{T}_{SV}(\mathbf{u}_h^n) := \mathbf{u}_h^{n+1}$  be the low-order Saint-Venant scheme defined by (5.2.6)–(5.3.21).

**Definition 5.3.2** (Hyperbolic Serre low-order scheme). Let  $\mathbf{T}_{HS} : \mathbf{u}_h^h \mapsto \mathbf{T}_{HS}(\mathbf{u}_h^n) := \mathbf{u}_h^{n+1}$  be the low-order Hyperbolic Serre scheme defined by (5.2.7)–(5.3.21).

### 5.3.8 Higher-order time stepping

To preserve the invariant domain preserving property while attaining higher-order accuracy in time, we use the the strong stability preserving (SSP) discretization method. In particular, we use SSPRK(3,3): a third-order accurate, three-stage explicit Runge-Kutta method that is SSP. For an overview of higher-order time stepping schemes, we refer the reader to Ern and Guermond [18] and the references therein.

The SSPRK(3,3) method for solving the problem  $\partial_t \mathbf{u} = L(t, \mathbf{u})$  is carried out as follows:

$$\begin{aligned} \mathbf{w}^{(1)} &:= \mathbf{u}^n + \tau L(t^n, \mathbf{u}^n), \\ \mathbf{w}^{(2)} &:= \frac{3}{4} \mathbf{u}^n + \frac{1}{4} \left( \mathbf{w}^{(1)} + \tau L(t^n + \tau, \mathbf{w}^{(1)}) \right), \\ \mathbf{w}^{(3)} &:= \frac{1}{3} \mathbf{u}^n + \frac{2}{3} \left( \mathbf{w}^{(2)} + \tau L(t^n + \frac{1}{2} \tau, \mathbf{w}^{(2)}) \right) \end{aligned}$$

Then, the final update at  $t^{n+1}$  is given by  $\mathbf{u}^{n+1} := \mathbf{w}^{(3)}$ .

## 5.4 Well-balancing and invariant domain preserving properties

In this section, we show that the low-order algorithm (5.3.21) is well-balanced and invariant domain preserving (i.e., positivity-preserving) for each model. We begin by recalling the precise definitions of the respective properties for each model.

**Definition 5.4.1** (Exact rest for Saint-Venant). A numerical state  $(h_h, \mathbf{q}_h)^\top$  is said to be exactly at rest if  $\mathbf{q}_h = \mathbf{0}$  for all  $i \in \mathcal{V}$ , and if the approximate water height  $h_h$  and the approximate bathymetry map  $z_h$  satisfy the following alternative for all  $i \in \mathcal{V}$ : for all  $j \in \mathcal{I}(i)$ , either  $H_j = H_i = 0$  or  $H_j + Z_j = H_i + Z_i$ .

**Definition 5.4.2** (Exact rest for Hyperbolic Serre). A numerical state  $(h_h, \mathbf{q}_h, q_{1,h}, q_{2,h}, q_{3,h})^\top$  is said to be exactly at rest if  $\mathbf{q}_h = \mathbf{0}$ ,  $q_{2,h} = 0$ ,  $q_{3,h} = 0$ ,  $Q_{1,i} = H_i^2$ , for all  $i \in \mathcal{V}$ , and if the approximate water height  $h_h$  and the approximate bathymetry map  $z_h$  satisfy the following alternative for all  $i \in \mathcal{V}$ : for all  $j \in \mathcal{I}(i)$ , either  $H_j = H_i = 0$  or  $H_j + Z_j = H_i + Z_i$ .

**Definition 5.4.3** (Exactly well-balanced). A mapping  $\mathbf{T} : \mathbf{P}(\mathcal{T}_h) \rightarrow \mathbf{P}(\mathcal{T}_h)$  is said to be an exactly well-balanced scheme if  $\mathbf{T}(\mathbf{u}_h) = \mathbf{u}_h$  when  $\mathbf{u}_h$  is an exact rest state.

*Remark 5.4.1* (Reconstructed star states at rest). When the numerical state is at rest, we have that by definition:

$$\mathbf{U}_j^{*,i,n} - \mathbf{U}_i^{*,j,n} = \mathbf{0}.$$

We show this holds for each component for the Hyperbolic Serre model using Definition 5.3.7 and Definition 5.4.2.

Assume the numerical state is at rest as defined by Definition 5.4.2. To avoid the trivial case of  $H_j^n = H_i^n = 0$ , we further assume that  $H_i^n \geq \delta H_{0,\max}$  for all  $i \in \mathcal{V}$  (which implies  $H_i^\delta = H_i^n$ ). Fix  $i \in \mathcal{V}$ . Without loss of generality, assume that  $Z_i > Z_j$  for a fixed  $j \in \mathcal{I}(i)$ . Then,  $H_i^{*,j,n} = H_i^n$  and  $H_j^{*,i,n} = \max(0, H_j^n + Z_j - Z_i)$ . Note that  $H_i^n + Z_i = H_j^n + Z_j \implies H_i^n = H_j^n + Z_j - Z_i$  and since  $H_i^n \geq \delta H_{0,\max}$ , we have that  $H_j^n + Z_j - Z_i > 0$ . Thus  $H_j^{*,i,n} = H_j^n + Z_j - Z_i$ . Then, for the water

depth component  $h$  we have that

$$\begin{aligned}
\frac{H_j^{*,i,n}}{H_j^\delta} H_j^n - \frac{H_i^{*,j,n}}{H_i^\delta} H_i^n &= \frac{H_j^{*,i,n}}{H_j^n} H_j^n - \frac{H_i^{*,j,n}}{H_i^n} H_i^n \\
&= H_j^{*,i,n} - H_i^{*,j,n} \\
&= H_j^n + Z_j - Z_i - H_i^n \\
&= H_j^n + Z_j - (H_i^n + Z_i) \\
(\text{Def. 5.3.7}) &= 0.
\end{aligned}$$

Since the flow is at rest,  $\mathbf{Q}_i^n = \mathbf{0}$  for all  $i \in \mathcal{V}$  and consequently  $\frac{H_j^{*,i,n}}{H_j^\delta} \mathbf{Q}_j^n - \frac{H_i^{*,j,n}}{H_i^\delta} \mathbf{Q}_i^n = \mathbf{0}$ . Similarly,  $\mathbf{Q}_{2,i}^n = \mathbf{0}$  and  $\mathbf{Q}_{3,i}^n = \mathbf{0}$  for all  $i \in \mathcal{V}$  so  $\frac{H_j^{*,i,n}}{H_j^\delta} \mathbf{Q}_{2,j}^n - \frac{H_i^{*,j,n}}{H_i^\delta} \mathbf{Q}_{2,i}^n = \mathbf{0}$  and  $\frac{H_j^{*,i,n}}{H_j^\delta} \mathbf{Q}_{3,j}^n - \frac{H_i^{*,j,n}}{H_i^\delta} \mathbf{Q}_{3,i}^n = \mathbf{0}$ . We have left to show that the  $q_1$  component of the star state difference vanishes.

Since the flow is at rest, we have that  $Q_{1,i}^n = (H_i^n)^2$  for all  $i \in \mathcal{V}$ . Then, we see that:

$$\begin{aligned}
\left(\frac{H_j^{*,i,n}}{H_j^\delta}\right)^2 Q_{1,j}^n - \left(\frac{H_i^{*,j,n}}{H_i^\delta}\right)^2 Q_{1,i}^n &= \left(\frac{H_j^{*,i,n}}{H_j^n}\right)^2 Q_{1,j}^n - \left(\frac{H_i^{*,j,n}}{H_i^n}\right)^2 Q_{1,i}^n \\
&= \frac{(H_j^n + Z_j - Z_i)^2}{(H_j^n)^2} Q_{1,j}^n - \frac{(H_i^n)^2}{(H_i^n)^2} Q_{1,i}^n \\
(\text{Def. 5.3.7}) &= \frac{(H_i^n + Z_i - Z_i)^2}{(H_j^n)^2} (H_j^n)^2 - (H_i^n)^2 \\
&= \frac{(H_i^n)^2}{(H_j^n)^2} (H_j^n)^2 - (H_i^n)^2 \\
&= (H_i^n)^2 - (H_i^n)^2 \\
&= 0.
\end{aligned}$$

**Definition 5.4.4** (Positivity-preserving). Let us denote  $h_h(\mathbf{u}_h) = \sum_{i \in \mathcal{V}} H_i(\mathbf{u}_h) \varphi_i$  the water height of  $\mathbf{u}_h$  for any  $\mathbf{u}_h \in \mathbf{P}(\mathcal{T}_h)$ . A mapping  $\mathbf{T} : \mathbf{P}(\mathcal{T}_h) \rightarrow \mathbf{P}(\mathcal{T}_h)$  is said to be a positivity-preserving scheme if  $H_i(\mathbf{u}_h) \geq 0$ , for all  $i \in \mathcal{V}$ , implies that  $H_i(\mathbf{T}(\mathbf{u}_h)) \geq 0$  for all  $i \in \mathcal{V}$ .

**Proposition 5.4.2** (Saint-Venant IDP and WB). Let  $\mathbf{T}_{SV} : \mathbf{u}_h^h \mapsto \mathbf{T}_{SV} := \mathbf{u}_h^{n+1}$  be the Saint-Venant scheme. (i) If  $\mathbf{S}_i^n = \mathbf{S}_F(\mathbf{U}_i^n)$  is just the contribution of the Gauckler-Manning friction source, the



scheme is exactly well-balanced; (ii) The scheme is positivity-preserving if the time step satisfies the following restriction  $\tau_n(\chi \frac{\sqrt{gH_0}}{\varepsilon_i} + \frac{1}{m_i} \sum_{j \in \mathcal{I}^*(i)} d_{ij}^n) \leq 1$ .

*Proof.* (i) Since  $\mathbf{u}_h^n$  is exactly at rest, we have that  $\mathbf{V}_i^n = \mathbf{V}_j^n = 0$  for all  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$ , thus  $\mu_{ij}^{L,n} = 0$ . We also have that  $\mathbf{U}_i^{*,j,n} = \mathbf{U}_j^{*,i,n}$  as a consequence of the definition (5.3.6). Since  $\mathbf{q}_h = 0$ , we have that  $\mathbf{S}_i^n = \mathbf{0}$ . Thus, the viscosity and source terms vanish and we have the leftover terms:

$$\begin{aligned} \frac{m_i}{\tau_n} (H_i^{L,n+1} - H_i^n) &= 0, \\ \frac{m_i}{\tau_n} (\mathbf{Q}_i^{L,n+1} - 0) &= \sum_{j \in \mathcal{I}^*(i)} (gH_i^n (H_j^n + Z_j) \mathbf{c}_{ij}). \end{aligned}$$

Assume  $H_i = 0$  for all  $i \in \mathcal{V}$ . Then  $H_i^{L,n+1} = 0$  and  $\mathbf{Q}_i^{L,n+1} = \mathbf{0}$ . If  $H_i \neq 0$ , then by assumption  $H_j^{n+1} + Z_j = H_i^n + Z_i$  for all  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$ . This gives:

$$\begin{aligned} \frac{m_i}{\tau_n} (\mathbf{Q}_i^{L,n+1} - 0) &= \sum_{j \in \mathcal{I}^*(i)} (gH_i^n (H_j^n + Z_j) \mathbf{c}_{ij}) \\ &= \sum_{j \in \mathcal{I}^*(i)} (gH_i^n (H_i^n + Z_i) \mathbf{c}_{ij}) \\ &= (gH_i^n (H_i^n + Z_i)) \underbrace{\sum_{j \in \mathcal{I}^*(i)} \mathbf{c}_{ij}}_{=0} \\ &= 0. \end{aligned}$$

Thus, since  $u_h$  is at exact rest and  $H_i^{L,n+1} = H_i^n$ ,  $\mathbf{Q}_i^{L,n+1} = \mathbf{0}$ , the scheme is well-balanced with respect to rest states.

(ii) Referring to (5.3.12), we recall that  $S_h^n = -\frac{\sqrt{gH_0}}{\varepsilon_i} (H_i^n - h_{\text{wave}}(\mathbf{a}_i, t^n)) G(\frac{\mathbf{a}_i - \mathbf{x}_{\min}}{L_{\text{gen}}})$  is the source in the mass balance equation, where  $h_{\text{wave}}(\mathbf{x}, t) \geq 0$  for all  $t$  and all  $\mathbf{x} \in D$ , and  $G \in [0, 1]$ . Fixing  $i \in \mathcal{V}$  and assuming  $H_j^n \geq 0$  for all  $j \in \mathcal{I}(i)$ , the water height update in (5.3.21) can be arranged

as follows:

$$\begin{aligned} H_i^{L,n+1} \geq H_i^n - \chi \frac{\tau_n \sqrt{gH_0}}{\mathcal{E}_i} H_i^n - \frac{1}{m_i} \sum_{j \in \mathcal{I}^*(i)} \left( \mu_{ij}^{L,n} H_i^n + (d_{ij}^{L,n} - \mu_{ij}^{L,n}) H_i^{*,j,n} \right) \\ + \frac{1}{m_i} \sum_{j \in \mathcal{I}^*(i)} \left( (\mu_{ij}^{L,n} - \mathbf{v}_j^n \cdot \mathbf{c}_{ij}) H_j^n + (d_{ij}^{L,n} - \mu_{ij}^{L,n}) H_j^{*,i,n} \right) \end{aligned}$$

Since by definition  $d_{ij}^{L,n} - \mu_{ij}^{L,n} \geq 0$ ,  $\mu_{ij}^{L,n} \geq 0$ ,  $H_i^n \geq H_i^{*,j,n} \geq 0$ ,  $H_j^n \geq H_j^{*,i,n} \geq 0$ , we have the following inequality

$$H_i^{L,n+1} \geq H_i^n \left( 1 - \chi \frac{\tau_n \sqrt{gH_0}}{\mathcal{E}_i} - \frac{\tau_n}{m_i} \sum_{j \in \mathcal{I}^*(i)} d_{ij}^{L,n} \right) + \sum_{j \in \mathcal{I}^*(i)} \left( (\mu_{ij}^{L,n} - \mathbf{v}_j^n \cdot \mathbf{c}_{ij}) H_j^n \right).$$

The conclusion follows from the condition on  $\tau_n$  and the definition (5.3.14a).  $\square$

**Proposition 5.4.3** (Hyperbolic Serre IDP and WB). *Let  $\mathbf{T}_{HS} : \mathbf{u}_h^h \mapsto \mathbf{T}_{HS}(\mathbf{u}_h^n) := \mathbf{u}_h^{n+1}$  be the low-order Hyperbolic Serre scheme. (i) If  $\mathbf{S}_i^n = \mathbf{S}_F(\mathbf{U}_i^n)$  is just the contribution of the Gauckler-Manning friction source, the scheme is exactly well-balanced; (ii) The scheme is positivity-preserving if the time step satisfies the following restriction  $\tau_n \left( \chi \frac{\sqrt{gH_0}}{\mathcal{E}_i} + \frac{1}{m_i} \sum_{j \in \mathcal{I}^*(i)} d_{ij}^n \right) \leq 1$ .*

*Proof.* (i) Since  $\mathbf{u}_h^n$  is exactly at rest,  $\mu_{ij}^{L,n} = 0$  for all  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$ . We also have that  $\mathbf{U}_i^{*,j,n} = \mathbf{U}_j^{*,i,n}$  as a consequence of the definition (5.3.7). Note that it is important that  $\mathbf{Q}_{1,i}^{*,j} := \left( \frac{H_i^{*,j,n}}{H_i^n} \right)^2 \mathbf{Q}_{1,i}$  for this to hold. Then since  $\mathbf{q}_h = 0$  and  $q_{2,h} = 0$ , we have that  $\mathbf{S}_i^n = \mathbf{0}$ ,  $\mathbf{R}_1(\mathbf{U}_i, (\nabla Z)_i) = \mathbf{R}_3(\mathbf{U}_i, (\nabla Z)_i) = 0$  for all  $i \in \mathcal{V}$ . Note that  $\mathbf{R}_2(\mathbf{U}_i) = 0$  since  $\mathbf{Q}_{1,i}^n = (H_i^n)^2$ . Thus, we have that  $\mathbf{R}_i^n = \mathbf{0}$  for all  $i \in \mathcal{V}$ . The rest of the proof is exactly the same as Proposition 5.4.2.

(ii) The proof is exactly the same as Proposition 5.4.2.  $\square$

## 5.5 Local auxiliary states and bounds

We now define auxiliary states and extract exact local bounds that will be useful when limiting the yet to be defined high-order solution which might not be positivity-preserving. The key idea behind defining the exact local bounds is noticing that the low-order update (5.3.21) (for both models) can be rewritten as a convex combination of auxiliary states.

Let us recall the hydrostatic contribution  $P_{ij}^{SW}$  to the numerical flux  $\mathbf{F}_{ij}^n$ :

$$\sum_{j \in \mathcal{I}(i)} P_{ij}^{SW} = \sum_{j \in \mathcal{I}(i)} g H_i^n (H_j^n + Z_j) \mathbf{c}_{ij}.$$

We want to show that this can be re-written using the numerical conservative hydrostatic flux  $(\nabla \frac{1}{2} g h^2)$  and a modified source term. We use the algebraic property  $ab = -\frac{1}{2}(a-b)^2 + \frac{1}{2}a^2 + \frac{1}{2}b^2$  and the partition of unity properties for  $\mathbf{c}_{ij}$  to rewrite the above definition:

$$\begin{aligned} \sum_{j \in \mathcal{I}(i)} P_{ij}^{SW} &= \sum_{j \in \mathcal{I}(i)} g H_i^n (H_j^n + Z_j) \mathbf{c}_{ij} \\ &= \sum_{j \in \mathcal{I}(i)} \left[ -\frac{1}{2}g(H_i^n - H_j^n)^2 + \frac{1}{2}g(H_i^n)^2 + \frac{1}{2}g(H_j^n)^2 + g H_i^n Z_j \right] \mathbf{c}_{ij} \\ &= \sum_{j \in \mathcal{I}(i)} \left( \frac{1}{2}g(H_j^n)^2 - \frac{1}{2}g(H_i^n)^2 \right) \mathbf{c}_{ij} - \left[ \frac{1}{2}g(H_j^n - H_i^n)^2 + g H_i^n Z_j \right] \mathbf{c}_{ij} \end{aligned}$$

Let  $\mathbf{e}_q$  denote the characteristic vector for the momentum equations. That is,  $\mathbf{e}_q = (0, \mathbb{1}_d)^\top$  for the Saint-Venant model and  $\mathbf{e}_q = (0, \mathbb{1}_d, 0, 0, 0)^\top$  for the Hyperbolic Serre model where  $d$  is the spatial dimension. Then, we have the following lemma.

**Lemma 5.5.1** (Generic convex combination). *Let  $\mathbf{W}_i^{L,n+1} := \mathbf{U}_i^{L,n+1} - \tau_n \widetilde{\mathbf{R}}_i^n$ , with the modified source given by*

$$\widetilde{\mathbf{R}}_i^n := \mathbf{R}_i^n + \mathbf{S}_i^n + \sum_{j \in \mathcal{I}(i)} \left( \frac{1}{2}g(H_j^n - H_i^n)^2 - g H_i^n Z_j \right) \mathbf{c}_{ij} \mathbf{e}_q. \quad (5.5.1)$$

Assume  $1 - \frac{2\tau_n}{m_i} \sum_{j \in \mathcal{I}^*(i)} d_{ij}^{L,n} \geq 0$ . (i) *Then the following convex combination holds true:*

$$\mathbf{W}_i^{L,n+1} = \mathbf{U}_i^n \left( 1 - \frac{\tau_n}{m_i} \sum_{j \in \mathcal{I}^*(i)} 2d_{ij}^{L,n} \right) + \frac{\tau_n}{m_i} \sum_{j \in \mathcal{I}^*(i)} 2d_{ij}^{L,n} \left( \overline{\mathbf{U}}_{ij}^n + \widetilde{\mathbf{U}}_{ij}^n \right). \quad (5.5.2)$$

with the auxiliary states defined by

$$\overline{\mathbf{U}}_{ij}^n = -\frac{c_{ij}}{2d_{ij}^{L,n}} \cdot (\mathbb{f}(\mathbf{U}_j^n) - \mathbb{f}(\mathbf{U}_i^n)) + \frac{1}{2}(\mathbf{U}_j^n + \mathbf{U}_i^n), \quad (5.5.3a)$$

$$\widetilde{\mathbf{U}}_{ij}^n = \frac{d_{ij}^{L,n} - \mu_{ij}^{L,n}}{2d_{ij}^{L,n}} (\mathbf{U}_j^{*,i,n} - \mathbf{U}_j^n - (\mathbf{U}_i^{*,j,n} - \mathbf{U}_i^n)). \quad (5.5.3b)$$

(ii) Furthermore,  $\overline{\mathbf{H}}_{ij}^n + \widetilde{\mathbf{H}}_{ij}^n \geq 0$  for all  $j \in \mathcal{I}(i)$ .

*Remark 5.5.2* (Generality). Note that the convex combination is again written independently of the mathematical model of interest. Substituting the respective quantities for each model gives the appropriate convex scheme.

Notice that the quantity  $\mathbf{W}_i^{L,n+1}$  is an update corresponding to solving the hyperbolic system without sources. The ‘‘source removing’’ concept is used in §5.7 to perform the convex limiting technique. We now define the bounds that we use to limit the provisional higher-order solution. The strategy that we propose consists of enforcing bounds that are naturally satisfied by the low-order update (5.5.2). Let us define the kinetic energy functional  $\psi(\mathbf{u}) := \frac{1}{2} \frac{1}{h(\mathbf{u})} \|\mathbf{q}(\mathbf{u})\|_{\ell^2}^2$ . Then we have the following local in space and time bounds for each model:

*Remark 5.5.3* (Saint-Venant bounds). The local bounds for the Saint-Venant model are given by:

$$\mathbf{h}_i^{n,\min} := \min_{j \in \mathcal{I}(i)} (\overline{\mathbf{H}}_{ij}^n + \widetilde{\mathbf{H}}_{ij}^n), \quad \mathbf{h}_i^{n,\max} := \max_{j \in \mathcal{I}(i)} (\overline{\mathbf{H}}_{ij}^n + \widetilde{\mathbf{H}}_{ij}^n) \quad (5.5.4a)$$

$$\mathbf{K}_i^{n,\max} := \max_{j \in \mathcal{I}(i)} \psi(\overline{\mathbf{U}}_{ij}^n + \widetilde{\mathbf{U}}_{ij}^n). \quad (5.5.4b)$$

*Remark 5.5.4* (Hyperbolic Serre bounds). The local bounds for the Hyperbolic Serre model are given by:

$$\mathbf{h}_i^{n,\min} := \min_{j \in \mathcal{I}(i)} (\overline{\mathbf{H}}_{ij}^n + \widetilde{\mathbf{H}}_{ij}^n), \quad \mathbf{h}_i^{n,\max} := \max_{j \in \mathcal{I}(i)} (\overline{\mathbf{H}}_{ij}^n + \widetilde{\mathbf{H}}_{ij}^n), \quad (5.5.5a)$$

$$q_{1,i}^{n,\min} := \min_{j \in \mathcal{I}(i)} (\overline{\mathbf{Q}}_{1,ij}^n + \widetilde{\mathbf{Q}}_{1,ij}^n), \quad q_{1,i}^{n,\max} := \max_{j \in \mathcal{I}(i)} (\overline{\mathbf{Q}}_{1,ij}^n + \widetilde{\mathbf{Q}}_{1,ij}^n), \quad (5.5.5b)$$

$$\mathbf{K}_i^{n,\max} := \max_{j \in \mathcal{I}(i)} \psi(\overline{\mathbf{U}}_{ij}^n + \widetilde{\mathbf{U}}_{ij}^n). \quad (5.5.5c)$$

Let us expand on the relationship between the local bounds defined above and the convex combination update (5.5.2) with an example. We denote the components of the low-order solutions without sources,  $\mathbf{W}^L$ , as follows:

$$\mathbf{W}^L := \mathbf{U}(\mathbf{W}^L) = (\mathbf{H}(\mathbf{W}^L), \mathbf{Q}(\mathbf{W}^L) \dots)^\top.$$

Using the properties of convexity, we can extract the following inequality on the water height update  $\mathbf{H}(\mathbf{W}_i^{L,n+1})$  from the convex scheme (5.5.2): or equivalently:

$$\min_{j \in \mathcal{I}(i)} (\overline{\mathbf{H}}_{ij}^n + \widetilde{\mathbf{H}}_{ij}^n) \leq \mathbf{H}_i^{L,n+1} - \tau_n S_h^n \leq \max_{j \in \mathcal{I}(i)} (\overline{\mathbf{H}}_{ij}^n + \widetilde{\mathbf{H}}_{ij}^n),$$

$$h_i^{n,\min} \leq \mathbf{H}(\mathbf{W}_i^{L,n+1}) \leq h_i^{n,\max}.$$

Thus, these local bounds arise naturally from the low-order numerical scheme. More precise statements are made in §5.7.1.

## 5.6 Provisional high-order method

In this section, we introduce a provisional higher-order method (second-order accurate in space) that may violate the invariant domain preserving property (i.e.,  $h > 0$ ). The two key ideas are as follows:

- (i) we reduce numerical dispersive errors induced by the lumped mass matrix,
- (ii) we define higher-order graph viscosity coefficients  $d_{ij}^{H,n}, \mu_{ij}^{H,n}$  via the estimation of an entropy residual/commutator using the energy results of both models.

### 5.6.1 Wave generation

If active, the wave generation mechanism must be modified in the high-order method since this source can potentially create dry states if the amplitude of the wave generated is larger than the mean water depth. We formalize this by setting  $h_{\text{wave}}^{\min} := \min_{\mathbf{x} \in D, t > 0} h_{\text{wave}}(\mathbf{x}, t)$  and by introducing the cut-off function  $\chi \in C(\mathbb{R}; [0, 1])$  defined by  $\chi(\xi) = 1$  if  $\xi \leq \frac{1}{2}$ ,  $\chi(\xi) := 4(\xi - 1)^2(4\xi - 1)$  if

$\frac{1}{2} \leq \xi \leq 1$ , and  $\chi(\xi) = 0$  otherwise. We then redefine the source term  $\mathbf{S}_i^n$  used in the low-order approximation (5.3.21) by setting

$$\mathbf{S}_i^n := \mathbf{S}_F(\mathbf{U}_i^n) + \chi\left(\frac{h_i^{i,\max} - h_i^{i,\min}}{h_{\text{wave}}^{\min}}\right) \mathbf{S}_G(\mathbf{U}_i^n) + \mathbf{S}_A(\mathbf{U}_i^n). \quad (5.6.1)$$

We also use this definition for the high-order update (see (5.6.8)). For most realistic applications, the amplitude of the incoming waves is of reasonable size and  $h_i^{n,\max} - h_i^{n,\min}$  is a priori small compared to  $h_{\text{wave}}^{\min}$  and the cut-off is therefore inactive. In particular, it is never active in the simulations reported in Chapter 6. This cut-off is necessary for theoretical purposes (see Theorem 5.7.8).

### 5.6.2 Commutator-based entropy viscosity

We present the definition of the higher-order artificial viscosity coefficients  $d_{ij}^{H,n}, \mu_{ij}^{H,n}$  following the method introduced in Guermond et al. [28]. The key idea consists of measuring the smoothness of an entropy by measuring how well a chain rule is satisfied by the discretization described above.

Let  $(E(\mathbf{u}), \mathbf{F}(\mathbf{u}))$  be any entropy pair that satisfies the following relation:

$$\nabla \cdot (\mathbf{F}(\mathbf{u})) = (\nabla E(\mathbf{u}))^\top \nabla \cdot (\mathbb{f}(\mathbf{u})), \quad (5.6.2)$$

coupled with a respective flux function  $\mathbb{f}(\mathbf{u})$ . We want to estimate the entropy production by inserting the approximate solution in (5.6.2). For all  $i \in \mathcal{V}$  we define the entropy commutator as follows:

$$C_i^n := \sum_{j \in \mathcal{I}(i)} \mathbf{c}_{ij} \cdot \left( \mathbf{F}(\mathbf{U}_j^n) - (\nabla E(\mathbf{U}_i^n))^\top \mathbb{f}(\mathbf{U}_j^n) \right) \quad (5.6.3)$$

This quantity measures how well the finite element approximation satisfies the chain rule (5.6.2). Notice that when the approximate solution  $\mathbf{u}_h^n$  is smooth, the quantity  $C_i^n$  is as small as the truncation error provided by the finite element setting. For piecewise linear elements on unstructured meshes  $C_i^n$  scales like  $\mathcal{O}(h)$  where  $h$  is the mesh-size. We define the normalized entropy resid-

ual/commutator to be

$$R_i^n := \frac{|C_i^n|}{D_i^n}, \quad (5.6.4a)$$

$$D_i^n := \left| \sum_{i \in \mathcal{I}(i)} \mathbf{c}_{ij} \cdot \mathbf{F}(\mathbf{U}_j^n) \right| + \left| \sum_{i \in \mathcal{I}(i)} \mathbf{c}_{ij} \cdot ((\nabla E(\mathbf{U}_i^n))^\top \mathbb{f}(\mathbf{U}_j^n)) \right|, \quad (5.6.4b)$$

where  $D_i^n$  is the rescaling factor. We then define the higher-order graph viscosity (or entropy viscosity) as follows:

$$d_{ij}^{\text{H},n} = d_{ij}^{\text{L},n} \max(R_i^n, R_j^n), \quad d_{ii}^{\text{H},n} := - \sum_{j \in \mathcal{I}^*(i)} d_{ij}^{\text{H},n} \quad (5.6.5)$$

$$\mu_{ij}^{\text{H},n} = \mu_{ij}^{\text{L},n} \max(R_i^n, R_j^n), \quad \mu_{ii}^{\text{H},n} := - \sum_{j \in \mathcal{I}^*(i)} \mu_{ij}^{\text{H},n}. \quad (5.6.6)$$

Notice that  $R_i^n \in [0, 1]$ . Denoting by  $\text{diam}(D)$  the diameter of  $D$ , it is argued in [31] that  $R_i^n = \mathcal{O}(h/\text{diam}(D))$  when the solution is smooth. Thus, by making the high-order graph viscosities proportional to the entropy production, (5.6.5)-(5.6.6), we have  $d_{ij}^{\text{H},n} \sim d_{ij}^{\text{L},n}$  when the entropy production is large, for instance in shock regions, and  $d_{ij}^{\text{H},n} \sim \mathcal{O}\left(\frac{h}{\text{diam}(D)}\right) d_{ij}^{\text{L},n}$  in regions where the approximate solution is smooth.

*Remark 5.6.1* (The Hyperbolic Serre entropy pair). Recall that it was shown in Proposition 4.5.8 that the Hyperbolic Serre energy functional and energy flux defined in (4.5.8a) do not satisfy the entropy chain rule (5.6.2), but instead satisfy an equation of the form:

$$\nabla \cdot (\mathbf{F}(\mathbf{u})) = (\nabla E(\mathbf{u}))^\top (\nabla \cdot \mathbb{f}(\mathbf{u}) - \mathbf{S}).$$

*Remark 5.6.2* (The Saint-Venant entropy pair). Let  $\mathbf{u}_{\text{SV}} := (h, \mathbf{q})^\top$  be the Saint-Venant conserved

variable. Recall the following entropy pair for the Saint-Venant model defined in (2.2.3):

$$\begin{aligned} E_{\text{SV}}(\mathbf{u}) &:= \frac{1}{2}gh^2 + \frac{1}{2}hv^2, \\ \mathbf{F}_{\text{SV}}(\mathbf{u}) &:= \mathbf{v}(E_{\text{SV}}(\mathbf{u}) + \frac{1}{2}gh^2). \end{aligned}$$

This pair satisfies the chain rule  $\nabla \cdot (\mathbf{F}_{\text{SV}}(\mathbf{u}_{\text{SV}})) = (\nabla E_{\text{SV}}(\mathbf{u}_{\text{SV}}))^T \nabla \cdot (\mathbb{f}_{\text{SV}}(\mathbf{u}_{\text{SV}}))$  where  $\mathbb{f}_{\text{SV}}(\mathbf{u}_{\text{SV}})$  is the Saint-Venant hyperbolic flux.

*Remark 5.6.3 (Comparison).* In Chapter 6, we use the Saint-Venant entropy pair for the Saint-Venant high-order method and in some of the numerical illustrations for the Hyperbolic Serre model as well. We also show that the convergence behavior of the numerical method with either entropy pair, (4.5.10) or (2.2.3), is similar.

### 5.6.3 Consistent mass matrix

Numerical dispersion errors induced by the lumped mass matrix can be significantly reduced by using the consistent mass matrix for the discretization of the time derivative term  $\mathbf{u}_t$  (at least for piecewise linear approximation). For more details on this issue, we refer the reader to Guermond and Pasquetti [29] and the references therein.

We now replace the lumped mass matrix in (5.3.21) with the consistent mass matrix defined in (5.2.3) and the low-order graph viscosity coefficients in (5.3.21) with the entropy-viscosity coefficients (5.6.5)-(5.6.6). Then, the provisional higher-order update (for either the Saint-Venant or Hyperbolic Serre model) is given as follows:

$$\begin{aligned} \sum_{j \in \mathcal{I}(i)} m_{ij} \frac{\tilde{\mathbf{U}}_j^{\text{H},n+1} - \mathbf{U}_j^n}{\tau_n} &= m_i (\mathbf{R}_i^n + \mathbf{S}_i^n) + \sum_{j \in \mathcal{I}(i)} -\mathbf{F}_{ij}^n \\ &+ \sum_{j \in \mathcal{I}^*(i)} \left( (d_{ij}^{\text{H},n} - \mu_{ij}^{\text{H},n}) (\mathbf{U}_j^{*,i,n} - \mathbf{U}_i^{*,j,n}) + \mu_{ij}^{\text{H},n} (\mathbf{U}_j^n - \mathbf{U}_i^n) \right). \end{aligned} \quad (5.6.8)$$

We see that finding  $\tilde{\mathbf{U}}^{\text{H},n+1}$  in (5.6.8) requires the inversion of the consistent mass matrix at every time step. Since this may be computationally costly, we follow the ideas of [29] and Maier



and Kronbichler [49, Sec. (3.4)], and approximate the inverse of the mass matrix with a Neumann series. We do this as follows. We denote by  $\tilde{\mathbf{S}}_i^n$  the right-hand side in (5.6.8) and rewrite (5.6.8) as follows:

$$\sum_{j \in \mathcal{I}(i)} \frac{m_{ij}}{m_j} \frac{m_j}{\tau_n} (\tilde{\mathbf{U}}_j^{\text{H},n+1} - \mathbf{U}_j^n) = \tilde{\mathbf{S}}_i^n. \quad (5.6.9)$$

We then approximate the inverse  $(\frac{m_{ij}}{m_j})^{-1}$  with the first-order approximation of its Neumann series representation:

$$\left(\frac{m_{ij}}{m_j}\right)^{-1} = \left(\delta_{ij} - \left(\delta_{ij} - \frac{m_{ij}}{m_j}\right)\right)^{-1} \approx \delta_{ij} + \left(\delta_{ij} - \frac{m_{ij}}{m_j}\right) = \delta_{ij} + b_{ij}.$$

Then, using that  $\sum_{j \in \mathcal{I}(i)} b_{ji} = 0$  (by the partition of unity), we infer the following new expression for the provisional higher-order update  $\mathbf{U}^{\text{H},n+1}$ :  $\frac{m_i}{\tau_n} (\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^n) = \tilde{\mathbf{S}}_i^n + \sum_{j \in \mathcal{I}(i)} (b_{ij} \tilde{\mathbf{S}}_j^n - b_{ji} \tilde{\mathbf{S}}_i^n)$ . Replacing the definition of  $\tilde{\mathbf{S}}_i^n$  therein gives

$$\begin{aligned} \frac{m_i}{\tau_n} (\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^n) &= m_i (\mathbf{R}_i^n + \mathbf{S}_i^n) + \sum_{j \in \mathcal{I}(i)} -\mathbf{F}_{ij}^n \\ &+ \sum_{j \in \mathcal{I}^*(i)} \left( b_{ij} \tilde{\mathbf{S}}_j^n - b_{ji} \tilde{\mathbf{S}}_i^n + (d_{ij}^{\text{H},n} - \mu_{ij}^{\text{H},n}) (\mathbf{U}_j^{*,i,n} - \mathbf{U}_i^{*,j,n}) + \mu_{ij}^{\text{H},n} (\mathbf{U}_j^n - \mathbf{U}_i^n) \right). \end{aligned} \quad (5.6.10)$$

*Remark 5.6.4 (Generality).* The details described above for the provisional high-order method were independent of the mathematical model.

#### 5.6.4 Loss of positivity

It is proved in Guermond et al. [30, Theorem 3.2] that the presence of the consistent mass matrix in any scheme that uses continuous finite elements based on artificial viscosity (such as (5.6.8)) and explicit time stepping violates the maximum principle for scalar conservation laws. A consequence of this result is that the scheme (5.6.8) is non-positivity-preserving regardless of the definition of the artificial viscosity coefficients. It is also observed numerically in [28] that the use of the higher-order entropy viscosity coefficients (5.6.5)-(5.6.6) in (5.3.21) can cause the scheme to be

non-positivity-preserving as well. We correct the loss of positivity in the following section.

## 5.7 Convex limiting with sources

In this section, we describe the convex limiting technique for both models that is used to make the higher-order methods described above invariant domain preserving (i.e., positivity-preserving). For the work shown here, we build off the ideas presented in Guermond et al. [31, 32], but give an emphasis on how to apply the convex limiting methodology to a hyperbolic system with source terms since it not well documented in the literature.

### 5.7.1 Quasiconcave functionals and bounds

We now give some definitions and results that will illustrate the notion of *convex* limiting. Let  $m$  be again the problem system size (e.g.,  $m := d + 1$  for Saint-Venant or  $m := d + 4$  for the hyperbolic Serre model).

**Definition 5.7.1** (Quasiconcavity). Given a convex set  $\mathcal{B} \subset \mathbb{R}^m$ , we say that a function  $\Psi : \mathcal{B} \rightarrow \mathbb{R}$  is quasiconcave if every upper level set of  $\Psi$  is *convex*; that is, the set  $L_\lambda(\Psi) := \{\mathbf{u} \in \mathcal{B} \mid \Psi(\mathbf{u}) \geq \lambda\}$  is convex for any  $\lambda \in \mathbb{R}$ .

**Definition 5.7.2.** A function  $\Psi : \mathcal{A} \rightarrow \mathbb{R}$  is quasiconvex if  $-\Psi$  is quasiconcave.

**Lemma 5.7.1.** Let  $\mathcal{B} := \{\mathbf{u} \in \mathbb{R}^m \mid h > 0\} \subset \mathbb{R}^{d+4}$ . Let  $\Psi : \mathcal{B} \rightarrow \mathbb{R}$  and assume that the product  $h\Psi$  is concave. Then the function  $\Psi$  is quasiconcave.

*Proof.* This is a special case of the result in Guermond et al. [32, Lem. 7.4]. □

The key idea behind the convex limiting technique is to correct (i.e., limit) the provisional high-order update (5.6.10) so that it satisfies the same (yet to be defined) quasiconcave constraints as the low-order update (5.7.8). Let  $\mathcal{L}$  be an index set denoting the total number of quasiconcave functionals numbered from 1 to  $\hat{\ell}$ .

*Remark 5.7.2* (Quasiconcave functionals for Saint-Venant). Recall that the conserved variable for the system (2.2.1) is  $\mathbf{u} := (h, \mathbf{q})^\top$  where  $h$  is the water depth and  $\mathbf{q}$  the momentum. The functional

$\Psi_1 : \mathbb{R}^{d+1} \ni (\mathbf{h}, \mathbf{q})^\top \mapsto \mathbf{h} \in \mathbb{R}$  is linear, hence concave, hence quasiconcave; this functional is also well-defined over  $\mathbb{R}^{d+1}$ . Let us set

$$\mathcal{A} := \{\mathbf{u} \in \mathbb{R}^{d+1} \mid \mathbf{h} > 0\} \quad (5.7.1)$$

Observe that  $\mathcal{A}$  is convex. Then, another important example is the (negative) kinetic energy  $\Psi_2 : \mathcal{A} \ni (\mathbf{h}, \mathbf{q})^\top \mapsto -\frac{1}{2\mathbf{h}}\|\mathbf{q}\|_{\ell^2}^2$ . Since the function  $\mathbf{h}\Psi_2 := -\frac{1}{2}\|\mathbf{q}\|_{\ell^2}^2$  is concave, using Lemma 5.7.1 we conclude the (negative) kinetic energy is quasiconcave. We are going to use the functionals  $\Psi_1$  and  $\Psi_2$  and bounds defined in (5.5.4) to enforce positivity of the water height, positivity and a local maximum principle on the kinetic energy. Letting  $\mathcal{L} := \{1:3\}$  and  $\mathcal{A} := \{\mathbf{u} \in \mathbb{R}^{d+1} \mid \mathbf{h} > 0\} \subset \mathbb{R}^{d+1}$ , we are going to work with the family of quasiconcave functionals  $\{\Psi_l^{i,n}\}_{l \in \mathcal{L}}, \Psi_l^{i,n} : \mathcal{A} \rightarrow \mathbb{R}$  defined as follows:

$$\Psi_1^{i,n}(\mathbf{u}) = \mathbf{h} - \mathbf{h}_i^{n,\min}, \quad \Psi_2^{i,n}(\mathbf{u}) = \mathbf{h}_i^{n,\max} - \mathbf{h}, \quad (5.7.2a)$$

$$\Psi_2^{i,n}(\mathbf{u}) = \mathbf{K}_i^{n,\max} - \frac{1}{2\mathbf{h}}\|\mathbf{q}\|_{\ell^2}^2. \quad (5.7.2b)$$

*Remark 5.7.3 (Quasiconcave functionals for Hyperbolic Serre).* Recall that the conserved variable for the system (4.4.1) is  $\mathbf{u} := (\mathbf{h}, \mathbf{q}, q_1, q_2, q_3)^\top$  where  $\mathbf{h}$  is the water height,  $\mathbf{q}$  the momentum,  $q_1, q_2, q_3$  the auxiliary variables which are thought of ansatz to  $\mathbf{h}^2$ ,  $\mathbf{h}\dot{\mathbf{h}} + \frac{3}{2}q_3$  and  $\mathbf{q} \cdot \nabla z$ , respectively. The functional  $\Psi_1 : \mathbb{R}^{d+4} \ni (\mathbf{h}, \mathbf{q}, q_1, q_2, q_3)^\top \mapsto \mathbf{h} \in \mathbb{R}$  is linear, hence concave, hence quasiconcave; this functional is also well-defined over  $\mathbb{R}^{d+4}$ . The functional  $\Psi_2 : \mathbb{R}^{d+4} \ni (\mathbf{h}, \mathbf{q}, q_1, q_2, q_3)^\top \mapsto q_1 \in \mathbb{R}$  is also linear, hence quasiconcave. Let us set

$$\mathcal{A} := \{\mathbf{u} \in \mathbb{R}^{d+4} \mid \mathbf{h} > 0\} \quad (5.7.3)$$

Observe that  $\mathcal{A}$  is convex. Similarly as above, another example is the (negative) kinetic energy  $\Psi_3 : \mathcal{A} \ni (\mathbf{h}, \mathbf{q}, q_1, q_2, q_3)^\top \mapsto -\frac{1}{2\mathbf{h}}\|\mathbf{q}\|_{\ell^2}^2$ . We use the above functionals and bounds defined in (5.5.5) to enforce positivity of the water height, positivity of the auxiliary variable  $q_1$ , and a local maximum

principle on the kinetic energy. Letting  $\mathcal{L} := \{1:5\}$  and  $\mathcal{A} := \{\mathbf{u} \in \mathbb{R}^{d+4} \mid \mathbf{h} > 0\} \subset \mathbb{R}^{d+4}$ , we are going to work with the family of quasiconcave functionals  $\{\Psi_l^{i,n}\}_{l \in \mathcal{L}}$ ,  $\Psi_l^{i,n} : \mathcal{A} \rightarrow \mathbb{R}$  defined as follows:

$$\Psi_1^{i,n}(\mathbf{u}) = \mathbf{h} - \mathbf{h}_i^{n,\min}, \quad \Psi_2^{i,n}(\mathbf{u}) = \mathbf{h}_i^{n,\max} - \mathbf{h}, \quad (5.7.4a)$$

$$\Psi_3^{i,n}(\mathbf{u}) = q_1 - q_{1,i}^{n,\min}, \quad \Psi_4^{i,n}(\mathbf{u}) = q_{1,i}^{\max} - q_1, \quad (5.7.4b)$$

$$\Psi_5^{i,n}(\mathbf{u}) = \mathbf{K}_i^{n,\max} - \frac{1}{2\mathbf{h}} \|\mathbf{q}\|_{\ell^2}^2. \quad (5.7.4c)$$

The following result is essential for the rest of the convex limiting argumentation.

**Lemma 5.7.4.** *Let  $n \geq 0$ ,  $i \in \mathcal{V}$ , and assume that  $\frac{2\tau_n}{m_i} \sum_{j \in \mathcal{I}^*(i)} d_{ij}^{L,n} \leq 1$ . Then, the low-order update  $\mathbf{W}_i^{L,n+1}$  computed by (5.5.2) is in  $\mathcal{A}$  and satisfies the following constraints for all  $l \in \mathcal{L}$ :*

$$\Psi_l^{i,n}(\mathbf{W}_i^{L,n+1}) \geq 0. \quad (5.7.5)$$

*Proof.* Under the assumption  $1 - \frac{2\tau_n}{m_i} \sum_{j \in \mathcal{I}^*(i)} d_{ij}^{L,n} \geq 0$ ,  $\mathbf{W}_i^{L,n+1}$  is a convex combination of the auxiliary states (5.5.3a)-(5.5.3b); thus, by Lemma 5.5.1, the update  $\mathbf{W}_i^{L,n+1}$  is in  $\mathcal{A}$ . The constraints  $\Psi_l^{i,n}(\mathbf{W}_i^{L,n+1}) \geq 0$  are a consequence of the convex combination (5.5.2), the definitions of the bounds in §5.7.1 and quasiconcavity.  $\square$

The limiting is done sequentially: First we limit  $\mathbf{W}_i^{H,n+1}$  with respect to  $\Psi_1^{i,n}$  and construct a  $\mathbf{W}_i^{1,n+1}$  so that  $\Psi_1^{i,n}(\mathbf{W}_i^{1,n+1}) \geq 0$ . This guarantees positivity of the water height and must be computed before limiting with respect to the other quantities  $(\Psi_l^{i,n})_{l>1}$ . This is explained in more detail below.

## 5.7.2 Limiting process

We discuss in this section the proposed convex limiting methodology. The main idea going forward is that we apply the limiting process to the solution without sources  $\mathbf{W}^n$  and then “put back” the sources after enforcing all the quasiconcave constraints.

Let  $\mathbf{W}^{\text{H},n+1} := \mathbf{U}^{\text{H},n+1} - \tau_n(\mathbf{R}_i^n + \mathbf{S}_i^n)$  be the provisional high-order update without sources. Our goal is to construct the final update  $\mathbf{U}_i^{n+1}$  so that  $\mathbf{W}_i^{n+1} := \mathbf{U}_i^{n+1} - \tau_n(\mathbf{R}_i^n + \mathbf{S}_i^n)$  satisfies all the constraints  $\Psi_l^{i,n}(\mathbf{W}_i^{n+1}) \geq 0$ ,  $l \in \mathcal{L}$ , defined in §5.7.1. For this purpose, we also define the low-order update without sources,  $\mathbf{W}^{\text{L},n+1} := \mathbf{U}^{\text{L},n+1} - \tau_n(\mathbf{R}_i^n + \mathbf{S}_i^n)$ . Proceeding as in the Flux-Corrected-Transport methodology (see Zalesak [67], [28] and references therein), we now compute the difference  $\mathbf{W}^{\text{H},n+1} - \mathbf{W}^{\text{L},n+1}$  by subtracting (5.3.21) from (5.6.10). This gives

$$m_i(\mathbf{W}_i^{\text{H},n+1} - \mathbf{W}_i^{\text{L},n+1}) = \sum_{j \in \mathcal{I}^*(i)} \mathbf{A}_{ij}^n, \quad (5.7.6)$$

with the  $\mathbb{R}^m$ -valued coefficients  $\mathbf{A}_{ij}^n$  defined by

$$\begin{aligned} \mathbf{A}_{ij}^n := \tau_n \left[ b_{ij} \tilde{\mathbf{S}}_j^n - b_{ji} \tilde{\mathbf{S}}_i^n + ((d_{ij}^{\text{H},n} - \mu_{ij}^{\text{H},n}) - (d_{ij}^{\text{L},n} - \mu_{ij}^{\text{L},n}))(\mathbf{U}_j^{*,i,n} - \mathbf{U}_i^{*,j,n}) \right. \\ \left. + (\mu_{ij}^{\text{H},n} - \mu_{ij}^{\text{L},n})(\mathbf{U}_j^n - \mathbf{U}_i^n) \right]. \quad (5.7.7) \end{aligned}$$

Notice that  $\mathbf{A}_{ij}^n = -\mathbf{A}_{ji}^n$ , which implies global mass conservation

$$\sum_{i \in \mathcal{V}} m_i \mathbf{W}_i^{\text{H},n+1} = \sum_{i \in \mathcal{V}} m_i \mathbf{W}_i^{\text{L},n+1},$$

which is equivalent to

$$\sum_{i \in \mathcal{V}} m_i \mathbf{U}_i^{\text{H},n+1} = \sum_{i \in \mathcal{V}} m_i \mathbf{U}_i^{\text{L},n+1}.$$

That is to say, the high-order solution and low-order solution have the same mass whether the source term is present or not.

Using (5.7.6), we introduce the final limited update as follows:

$$\mathbf{W}_i^{n+1} = \sum_{j \in \mathcal{I}^*(i)} \theta_j \left( \mathbf{W}_i^{\text{L},n+1} + \ell_{ij} \mathbf{P}_{ij}^n \right), \quad \text{with} \quad \mathbf{P}_{ij}^n := \frac{1}{m_i \theta_j} \mathbf{A}_{ij}^n, \quad (5.7.8)$$

where  $\{\theta_j\}_{j \in \mathcal{I}^*(i)}$  is any set of strictly positive coefficients adding up to 1. In the computations

reported below, we take  $\theta_j := \frac{1}{\text{card}(\mathcal{I}(j)) - 1}$ . The parameter  $\ell_{ij} \in [0, 1]$ , which we call the limiter, is defined to be symmetric  $\ell_{ij} = \ell_{ji}$  to preserve the mass conservation property mentioned above. Note that  $\mathbf{W}_i^{n+1} = \mathbf{W}_i^{\text{L},n+1}$  if  $\ell_{ij} = 0$  (i.e.,  $\mathbf{U}_i^{n+1} = \mathbf{U}_i^{\text{L},n+1}$ ) and  $\mathbf{W}_i^{n+1} = \mathbf{W}_i^{\text{H},n+1}$  if  $\ell_{ij} = 1$ . The key idea is to find a set of limiters  $\ell_{ij} \in [0, 1]$  as large as possible so that  $\Psi_i^{l,n}(\mathbf{W}_i^{n+1}) \geq 0$  for all  $l \in \mathcal{L}$ . Notice that this optimization program is possible since  $\ell_{ij} = 0$  is in the feasible set owing to Lemma 5.7.4. The following lemma proved in Guermond et al. [31, Lem. 4.4] is paramount for the convex limiting technique and sums up how to efficiently find the limiting parameters  $\ell_{ij}$ .

**Lemma 5.7.5.** *Let  $\mathcal{A} \subset \mathbb{R}^m$  and  $\Psi \in C^0(\mathcal{A}; \mathbb{R})$  be such that  $\{\mathbf{u} \in \mathcal{A} \mid \Psi(\mathbf{u}) \geq 0\}$  is convex. Let  $i \in \mathcal{V}$  and  $j \in \mathcal{I}(i)$ . Assume that  $\mathbf{W}_i^{\text{L},n+1} \in \mathcal{A}$ ,  $\Psi(\mathbf{W}_i^{\text{L},n+1}) \geq 0$ , and  $\Psi(\mathbf{W}_i^{\text{L},n+1} + \mathbf{P}_{ij}^n) < 0$  (otherwise there is nothing to limit), then*

(i) *There is a unique  $\ell_j^i \in [0, 1]$  such that*

$$\Psi(\mathbf{W}_i^{\text{L},n+1} + \ell_j^i \mathbf{P}_{ij}^n) = 0, \quad (5.7.9)$$

$\Psi(\mathbf{W}_i^{\text{L},n+1} + \ell \mathbf{P}_{ij}^n) \geq 0$  for all  $\ell \in [0, \ell_j^i]$ , and  $\Psi(\mathbf{W}_i^{\text{L},n+1} + \ell \mathbf{P}_{ij}^n) < 0$  for all  $\ell \in (\ell_j^i, 1]$ .

(ii) *Setting  $\ell_{ij} = \min(\ell_j^i, \ell_i^j)$ , we have  $\Psi(\mathbf{W}_i^{\text{L},n+1} + \ell_{ij} \mathbf{P}_{ij}^n) \geq 0$  and  $\ell_{ij} = \ell_{ji}$ .*

(iii) *Let  $\mathbf{W}_i^{n+1}$  be defined by (5.7.8), then  $\Psi(\mathbf{W}_i^{n+1}) \geq 0$ .*

### 5.7.3 Application to the system (4.4.1)

We now illustrate Lemma 5.7.5 with  $\Psi := \Psi_l^{i,n}$ ,  $l \in \mathcal{L}$ , defined in (5.7.4) and  $\mathcal{A} := \{\mathbf{u} \in \mathbb{R}^{d+4} \mid \mathbf{h} > 0\}$ . For the sake of brevity, we restrict ourselves to discussing the application of the limiting technique only to the Hyperbolic Serre model since application to Saint-Venant model is similar. The limiting is implemented by traversing  $\mathcal{L}$  from the smallest index to the largest one.

We begin with the limiting of the water height. To avoid divisions by zero, we introduce the small parameter  $\delta \mathbf{h}_i^{n,\max}$  where  $\delta := 10^{-14}$  for all  $i \in \mathcal{V}$ . Let us denote the  $\mathbf{h}$ -component of  $\mathbf{P}_{ij}$  by

$\mathbf{P}_{ij}^h$ . Then we set:

$$\ell_j^{i,h} = \begin{cases} \min \left( \frac{|\mathbf{h}_i^{n,\min} - \mathbf{H}(\mathbf{W}_i^{L,n+1})|}{|\mathbf{P}_{ij}^h| + \delta \mathbf{h}_i^{n,\max}}, 1 \right), & \text{if } \mathbf{H}(\mathbf{W}_i^{L,n+1}) + \mathbf{P}_{ij}^h < \mathbf{h}_i^{n,\min}, \\ 1, & \mathbf{h}_i^{n,\min} \leq \mathbf{H}(\mathbf{W}_i^{L,n+1}) + \mathbf{P}_{ij}^h \leq \mathbf{h}_i^{n,\max}, \\ \min \left( \frac{|\mathbf{h}_i^{n,\max} - \mathbf{H}(\mathbf{W}_i^{L,n+1})|}{|\mathbf{P}_{ij}^h| + \delta \mathbf{h}_i^{n,\max}}, 1 \right), & \text{if } \mathbf{h}_i^{n,\max} < \mathbf{H}(\mathbf{W}_i^{L,n+1}) + \mathbf{P}_{ij}^h. \end{cases} \quad (5.7.10)$$

This guarantees that  $\Psi_1(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij}) \geq 0$  and  $\Psi_2(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij}) \geq 0$  for all  $\ell \in [0, \ell_j^{i,h}]$ . This enforces a local minimum principle and a local maximum principle on the water height. As a corollary this also enforces positivity of the water height  $\mathbf{H}_i^{n+1}$ .

We proceed similarly to limit  $q_1$  since the functionals  $\Psi_3$  and  $\Psi_4$  are linear. Denoting the  $q_1$ -component of  $\mathbf{P}_{ij}$  by  $\mathbf{P}_{ij}^{q_1}$ , for  $\ell_j^{i,q_1} \in [0, \ell_j^{i,h}]$ , we set

$$\ell_j^{i,q_1} = \begin{cases} \min \left( \frac{|\mathbf{q}_{1,i}^{n,\min} - \mathbf{Q}_1(\mathbf{W}_i^{L,n+1})|}{|\mathbf{P}_{ij}^{q_1}| + \delta \mathbf{q}_{1,i}^{n,\max}}, 1 \right), & \text{if } \mathbf{Q}_1(\mathbf{W}_i^{L,n+1}) + \mathbf{P}_{ij}^{q_1} < \mathbf{q}_{1,i}^{n,\min}, \\ 1, & \mathbf{q}_{1,i}^{n,\min} \leq \mathbf{Q}_1(\mathbf{W}_i^{L,n+1}) + \mathbf{P}_{ij}^{q_1} \leq \mathbf{q}_{1,i}^{n,\max}, \\ \min \left( \frac{|\mathbf{q}_{1,i}^{n,\max} - \mathbf{Q}_1(\mathbf{W}_i^{L,n+1})|}{|\mathbf{P}_{ij}^{q_1}| + \delta \mathbf{q}_{1,i}^{n,\max}}, 1 \right), & \text{if } \mathbf{q}_{1,i}^{n,\max} < \mathbf{Q}_1(\mathbf{W}_i^{L,n+1}) + \mathbf{P}_{ij}^{q_1}. \end{cases} \quad (5.7.11)$$

This guarantees that  $\Psi_3(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij}) \geq 0$  and  $\Psi_4(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij}) \geq 0$  for all  $\ell \in [0, \ell_j^{i,q_1}]$ . This enforces a local minimum principle and a local maximum principle on  $q_1$ . As a corollary this also enforces positivity of  $\mathbf{Q}_{1,i}^{n+1}$ .

*Remark 5.7.6* (FCT limiting on linear functionals). It is also possible to use the FCT methodology for limiting the linear functionals  $\Psi_1, \dots, \Psi_4$ . We refer the reader to [28] where this is shown for the Saint-Venant model.

We now move on to the kinetic energy functional  $\Psi_5^{i,n}$ . We seek an  $\ell_j^{i,K} \in [0, \ell_j^{i,q_1}]$  such that  $\Psi_5^{i,n}(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij}) \geq 0$  for all  $\ell \in [0, \ell_j^{i,K}]$ . Let us define the functional:  $\Phi(\mathbf{U}) := \mathbf{H} \Psi_5^{i,n}(\mathbf{U}) = \mathbf{H} \mathbf{K}_i^{n,\max} - \frac{1}{2} \|\mathbf{Q}\|_{\ell^2}^2$ . Notice that  $\Psi_5^{i,n}(\mathbf{U}) \geq 0$  iff  $\Phi(\mathbf{U}) \geq 0$  provided  $\mathbf{H} > 0$ . Hence, assuming that  $\Psi_5^{i,n}(\mathbf{W}_i^{L,n+1} + \mathbf{P}_{ij}) < 0$  (otherwise there is nothing to optimize), our optimization problem consists of finding the unique  $\ell \in [0, 1)$  such that  $\Phi(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij}) = 0$ . But  $\Phi(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij})$  is

a quadratic functional with respect to  $\ell$ :  $\Phi(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij}) = a\ell^2 + b\ell + c$ , where

$$a = -\frac{1}{2} \|\mathbf{P}_{ij}^q\|_{\ell^2}^2, \quad (5.7.12a)$$

$$b = K_i^{n,\max} \mathbf{P}_{ij}^h - \mathbf{Q}(\mathbf{W}_i^{L,n+1}) \cdot \mathbf{P}_{ij}^q, \quad (5.7.12b)$$

$$c = H(\mathbf{W}_i^{L,n+1}) K_i^{n,\max} - \frac{1}{2} \|\mathbf{Q}(\mathbf{W}_i^{L,n+1})\|_{\ell^2}^2. \quad (5.7.12c)$$

Let  $t_0$  be the smallest positive root of the equation  $at^2 + bt + c = 0$ , with the convention that  $t_0 = 1$  if the equation has no positive root. Then, we choose  $\ell_j^{i,K}$  to be such that

$$\ell_j^{i,K} = \min(t_0, \ell_j^{i,q_1}). \quad (5.7.13)$$

It is proved in [31] that the definition (5.7.13) guarantees that  $\Psi_5^{i,n}(\mathbf{W}_i^{L,n+1} + \ell \mathbf{P}_{ij}) \geq 0$  for all  $\ell \in [0, \ell_j^{i,K}]$ . This enforces a local maximum principle on the kinetic energy.

Finally, we set

$$\ell_{ij} = \min(\ell_j^{i,K}, \ell_i^{j,K}). \quad (5.7.14)$$

Then with the above definition and by Lemma 5.7.5, the update  $\mathbf{W}_i^{n+1}$  computed by (5.7.8) satisfies the following constraints  $\Psi_l^{i,n}(\mathbf{W}_i^{n+1}) \geq 0$  for all  $l \in \mathcal{L}$ . We now “put back” the sources to compute the final limited update  $\mathbf{U}_i^{n+1}$ :

$$\mathbf{U}_i^{n+1} = \tau_n (\mathbf{R}_i^n + \mathbf{S}_i^n) + \sum_{j \in \mathcal{I}^*(i)} \theta_j (\mathbf{W}_i^{L,n+1} + \ell_{ij} \mathbf{P}_{ij}^n). \quad (5.7.15)$$

*Remark 5.7.7* (Saint-Venant limiting). To apply the limiting process described above to the Saint-Venant model, one just has to skip the limiting on the  $q_1$  functionals (i.e.,  $\Psi_4$  and  $\Psi_5$ ).

**Theorem 5.7.8** (Invariant domain preserving property for limited scheme). *Let  $i \in \mathcal{V}$  and  $n \geq 0$ . Assume that  $\mathbf{U}_j^n \in \mathcal{A} := \{\mathbf{u} \in \mathbb{R}^{d+4} \mid \mathbf{h} > 0\}$  for all  $j \in \mathcal{I}(i)$ . Suppose that the time step  $\tau_n$  is small enough so that  $\tau_n \max \left( \frac{2}{m_i} \sum_{j \in \mathcal{I}^*(i)} d_{ij}^{L,n}, \frac{\sqrt{gH_0}}{\mathcal{E}_i} \right) \leq 1$ . Let  $\mathbf{W}_i^{n+1}$  be defined by (5.7.8) with the limiter  $\ell_{ij}$  given by (5.7.14). Then  $\mathbf{W}_i^{n+1} \in \mathcal{A}$ . Consequently, the full update  $\mathbf{U}_i^{n+1}$  defined*



by (5.7.15) is in  $\mathcal{A}$  as well.

*Proof.* By construction, the definition (5.7.14) along with Lemma 5.7.5 gives

$$\mathbf{H}(\mathbf{W}_i^{n+1}) := \mathbf{H}\left(\sum_{j \in \mathcal{I}^*(i)} \theta_j \left(\mathbf{W}_i^{L,n+1} + \ell_{ij} \mathbf{P}_{ij}^n\right)\right) \geq \mathbf{h}_i^{n,\min},$$

for all  $i \in \mathcal{V}$  and all  $j \in \mathcal{I}(i)$ . The goal is to show that the limited water height update  $\mathbf{H}_i^{n+1}$  stays positive with the contribution of the source  $\tau_n(\mathbf{R}_i^n + \mathbf{S}_i^n)$ . For  $i \in \mathcal{V}$ , consider the update for  $\mathbf{H}_i^{n+1}$ :

$$\begin{aligned} \mathbf{H}_i^{n+1} &= \tau_n \chi\left(\frac{\mathbf{h}_i^{i,\max} - \mathbf{h}_i^{i,\min}}{\mathbf{h}_{\text{wave}}^{\min}}\right) \left(\frac{\sqrt{g\mathbf{H}_0}}{\mathcal{E}_i} (\mathbf{h}_{\text{wave}}(\mathbf{a}_i, t^n) - \mathbf{H}_i^n) G\left(\frac{\mathbf{a}_i - x_{\min}}{L_{\text{gen}}}\right)\right) + \mathbf{H}(\mathbf{W}_i^{n+1}), \\ &\geq \tau_n \frac{\sqrt{g\mathbf{H}_0}}{\mathcal{E}_i} \chi\left(\frac{\mathbf{h}_i^{i,\max} - \mathbf{h}_i^{i,\min}}{\mathbf{h}_{\text{wave}}^{\min}}\right) G\left(\frac{\mathbf{a}_i - x_{\min}}{L_{\text{gen}}}\right) (\mathbf{h}_{\text{wave}}^{\min} - \mathbf{h}_i^{n,\max}) + \mathbf{h}_i^{n,\min}, \end{aligned}$$

where we used the fact that  $\mathbf{h}_i^{n,\max} \geq \mathbf{H}_i^n$ . If the cut-off function  $\chi$  is active (i.e., when  $\mathbf{h}_i^{n,\max} - \mathbf{h}_i^{n,\min} \geq \mathbf{h}_{\text{wave}}^{\min}$ ), then  $\chi = 0$  and  $\mathbf{H}_i^{n+1} \geq \mathbf{h}_i^{n,\min}$  and thus positive. If the cut-off function  $\chi$  is not active (i.e., when  $\mathbf{h}_i^{n,\max} - \mathbf{h}_i^{n,\min} < \mathbf{h}_{\text{wave}}^{\min}$ ), then

$$\begin{aligned} \mathbf{H}_i^{n+1} &\geq \tau_n \frac{\sqrt{g\mathbf{H}_0}}{\mathcal{E}_i} \chi\left(\frac{\mathbf{h}_i^{i,\max} - \mathbf{h}_i^{i,\min}}{\mathbf{h}_{\text{wave}}^{\min}}\right) G\left(\frac{\mathbf{a}_i - x_{\min}}{L_{\text{gen}}}\right) (\mathbf{h}_{\text{wave}}^{\min} - \mathbf{h}_i^{n,\max}) + \mathbf{h}_i^{n,\min} \\ &\geq \tau_n \frac{\sqrt{g\mathbf{H}_0}}{\mathcal{E}_i} \chi\left(\frac{\mathbf{h}_i^{i,\max} - \mathbf{h}_i^{i,\min}}{\mathbf{h}_{\text{wave}}^{\min}}\right) G\left(\frac{\mathbf{a}_i - x_{\min}}{L_{\text{gen}}}\right) (-\mathbf{h}_i^{n,\min}) + \mathbf{h}_i^{n,\min} \\ &= \mathbf{h}_i^{n,\min} \left(1 - \tau_n \frac{\sqrt{g\mathbf{H}_0}}{\mathcal{E}_i} \chi\left(\frac{\mathbf{h}_i^{i,\max} - \mathbf{h}_i^{i,\min}}{\mathbf{h}_{\text{wave}}^{\min}}\right) G\left(\frac{\mathbf{a}_i - x_{\min}}{L_{\text{gen}}}\right)\right) \\ &\geq \mathbf{h}_i^{n,\min} \left(1 - \tau_n \frac{\sqrt{g\mathbf{H}_0}}{\mathcal{E}_i}\right). \end{aligned}$$

Thus,  $\mathbf{H}_i^{n+1}$  is positive under the CFL condition.  $\square$

**Proposition 5.7.9** (Well-balancing property for limited scheme). *Let  $\mathbf{T}_{HS} : \mathbf{u}_h^h \mapsto \mathbf{T}_{HS}(\mathbf{u}_h^n) := \mathbf{u}_h^{n+1}$  be the high-order scheme defined by (5.7.15) for the Hyperbolic Serre model. This scheme is exactly well-balanced if  $\mathbf{S}_G \equiv \mathbf{0}$ .*

*Proof.* Assume that  $\mathbf{u}_h^n$  is exactly at rest, then one can verify that  $\mathbf{P}_{ij} = \mathbf{0}$  for all  $i \in \mathcal{V}$  and all  $j \in \mathcal{I}^*(i)$ . Hence,  $\mathbf{u}_h^{n+1}$  is equal to the low-order update. We conclude by invoking Proposition 5.4.3.

□

*Remark 5.7.10* (Iterative limiting). We note here that the limiting process described above can be iterated multiple times by observing from (5.7.6) that

$$\mathbf{W}_i^{\mathbf{H},n+1} = \mathbf{W}_i^{\mathbf{L},n+1} + \frac{1}{m_i} \sum_{j \in \mathcal{I}(i)} \ell_{ij} \mathbf{A}_{ij}^n + \frac{1}{m_i} \sum_{j \in \mathcal{I}(i)} (1 - \ell_{ij}) \mathbf{A}_{ij}^n.$$

Then, by setting  $\mathbf{W}^{(0)} := \mathbf{W}_i^{\mathbf{L},n+1}$  and  $\mathbf{A}_{ij}^{(0)} = \mathbf{A}_{ij}^n$ , the iterative limiting process is shown in Algorithm 1. In the numerical simulations reported in Chapter 6, we take  $k_{\max} = 2$ .

---

**Algorithm 1** Iterative limiting with sources

---

**Input:**  $\mathbf{W}_i^{\mathbf{L},n+1}$ ,  $\mathbf{A}_{ij}^n$ ,  $k_{\max}$

**Output:**  $\mathbf{U}^{n+1}$

Set  $\mathbf{W}^{(0)} := \mathbf{W}_i^{\mathbf{L},n+1}$  and  $\mathbf{A}_{ij}^{(0)} = \mathbf{A}_{ij}^n$

**for**  $k = 0$  to  $k_{\max} - 1$  **do**

    Compute limiter  $\ell^{(k)}$

    Update  $\mathbf{W}^{(k+1)} = \mathbf{W}^{(k)} + \frac{1}{m_i} \sum_{j \in \mathcal{I}(i)} \ell_{ij}^{(k)} \mathbf{A}_{ij}^{(k)}$

    Update  $\mathbf{A}_{ij}^{(k+1)} = (1 - \ell_{ij}^{(k)}) \mathbf{A}_{ij}^{(k)}$

**end**

$\mathbf{U}^{n+1} = \mathbf{W}^{(k_{\max})} + \tau_n(\mathbf{R}^n + \mathbf{S}^n)$

---

#### 5.7.4 Relaxation of the bounds

The methodology described above leads to second-order accuracy in the  $L^1$ -norm, but the bounds defined in (5.5.4) and (5.5.5) are too tight to make the method higher-order or even second-order in the  $L^\infty$ -norm in the presence of smooth extrema. A similar phenomena was observed in Khobalatte and Perthame [42] in the context of the Euler Equations for gas dynamics and explained in Guermond et al. [31, Sec. 4.7]. We want to avoid this reduction in accuracy since

we want to model smooth solitary waves and periodic waves with the Hyperbolic Serre system (2.3.1)–(2.3.2). To recover the full accuracy in the  $L^\infty$ -norm, one should *relax* the bounds (5.5.4) and (5.5.5) for smooth solutions. We briefly discuss the relaxation here and refer the reader to [32, Sec. 7.6] where this is discussed in detail. All the numerical results reported in Chapter 6 are done using this relaxation technique.

The relaxation is done as follows. Let  $\Psi_i^{\min}$  and  $\Psi_i^{\max}$  denote either a min or max bound defined in §5.7.1 and  $\Psi_l(\mathbf{u})$  the associated functional; for example,  $\Psi_i^{\min} = h_i^{n,\min}$  and  $\Psi_1(\mathbf{u}) = h$ . For each  $l \in \mathcal{L}$  and all  $i \in \mathcal{V}$ , we set

$$(\Delta^2 \Psi_l)_i = \frac{1}{\sum_{j \in \mathcal{I}^*(i)} \beta_{ij}} \sum_{j \in \mathcal{I}^*(i)} \beta_{ij} \left( \Psi_l(\mathbf{U}_j^n) - \Psi_l(\mathbf{U}_i^n) \right), \quad (5.7.16)$$

where the coefficients  $\beta_{ij}$  are meant to make the computation linearity-preserving (see: [32, Rem. 6.2]).

In the computations reported below, we take  $\beta_{ij} = \int_D \nabla \varphi_i \cdot \nabla \varphi_j \, dx$ . We then compute an average of the above quantity

$$\overline{(\Delta^2 \Psi_l)_i} = \frac{1}{2(\text{card}(\mathcal{I}(i)) - 1)} \sum_{j \in \mathcal{I}^*(i)} \left( \frac{1}{2} (\Delta^2 \Psi_l)_i + \frac{1}{2} (\Delta^2 \Psi_l)_j \right), \quad (5.7.17)$$

Finally, we define the relaxation of  $\Psi^{i,\min}$  and  $\Psi^{i,\max}$  by redefining them as follows:

$$\Psi^{i,\min} \leftarrow \max \left( (1 - s_i) \Psi^{i,\min}, \Psi^{i,\min} - \overline{(\Delta^2 \Psi_l)_i} \right), \quad (5.7.18a)$$

$$\Psi^{i,\max} \leftarrow \min \left( (1 + s_i) \Psi^{i,\max}, \Psi^{i,\max} + \overline{(\Delta^2 \Psi_l)_i} \right), \quad (5.7.18b)$$

where  $s_i = \left( \frac{m_i}{|D|} \right)^{\frac{3}{2d}} \in (0, 1)$ . Here  $|D|$  is the measure of the computational domain (surface area if  $d = 2$ , or length if  $d = 1$ ).

In the applications reported in Chapter 6 that involve periodic waves, we observed numerically that the above relaxation for the (negative) kinetic energy can still be too restrictive. To remedy

this, we consider the following less restrictive relaxation for limiting the (negative) kinetic energy:

$$\Psi^{i,\max} \leftarrow \min \left( (1 + s_i^{\frac{2}{3}}) \Psi^{i,\max}, 0 \right). \quad (5.7.19)$$

## 6. NUMERICAL ILLUSTRATIONS

### 6.1 Introduction

In this Chapter, we conclude the thesis by illustrating the performance of the hyperbolic relaxed Serre model (4.4.1) and the associated proposed convex limiting method described in Chapter 5 in both spatial dimensions ( $\mathbb{R}^d$  with  $d = \{1, 2\}$ ). We focus on numerical illustrations for the Hyperbolic Serre model (4.4.1), but consider tests with the Saint-Venant model as well. In particular, we verify the accuracy of the method using analytical solutions of the Serre model and then again with the method of manufactured solutions. We then show the method is well-balanced with and without the effects of topography. We perform several numerical tests involving the Riemann Problem for the hyperbolic relaxed system (4.4.1). We also reproduce several academic benchmarks proposed in the literature and then introduce some new benchmarks as well. We conclude the chapter by reproducing several laboratory experiments that validate the hyperbolic relaxed model.

The chapter is organized as follows. In Section 6.2, we discuss the implementation details for the numerical method. Then in Section 6.3, we verify the accuracy of the convex limiting numerical method (5.7.15). In Section 6.4, we verify the numerical method is indeed well-balanced for both models. In Section 6.5, we investigate the numerical solution of the Riemann problem for the Serre model and compare with the Saint-Venant solution. Then in Section 6.6, we reproduce different benchmarks introduced in the literature and propose new benchmarks with interesting properties. Finally, in Section 6.7, we reproduce several laboratory experiments seen in the literature.

### 6.2 Preliminaries

The simulations reported in the paper are done in  $\mathbb{R}^d$  with  $d = \{1, 2\}$  using continuous, linear finite elements. When  $d = 1$ , we use a uniform grid. Some two-dimensional tests are done with continuous  $\mathbb{P}_1$  finite elements on unstructured Delaunay meshes, and some tests are done using continuous  $\mathbb{Q}_1$  finite elements on quadrangular meshes. In all the tests, the relaxation is done with  $\bar{\lambda} = 1$  and  $\mathcal{E}_i = m_i^{\frac{1}{d}}$ , for all  $i \in \mathcal{V}$ .

For the numerical tests in  $\mathbb{R}^2$ , three different codes implementing the method described in the paper have been written to ensure reproducibility. The first code, henceforth referred to as TAMU, does not use any particular software and is written in Fortran 95/2003. All 1D results are done with the TAMU code. The second code has been written at the US Army Engineer Research and Development Center using the `Proteus` toolkit (the reader is referred to Kees and Farthing [41]). Both codes use continuous  $\mathbb{P}_1$  Lagrange elements on triangles and unstructured, non-nested, Delaunay meshes. The third code is `Ryujin` [49, 34], a high-performance finite-element solver based on the `deal.II` library Arndt et al. [3] and uses continuous  $\mathbb{Q}_1$  elements. The time stepping in all three codes is done with the third-order, three stage, strong stability preserving Runge-Kutta method, SSP RK(3,3).

### 6.3 Convergence tests

In this section we verify the accuracy of the method by performing several convergence tests with the TAMU code described in 6.2. Similar tests have been performed with the `Proteus` and `Ryujin`, but are omitted here for the sake of brevity. The goal of this section is to show that when the relaxation parameter  $\epsilon$  of the Hyperbolic Serre model (4.4.1) is chosen to be proportional to the mesh-size, the numerical solution of this model converges to that of the original Serre model (2.3.1).

#### 6.3.1 Solitary wave solution of Serre model (2.3.1)

The Serre model (2.3.1) admits an exact solution in the form of a solitary wave propagating over a flat bottom. The goal of this particular test is to show that solutions of the relaxed model (4.4.1) converge with a first-order rate with respect to the mesh-size to solutions of the Serre model (2.3.1). Note that the hyperbolic relaxed model (4.4.1) does *not* analytically support exact solitary waves; more on this is discussed in §6.3.2.

Let  $\tilde{h}(x, t)$  and  $\tilde{u}(x, t)$  be the water height and velocity of an exact solitary wave:

$$\tilde{h}(x, t) = h_0 + \frac{\alpha}{(\cosh(r(x - x_0 - ct)))^2}, \quad \tilde{u}(x, t) = c \frac{\tilde{h}(x, t) - h_0}{\tilde{h}(x, t)}, \quad (6.3.1)$$

I	$E_1$		$E_2$		$E_\infty$	
100	1.546E-04	Rate	3.796E-04	rate	8.383E-04	Rate
200	4.509E-05	1.78	7.910E-05	2.26	1.272E-04	2.72
400	3.141E-05	0.52	7.258E-05	0.12	3.025E-04	-1.25
800	2.317E-05	0.44	5.509E-05	0.4	2.396E-04	0.34
1600	1.369E-05	0.76	3.155E-05	0.8	1.376E-04	0.8
3200	7.533E-06	0.86	1.733E-05	0.86	7.551E-05	0.87
6400	3.815E-06	0.98	8.869E-06	0.97	3.871E-05	0.96

Table 6.1: Convergence rates for solitary wave solution  $T = 50$  s,  $CFL = 0.05$ .

with wave speed  $c = \sqrt{g(h_0 + \alpha)}$  and width  $r = \sqrt{\frac{3\alpha}{4h_0^2(h_0 + \alpha)}}$ . We initialize the water height and discharge by setting

$$h(x, 0) = \max\{\tilde{h}(x, 0) - z(x), 0\}, \quad q(x, 0) = h(x, 0)\tilde{u}(x, 0), \quad (6.3.2)$$

with  $z(x) \equiv 0$ . Note that the variables  $q_1, q_2$ , and  $q_3$  are initialized with (6.4).

We consider a 1D uniform grid on the domain  $D = (0, 1000 \text{ m})$ . We set  $h_0 = 10 \text{ m}$  and  $\alpha = 1 \text{ m}$ . The solitary wave is initiated at  $x_0 = 250 \text{ m}$ . The final time is set to  $T = 50 \text{ s}$  which allows the solitary wave to travel approximately  $519.4 \text{ m}$ . Dirichlet conditions are set at the boundaries for all variables. In Table 6.1, we show the numerical results for this problem. The number of grid points is shown in the leftmost column. The relative errors on the water height measured in the  $L^1$ -norm,  $E_1 := \|h - h_h\|_{L^1}/\|h\|_{L^1}$ ,  $L^2$ -norm,  $E_2 := \|h - h_h\|_{L^2}/\|h\|_{L^2}$ , and  $L^\infty$ -norm,  $E_\infty := \|h - h_h\|_{L^\infty}/\|h\|_{L^\infty}$ , are shown in the second, third and fourth columns. We observe that the method converges to the solution of the Serre model (2.3.1) with first-order rate with respect to the mesh-size. The first-order rate is a consequence of the relaxation parameter in the relaxed system (4.4.1) being proportional to the local mesh-size. In Figure 6.1, we plot the final solution at  $T = 50 \text{ s}$ . The solid blue profile represents the elevation  $h(x) + z(x)$ .

We now use this analytical solution to verify the accuracy of the numerical method using the different mathematical entropies defined by (4.5.10) and (2.2.3) for the entropy-viscosity graph coefficients (5.6.5) and (5.6.6). The goal of this test is to show that the numerical method 5.7.15

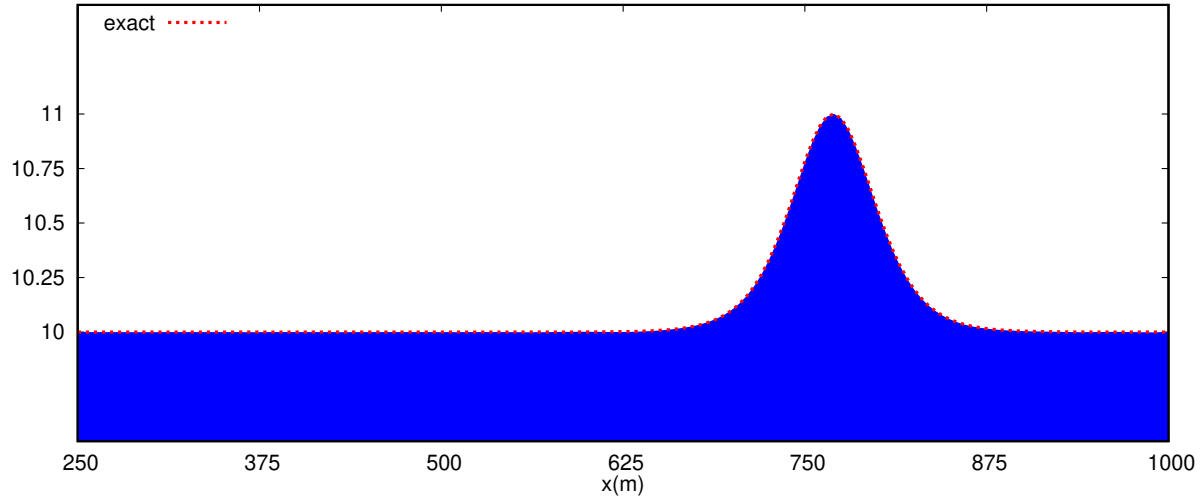


Figure 6.1: Computational solitary wave solution at  $T = 50$  s with  $I = 200$ .

behaves similarly with either entropy. We use the same computational set-up as above. In Table 6.2, we compare the results using (i) the Galerkin method (i.e., no artificial viscosity); (ii) the method (5.7.15) using the Shallow Water Equations entropy pair (2.2.3); (iii) the method (5.7.15) using the hyperbolic Serre entropy pair (4.5.10).

I	Galerkin		EV with (4.5.10)		EV with (2.2.3)	
100	2.80E-04	Rate	2.48E-04	rate	2.53E-04	Rate
200	4.24E-05	2.72	5.54E-05	2.16	5.64E-05	2.17
400	3.02E-05	0.49	3.74E-05	0.57	3.74E-05	0.59
800	2.32E-05	0.38	2.48E-05	0.59	2.48E-05	0.59
1600	1.39E-05	0.74	1.43E-05	0.79	1.43E-05	0.79
3200	7.67E-06	0.85	7.89E-06	0.86	7.89E-06	0.86
6400	3.84E-06	1.00	4.08E-06	0.95	4.08E-06	0.95

Table 6.2: Convergence table using  $\|h - h_h\|_{L^1} / \|h\|_{L^1}$  for solitary wave solution of Serre model (2.3.1).  $T = 50$  s, CFL = 0.05.

### 6.3.2 Method of manufactured solutions for Hyperbolic Serre model (4.4.1)

We now verify the second-order accuracy of the method (5.7.15). This is typically done by comparing a computational solution to a reference analytical solution and then looking at the be-



havior of the  $L^1$ ,  $L^2$ ,  $L^\infty$  errors as the mesh-size decreases. However, deriving exact solutions for the hyperbolic relaxation model (4.4.1) is a highly non-trivial task and thus no known solution exists. To test the accuracy of our numerical method for solving the hyperbolic model (4.4.1), we use the method of manufactured solutions. We discuss the details below.

Assume there are no external sources and assume the topography is flat (i.e.,  $\mathbf{S}(\mathbf{u}) = \mathbf{0}$  and  $\mathbf{R}(\mathbf{u}, \nabla z) = \mathbf{0}$ ). Let  $\tilde{h}(x, t)$  and  $\tilde{u}(x, t)$  be the same profiles defined by (6.3.1). We want to derive a source term  $\mathbf{S}_{\text{man}}(x, t)$  for the hyperbolic Serre model (4.4.1) such that

$$\partial_t \mathbf{u} + \nabla \cdot \mathbf{f}(\mathbf{u}) = \mathbf{S}_{\text{man}}(x, t),$$

holds exactly. To do so, we assume the following the profiles for the conserved variables  $\mathbf{u} = (\mathbf{h}, u, q_1, q_2, q_3)^\top$ :

$$\mathbf{h}(\mathbf{x}, t) = \tilde{h}(x, t), \quad q(x, t) = \tilde{h}(x, t)\tilde{u}(x, t), \quad (6.3.3a)$$

$$q_1(x, t) = \tilde{h}^2(x, t), \quad q_2(x, t) = -\tilde{h}^2(x, t)\partial_x(\tilde{u}(x, t)), \quad q_3(x, t) = 0. \quad (6.3.3b)$$

Substituting (6.3.3) into (4.4.1), we see that the manufactured source term should be of the form:

$$\mathbf{S}_{\text{man}}(x, t) = (0, q_{\text{man}}(x, t), 0, q_{2\text{man}}(x, t), 0)^\top, \quad (6.3.4a)$$

with

$$q_{\text{man}}(x, t) = 2r\alpha \left( \frac{2c^2\mathbf{h}_0^2 \sinh(2r(x - x_0 - ct))}{(\mathbf{h}_0 + 2\alpha + \mathbf{h}_0 \cosh(2r(x - x_0 - ct)))^2} - \frac{g}{\alpha} \tilde{h}(x, t)(\tilde{h}(x, t) - \mathbf{h}_0) \tanh(r(x - x_0 - ct)) \right), \quad (6.3.4b)$$

$$q_{2\text{man}}(x, t) = \frac{2c^2\mathbf{h}_0^2 r^2 \alpha (-3\mathbf{h}_0 - 4\alpha - 2\mathbf{h}_0 \cosh(2r(x - x_0 - ct)) + \frac{\mathbf{h}_0}{\alpha} \cosh(4r(x - x_0 - ct))) (\tilde{h}(x, t) - \mathbf{h}_0)}{(\mathbf{h}_0 + 2\alpha + \mathbf{h}_0 \cosh(2r(x - x_0 - ct)))^2} \quad (6.3.4c)$$

The profiles defined above were obtained using the `Mathematica` software [65].

The computational domain is set to  $D = (0, 1000 \text{ m})$ . We initialize the conserved variables with (6.3.3) at  $t = 0 \text{ s}$  with  $h_0 = 10 \text{ m}$ ,  $\alpha = 0.1h_0$ ,  $x_0 = 250 \text{ m}$ . The final time is set to  $T = 50 \text{ s}$  and  $\text{CFL} = 0.05$ . Dirichlet conditions are set at the boundaries for all variables. We show in Table 6.3 the numerical results obtained at  $T = 50 \text{ s}$ . We observe that that all the quantities converge with second-order rate with respect to the mesh-size as expected.

I	$E_1$		$E_2$		$E_\infty$	
		Rate		rate		Rate
100	1.76E-04		4.50E-04		1.35E-03	
200	4.98E-05	1.82	1.23E-04	1.86	5.30E-04	1.35
400	1.61E-05	1.63	3.92E-05	1.66	1.77E-04	1.58
800	4.30E-06	1.90	1.03E-05	1.93	4.68E-05	1.92
1600	1.11E-06	1.95	2.60E-06	1.99	1.19E-05	1.97
3200	3.10E-07	1.84	6.65E-07	1.97	3.05E-06	1.97

Table 6.3: Convergence rates using manufactured solution.  $T = 50 \text{ s}$ ,  $\text{CFL} = 0.05$

### 6.3.3 Steady-state solution with topography

We now verify the accuracy of the relaxation technique using the steady state solution described in Section 2.3.5.2. Recall that this is an analytical solution to the Serre model (2.3.1)–(2.3.2) with topography effects. In particular, the water height profile and bathymetry are given by:

$$h(x) = h_0 \left( 1 + \frac{a}{(\cosh(rx))^2} \right), \quad z(x) = -\frac{1}{2}(h(x) - h_0)$$

Here,  $h(x)$  describes a stationary solitary wave and the bathymetry is a depression (see Figure 2.2). We set  $h_0 = 1 \text{ m}$ ,  $a = 0.2$ ,  $g = 9.81 \text{ ms}^{-2}$ . This gives a discharge value of  $q = \sqrt{5.886} \text{ m}^2 \text{ s}^{-1}$  and coefficient  $r = \sqrt{0.5} \text{ m}^{-1}$  given by the expressions in (2.3.20). The simulations are done with  $D = (-10 \text{ m}, 15 \text{ m})$ . The discharge is enforced at the inflow boundary  $x = -10 \text{ m}$ . The water height is enforced at  $x = -10 \text{ m}$  and  $x = 15 \text{ m}$ . We show in Table 6.4 the numerical results obtained at  $t = 1000 \text{ s}$ . The number of grid points is shown in the leftmost column. The relative

errors on the water height measured in the  $L^1$ -norm,  $E_1 := \|\mathbf{h} - \mathbf{h}_h\|_{L^1} / \|\mathbf{h}\|_{L^1}$ , and the relative errors measured in the  $L^\infty$ -norm,  $E_2 := \|\mathbf{h} - \mathbf{h}_h\|_{L^\infty} / \|\mathbf{h}\|_{L^\infty}$ , are shown in the second and third columns. The relative  $L^1$ -norm of the difference between  $\mathbf{h}_h^2$  and  $q_{1h}$ ,  $E_3 := \|\mathbf{h}_h^2 - q_{1h}\|_{L^1} / \|q_{1h}\|_{L^1}$ , and the relative  $L^1$ -norm of the difference between  $q_h \partial_x z$  and  $q_{3h}$ ,  $E_4 := \|q_h \partial_x z - q_{3h}\|_{L^1} / \|q_h \partial_x z\|_{L^1}$ , are shown in the fourth and fifth columns. We observe that all the quantities converge with a first-order rate with respect to the mesh size. This is consistent since we chose the relaxation parameter in the relaxed system (4.4.1) to be proportional to the local mesh size which gives a first-order approximation of the fully coupled system (3.2.2).

I	$E_1$		$E_2$		$E_3$		$E_4$	
100	1.98E-03	Rate	7.55E-03	rate	8.54E-05	Rate	1.33E-01	Rate
200	1.09E-03	0.86	3.15E-03	1.26	3.83E-05	1.16	6.77E-02	0.97
400	4.23E-04	1.36	1.05E-03	1.58	1.76E-05	1.12	3.40E-02	0.99
800	1.73E-04	1.29	4.07E-04	1.37	8.51E-06	1.05	1.70E-02	1.00
1600	7.92E-05	1.13	1.82E-04	1.16	4.20E-06	1.02	8.51E-03	1.00
3200	3.62E-05	1.13	8.51E-05	1.10	2.09E-06	1.01	4.25E-03	1.00
6400	1.85E-05	0.97	4.31E-05	0.98	1.05E-06	1.00	2.13E-03	1.00

Table 6.4: Convergence rates table for steady state solution with topography.

## 6.4 Well-balancing tests

In this section, we verify that the scheme is well-balanced. To quantify the concept of well-balancing, we define the following error indicator for the hyperbolic Serre model:

$$\delta_\infty(t) := \frac{\|\mathbf{h}_h(t) - \mathbf{h}_0\|_{L^\infty(D)}}{H_0} + \frac{\|\mathbf{q}_h(t) - \mathbf{q}_0\|_{L^\infty(D)}}{H_0 \sqrt{gH_0}} + \frac{\|Q_{1,h}(t) - Q_{1,0}\|_{L^\infty(D)}}{H_0^2} + \frac{\|Q_{2,h}(t) - Q_{2,0}\|_{L^\infty(D)}}{H_0 \sqrt{gH_0}} + \frac{\|Q_{3,h}(t) - Q_{3,0}\|_{L^\infty(D)}}{H_0 \sqrt{gH_0}}, \quad (6.4.1)$$

where  $H_0$  is some reference water depth,  $\mathbf{h}_0, \mathbf{q}_0, Q_{1,0}, Q_{2,0}, Q_{3,0}$  are the initial states. The quantities  $\mathbf{h}_h(t), \mathbf{q}_h(t), Q_{1,h}(t), Q_{2,h}(t), Q_{3,h}(t)$  are the finite element approximations at time  $t$  for the respective conserved variables. We show that this quantity stays close to round-off error for the numerical

tests.

*Remark 6.4.1.* The error indicator for the Saint-Venant model is defined by:

$$\delta_\infty(t) := \frac{\|\mathbf{h}_h(t) - \mathbf{h}_0\|_{L^\infty(D)}}{H_0} + \frac{\|\mathbf{q}_h(t) - \mathbf{q}_0\|_{L^\infty(D)}}{H_0\sqrt{gH_0}} \quad (6.4.2)$$

For these tests, we consider the set up of the 1995 experiments by Briggs et al. [11] where the bathymetry is defined by a conical island (see: §6.7). The experimental domain is given by  $D = (0, 25 \text{ m}) \times (0, 30 \text{ m})$ . The experimental bathymetry is defined by:

$$z(\mathbf{x}) = \begin{cases} \min(0.625, 0.9 - r(\mathbf{x})/4), & r(\mathbf{x}) < 3.6 \\ 0, & \text{otherwise,} \end{cases}$$

where  $r(\mathbf{x})$  is the radius from the center of the island located at (12.96 m, 13.80 m).

We first consider the case where the initial profile is a complete wet state. The reference water depth is set to  $H_0 = 1 \text{ m}$  (above the island) and we define the initial water height to be  $h_0(\mathbf{x}) = H_0 - z(\mathbf{x})$  and initial flow rate  $\mathbf{q}_0 = 0 \text{ m/s}$ . The simulations are run until  $T = 50 \text{ s}$ . All the computations are done with  $\text{CFL} = 0.5$ . In Table 6.2a, we report the well-balancing quantity (6.4.1) for two different meshes and observe that indeed the values are near machine precision. We now consider the case where the initial state is a wet-dry state. We set the reference depth

$I$	$\delta_\infty(t)$
3587	2.5101E-13
14023	5.6512E-13

(a)  $H_0 = 1 \text{ m}$ .  $T = 50 \text{ s}$ .

$I$	$\delta_\infty(t)$
3587	7.7165E-12
14023	1.0692E-11

(b)  $H_0 = 0.32 \text{ m}$ .  $T = 50 \text{ s}$ .

Figure 6.2: Error tables for well-balancing tests using conical island topography.

to be  $H_0 = 0.32 \text{ m}$  so that the water elevation intersects the cone at  $r(\mathbf{x}) = 2.32 \text{ m}$ . To initiate this problem properly, we begin exactly at rest with respect to the mesh. That is to say, the mesh

is *aligned* with the initial data in regions where  $h + z$  is constant. We show this refinement in Figure 6.3c. In Table 6.2b, we report the well-balancing quantity (6.4.1) for two different meshes and see that it also stays close to machine precision 0. We also compute the  $\delta_\infty$  error indicator

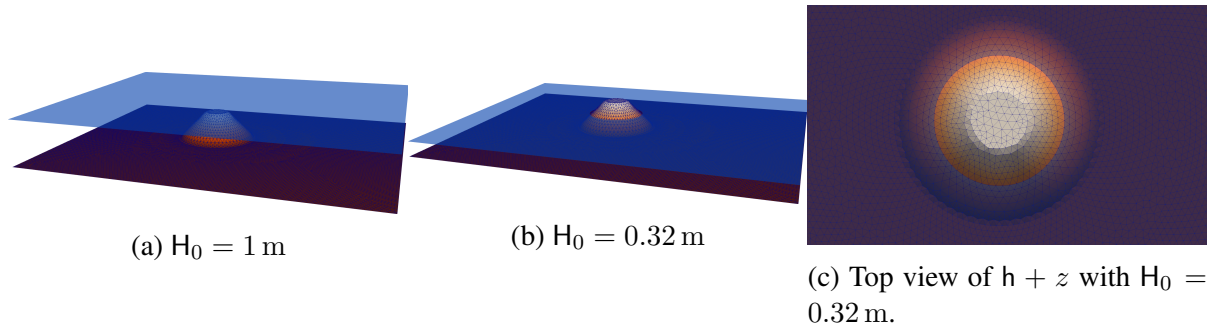


Figure 6.3: Figures for well-balancing tests with conical island topography.

for the Saint-Venant model. In Figure 6.4, we plot the  $\delta_\infty(t)$  indicator as a function of time over the interval  $t \in [0, 200 \text{ s}]$  for the Saint-Venant low/high-order schemes and the Hyperbolic Serre low/high-order schemes. We see that all quantities for each test stay close to machine precision 0 for long times.

### 6.5 Riemann problem for Serre Equations

In this section, we investigate numerically the solution to the Riemann Problem for the Serre model. It has been numerically observed that the solutions of the Riemann Problem can produce solutions in the form of dispersive shock waves which can have practical applications in studying undular bores. For a detailed study of the Riemann Problem for the Serre Equations, we refer the reader to El et al. [17] and Gavriluk et al. [23] and the references therein.

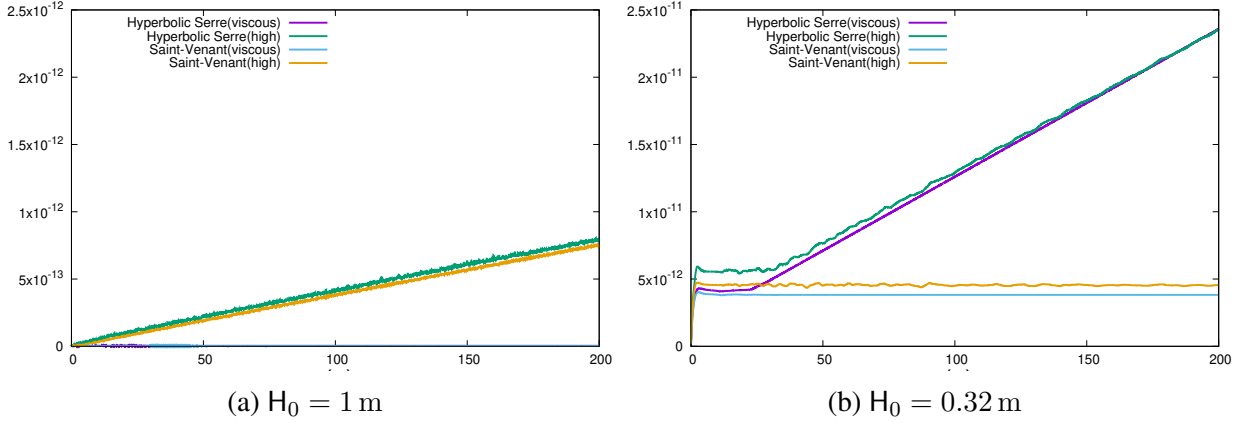


Figure 6.4: A plot of  $\delta_\infty(t)$  as a function of time for  $t \in [0, 200 \text{ s}]$  for the Hyperbolic Serre model and the Saint-Venant model for the low and high-order schemes.

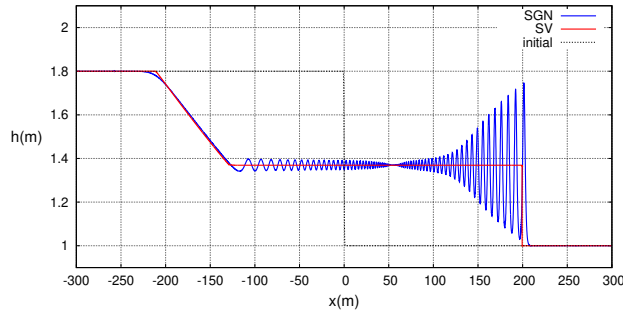
### 6.5.1 1D – Dam break over wet bed

We consider a one-dimensional dam-break problem introduced in El et al. [17, Sec. 7] with initial conditions:

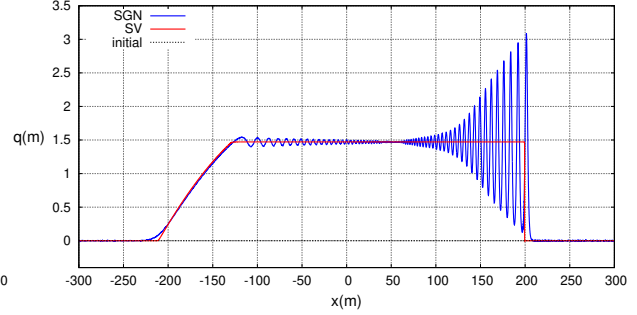
$$h(x, 0) = \begin{cases} 1.8 \text{ m}, & x < 0, \\ 1 \text{ m}, & x > 0, \end{cases}, \quad q(x, 0) = 0.$$

The computational domain is defined to be  $D = (-300, 300 \text{ m})$ . We set the mesh-size to be  $h = 0.05 \text{ m}$  corresponding to 12,000  $\mathbb{P}_1$  elements. The final time is set to  $T = 50 \text{ s}$  and the CFL number is chosen to be 0.05. In Figure 6.7, we show the computed water depth and momentum (respectively) of both the hyperbolic Serre–Green–Naghdi model (4.4.1) and the Saint-Venant model (2.2.1). The computed profiles for the hyperbolic Serre model (blue profiles) quantitatively look similar to that of El et al. [17, Fig. 8] and Tkachenko [60, Fig. 2.4]. We see that in Figure 6.5, the hyperbolic Serre model produces dispersive a shock wave traveling to the right and an expansion wave traveling to the left followed by a dispersive tail. This structure is fundamentally different than the Saint-Venant solution (red profiles) which produces a right-going shock and a left-going expansion wave.

We now modify the dam-break problem with a smaller jump in the initial water depth. That is,



(a) Case 1 – water depth



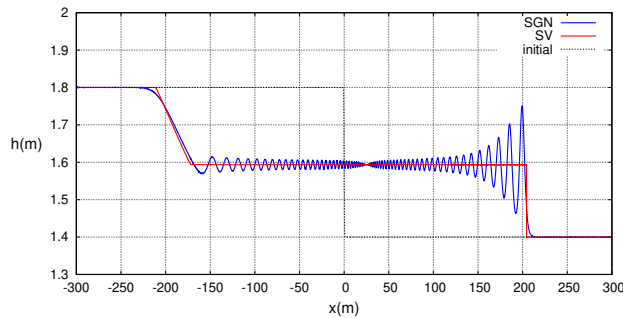
(b) Case 1 – momentum

Figure 6.5: Numerical solution to 1D dam-break problem with dry bed with  $\Delta h = 1.8$  m.

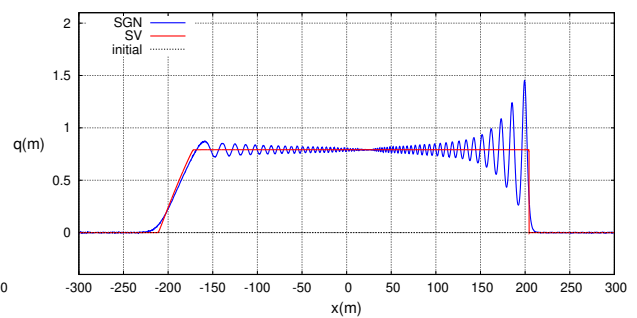
we set the initial conditions to be:

$$h(x, 0) = \begin{cases} 1.8 \text{ m}, & x < 0, \\ 1.4 \text{ m}, & x > 0, \end{cases}, \quad q(x, 0) = 0.$$

The computational set-up is the same as above. We show in Figure 6.6 the comparison of the water depth and flow discharge using both models. We see similar profiles as in Figure 6.5.



(a) Case 2 – water depth



(b) Case 2 – momentum

Figure 6.6: Numerical solution to 1D dam-break problem with a wet bed with  $\Delta h = 0.4$  m.

### 6.5.2 1D – Dam break over dry bed

We consider a one-dimensional dam-break over a dry bed. This test case highlights the ability of the proposed numerical methods to handle dry states. We set the initial conditions to be:

$$h(x, 0) = \begin{cases} 1.8 \text{ m}, & x < 0, \\ 0 \text{ m}, & x > 0, \end{cases}, \quad q(x, 0) = 0.$$

The computational domain is defined to be  $D = (-300, 300 \text{ m})$ . We set the mesh-size to be  $h = 0.05 \text{ m}$  corresponding to 12,000  $\mathbb{P}_1$  elements. The final time is set to  $T = 30 \text{ s}$  and the CFL number is chosen to be 0.05. In Figure 6.7, we show the computed water depth and momentum (respectively) of both the hyperbolic Serre–Green–Naghdi model (4.4.1) and the Saint-Venant model (2.2.1).

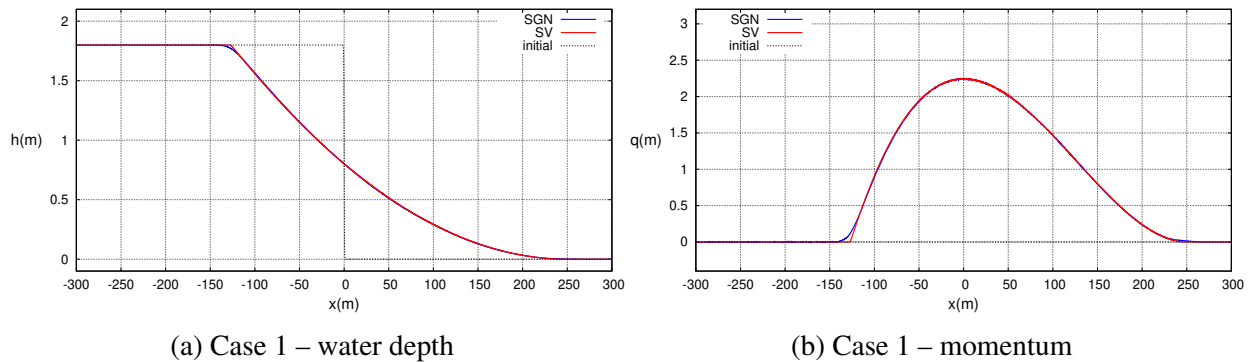


Figure 6.7: Numerical solution to 1D dam-break problem with a dry bed with  $\Delta h = 1.8 \text{ m}$ .

### 6.5.3 2D – Circular dam break

We now consider a two-dimensional circular dam-break problem. A similar test case was proposed in Tkachenko [60, Sec. 5.2]. Let the still water depth be  $H_0 = 1 \text{ m}$ . We initialize circular dam at rest with amplitude 0.8 m centered at  $(250 \text{ m}, 250 \text{ m})$  with radius  $R = 40 \text{ m}$ . That is to say, the initial condition is set to:



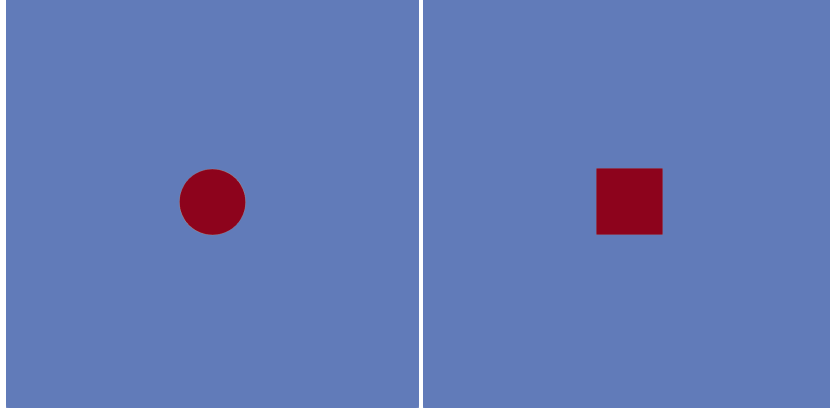


Figure 6.8: Initial profiles for circular and square Riemann Problems.

$$h(\mathbf{x}, 0) = \begin{cases} 1.8 \text{ m}, & \sqrt{(x - 250)^2 + (y - 250)^2} < R, \\ 1 \text{ m}, & \text{otherwise,} \end{cases}, \quad \mathbf{q}(\mathbf{x}, 0) = 0,$$

The initial condition is shown in Figure 6.8 (left). For this test case, we use the `Ryujin` software described in §6.2. The computational domain is set to  $D = (0, 500 \text{ m}) \times (0, 500 \text{ m})$  with 4,004,001  $\mathbb{Q}_1$  nodes. The final time is set to  $T = 50 \text{ s}$  with CFL number of 0.075. In Figure 6.9, we show the water depth profiles for both the hyperbolic Serre model (top) and the Saint-Venant model (bottom) at the times  $t = \{10, 30, 50\} \text{ s}$ .

#### 6.5.4 2D – Square dam break

We consider a two-dimensional square dam-break problem introduced in Tkachenko [60, Sec. 5.2]. Let the still water depth be  $H_0 = 1 \text{ m}$ . The square dam of side length  $d_s = 80 \text{ m}$  of amplitude 0.8 m is initialized as follows:

$$h(\mathbf{x}, 0) = \begin{cases} 1.8 \text{ m}, & |x - 250| \leq \frac{d_s}{2} \text{ and } |y - 250| \leq \frac{d_s}{2}, \\ 1 \text{ m}, & \text{otherwise,} \end{cases}, \quad \mathbf{q}(\mathbf{x}, 0) = 0,$$

The initial condition is shown in Figure 6.8 (right). For this test case, we use the `Ryujin` software described in §6.2. The computational domain is set to  $D = (0, 500 \text{ m}) \times (0, 500 \text{ m})$  with 4,004,001

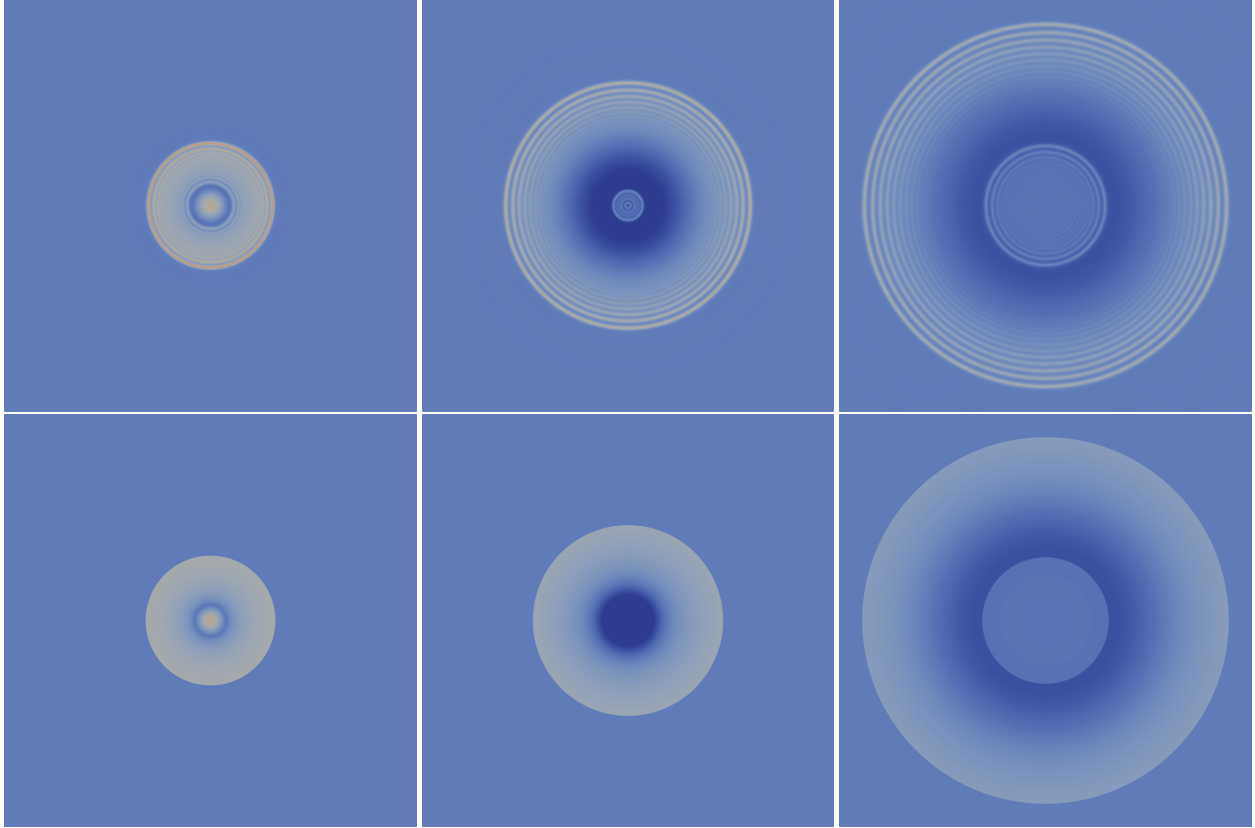


Figure 6.9: Circular dam break – Comparison with Hyperbolic Serre (top) and Saint-Venant (bottom) at  $t = \{10, 30, 50\}$ s.

$\mathbb{Q}_1$  nodes. The final time is set to  $T = 50$  s with CFL number of 0.075. In Figure 6.10, we show the water depth profiles for both the hyperbolic Serre model (top) and the Saint-Venant model (bottom) at the times  $t = \{10, 30, 50\}$ s.

## 6.6 Academic benchmarks

In this section, we reproduce several academic benchmarks seen in the literature. Since some of these benchmarks have not been reproduced using the Serre–Green–Naghdi equations, the goal is to highlight the computational results using the for such benchmarks. We also propose new synthetic benchmarks with interesting structures not seen in the literature.

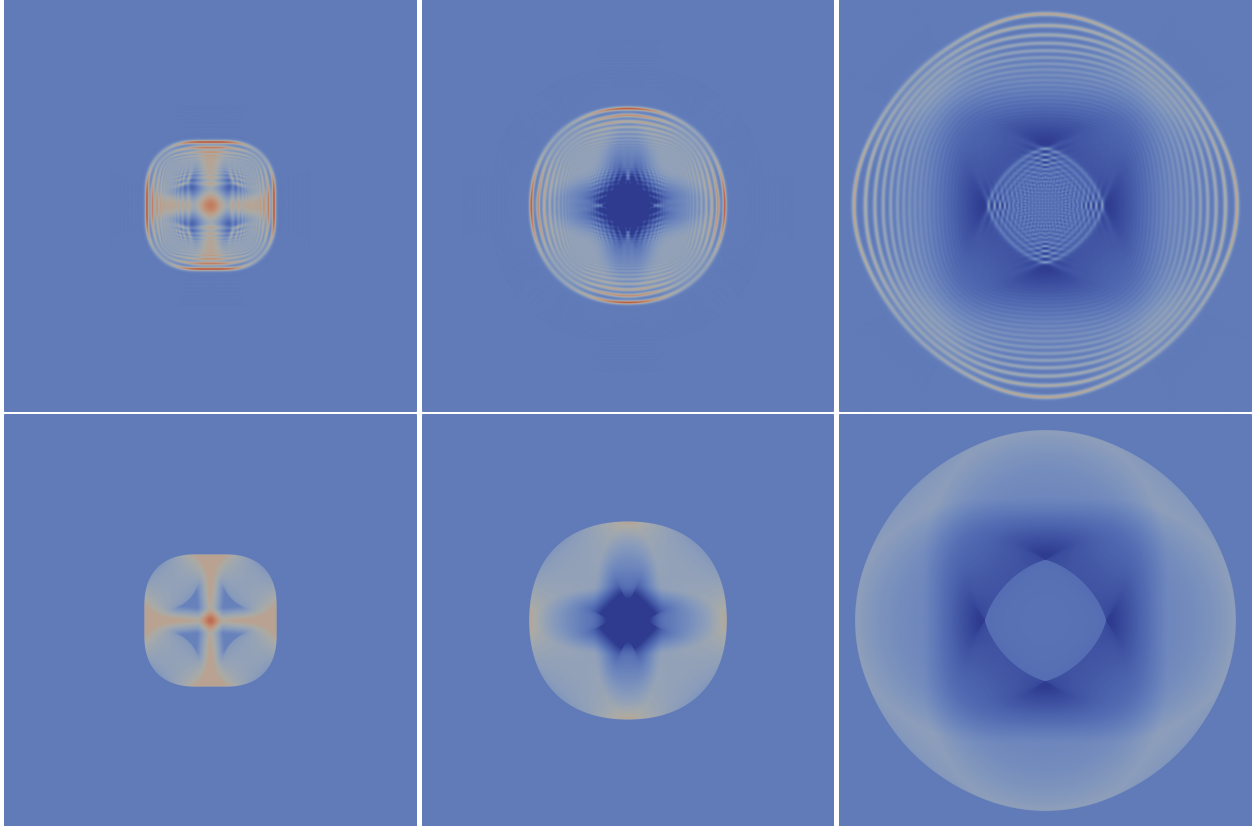


Figure 6.10: Square dam break – Comparison with Hyperbolic Serre (top) and Saint-Venant (bottom) at  $t = \{10, 30, 50\}$ s.

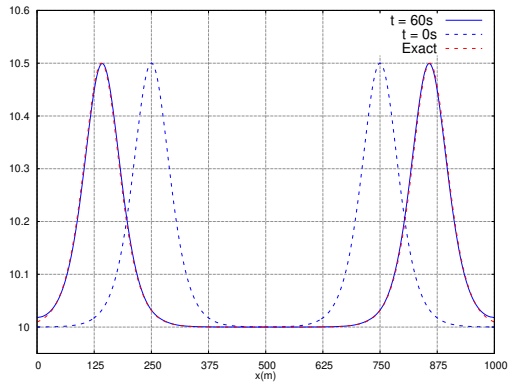
### 6.6.1 1D – Interaction of solitary waves

We now consider the collision of multiple solitary waves in one spatial dimension. It is known that the Serre Equations admit quasi-elastic collisions of solitary waves [16] when the amplitude is small (unlike the Korteweg–De Vries (KdV) equation which admits elastic collisions of solitary waves). We say a collision is elastic if the solitary waves are unchanged after the collision. For the Serre Equations, the collision of two solitary waves of small amplitude will produce two waves propagating in the opposite directions with similar shape and height with potentially a small phase shift and loss of amplitude (or none at all). When the amplitudes of the colliding waves are larger, there is a stronger interaction between the waves which leads to a larger phase shift and loss of amplitude and the creation of a dispersive tail. We call this type of collision an

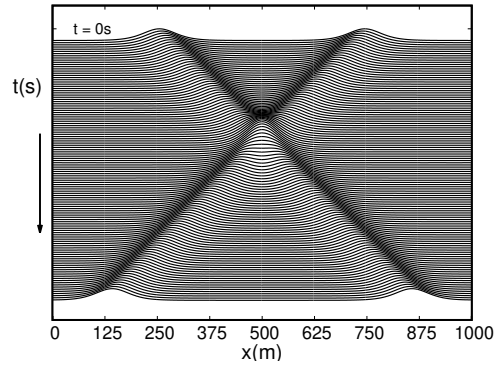
inelastic collision. For the study of such phenomena, we refer the reader to [16, 51] and references therein. Before discussing the different tests, note that for all computations in this section, the local mesh-size is set to  $h = 0.5$  m and the CFL number is set to 0.075.

The first test we consider is the collision of two solitary waves. In particular, we consider two different cases that show the phenomena of quasi-elastic collision (when the wave amplitude is small) and inelastic collision (when the wave amplitude is larger). The computational domain is set to  $D = (0, 1000$  m). To initialize the solitary waves, we use the profiles defined by (6.3.1). We set the reference depth to  $h_0 = 10$  m. For the quasi-elastic case, we set the amplitude to be  $\alpha = 0.05h_0 = 0.5$  m. The first solitary wave (which travels to the right) is initiated at  $x_1 = 250$  m and the second solitary wave (which travels to the left) is initiated at  $x_2 = 750$  m. For the second solitary wave, we switch the sign of the wave celerity in (6.3.1) (i.e.,  $c \rightarrow -c$ ) so that the wave travels to the left. The final time is set to  $T = 60$  s. We show in Figure 6.11 the results for the quasi-elastic collision. In particular, we compare the final numerical solution with the exact solution (at the final time) in Figure 6.11a and see that the profiles are close to each other. In Figure 6.11b, we give a space-time plot which highlights the trajectory of the waves. For the inelastic collision, we set the amplitude of both waves to be  $\alpha = 0.1h_0 = 1$  m. We show in Figure 6.12 the results for the inelastic collision. Again, we compare the numerical solution with the exact solution and see that the numerical solution has produced small dispersive tails and a slight phase shift was introduced.

The second test we consider is the collision of four solitary waves. Again we consider two cases that show the phenomena of quasi-elastic collision and inelastic collision. The computational domain is set to  $D = (0, 2000$  m). We set the reference depth to  $h_0 = 10$  m. For the quasi-elastic case, we set the amplitude of each wave to be  $\alpha = 0.05h_0 = 0.5$  m. The solitary waves are respectively initiated at  $x_1 = 250$  m,  $x_2 = 650$  m,  $x_3 = 1350$  m,  $x_4 = 1750$  m. The final time is set to  $T = 100$  s. We show in Figure 6.13 the results for the quasi-elastic case. We see again that the waves are hardly unchanged and closely match the final exact profiles. For the inelastic collision, we set amplitude of each wave (left to right respectively) to be  $\alpha_1 = 2$  m,  $\alpha_2 = 1$  m,  $\alpha_3 = 1$  m,  $\alpha_4 = 1$  m. It is clear from Figure 6.14 that the larger amplitudes introduce stronger

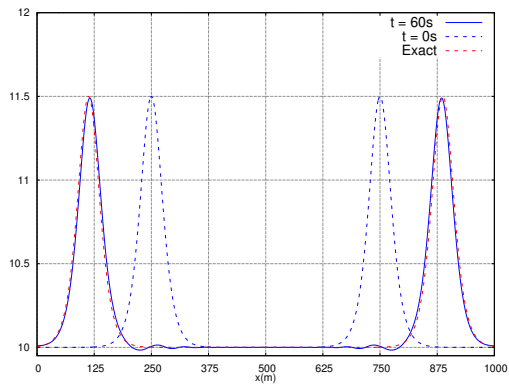


(a) Free surface elevation

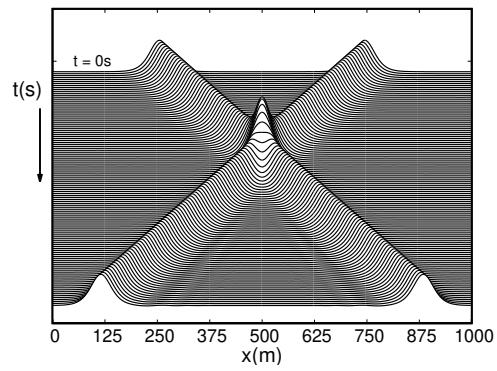


(b) Space-time plot

Figure 6.11: Two solitary waves – quasi-elastic collision

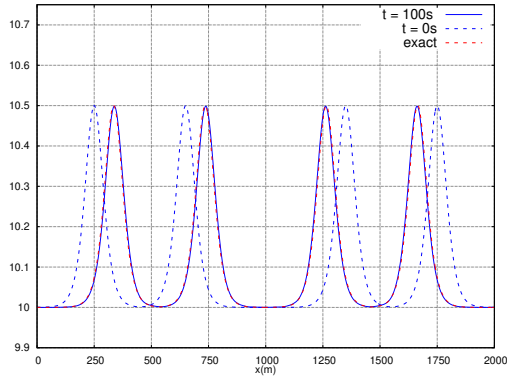


(a) Free surface elevation

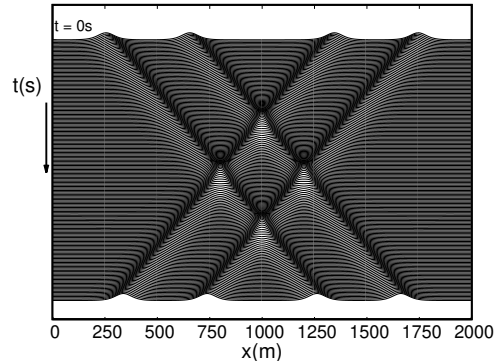


(b) Space-time plot

Figure 6.12: Two solitary waves – inelastic collision

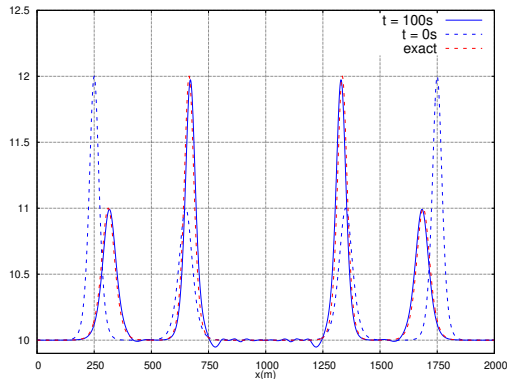


(a) Free surface elevation

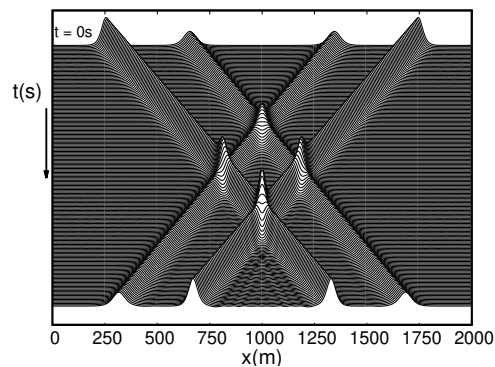


(b) Space-time plot

Figure 6.13: Four solitary waves – quasi-elastic collision



(a) Free surface elevation



(b) Space-time view

Figure 6.14: Four solitary waves – inelastic collision

interactions and affects the collision elasticity.

The next case we consider is the collision of 8 solitary waves. Since we have established the quasi-elasticity property of the Serre–Green–Naghdi equations, this test shown here is to highlight the robustness of the numerical method. The computational domain is set to  $D = (0, 4000 \text{ m})$ . We set the reference depth to  $h_0 = 10 \text{ m}$ . Letting the integer “1” denote the left-most solitary wave, the amplitudes of the solitary wave are as follows:  $\alpha_i = 2 \text{ m}$  for  $i = \{1, 4, 5, 8\}$  and  $\alpha_j = 1 \text{ m}$  for  $i = \{2, 3, 6, 7\}$ . The final time is set to  $T = 200 \text{ s}$ . We show the computations for this test in

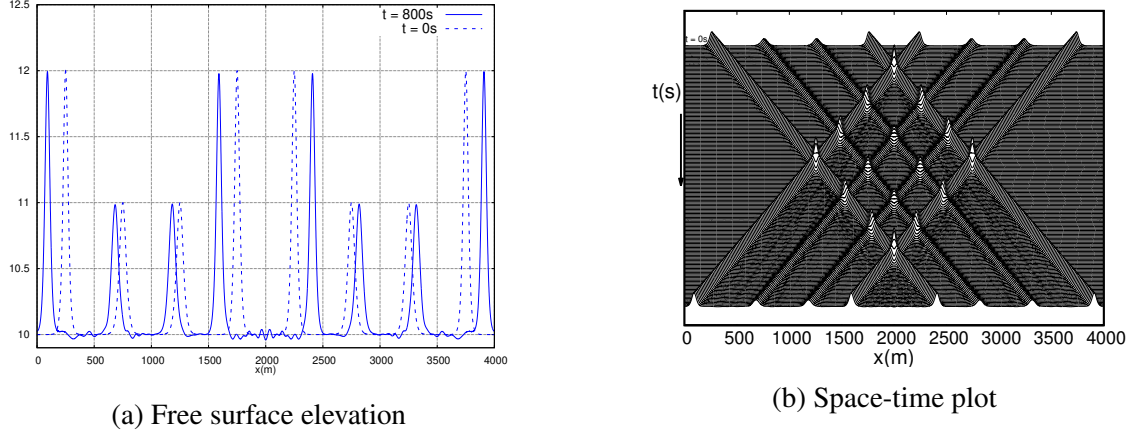


Figure 6.15: Collision of Eight solitary waves

Figure 6.15.

### 6.6.2 2D – Dam Break over three obstacles with friction

We now consider the test case of a dam break over a dry bottom with three conical obstacles introduced by Kawahara and Umetsu [38] (and reproduced by others: Guermond et al. [33], Huang et al. [37], etc.). This benchmark tests the complex wetting/drying process and the method's ability to handle the Gauckler-Manning friction source. Here the Manning's coefficient is set to  $n = 0.02 \text{ m}^{-1/3}\text{s}$ .

The domain is set to  $D = (0, 75 \text{ m}) \times (0, 30 \text{ m})$ . The bathymetry consisting of three conical obstacles is defined by  $z(\mathbf{x}) := \max\{0, z_1(\mathbf{x}), z_2(\mathbf{x}), z_3(\mathbf{x})\}$  where

$$z_1(\mathbf{x}) = 1 - \frac{1}{8} \sqrt{(x - 30)^2 + (y - 6)^2}, \quad (6.6.1a)$$

$$z_2(\mathbf{x}) = 1 - \frac{1}{8} \sqrt{(x - 30)^2 + (y - 24)^2}, \quad (6.6.1b)$$

$$z_3(\mathbf{x}) = 3 - \frac{3}{10} \sqrt{(x - 47.5)^2 + (y - 15)^2}. \quad (6.6.1c)$$

are three three obstacles. The initial state is set to

$$\mathbf{h}_0(\mathbf{x}) = \begin{cases} 1.875, & x \leq 16 \\ 0, & \text{otherwise,} \end{cases} \quad \mathbf{q}_0(\mathbf{x}) = \mathbf{0}.$$

For these computations, we use the `Ryujin` code described above. The mesh is composed of rectangular elements with 2,307,361  $\mathbb{Q}_1$  DOFS. The final time is set to  $T = 20$  s with CFL 0.075. In Figure 6.16, we show the computational free surface elevation  $h + z$  at several time snapshots. Then, in Figure 6.17, we show the comparison of the computations with the hyperbolic Serre model (4.4.1) and the Saint-Venant shallow water equations (2.2.1). We observe that more realistic structures are produced by the dispersive Serre model.

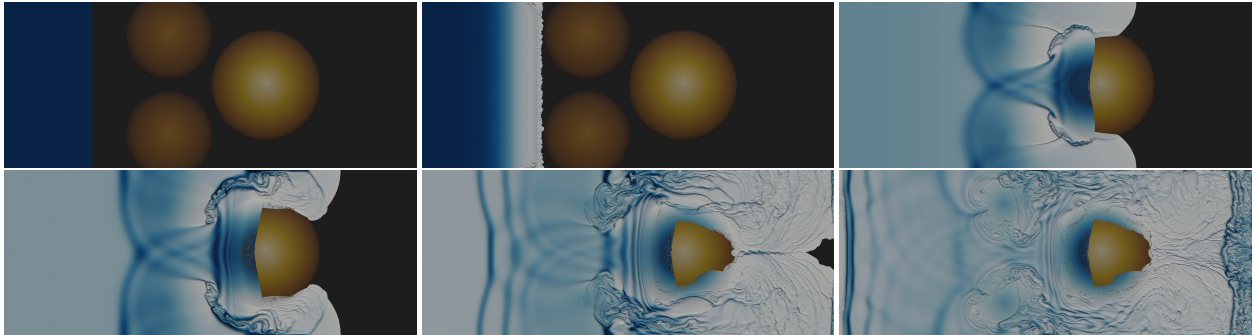


Figure 6.16: Dam break with bumps – Surface plot of the water elevation  $h + z$  at several time snapshots.

### 6.6.3 2D – Circular Dam Break with divergence-free velocity

We now reproduce a modified version of the circular dam break problem introduced in Section 6.5.3. Instead of initiating the dam at rest, we introduce a divergent-free velocity field so that circular dam is initially “spinning”. More precisely, we initialize the circular dam with amplitude 0.8 m centered at  $(x_c, y_c) = (250 \text{ m}, 250 \text{ m})$  with radius  $R = 40$  m and set the initial conditions to



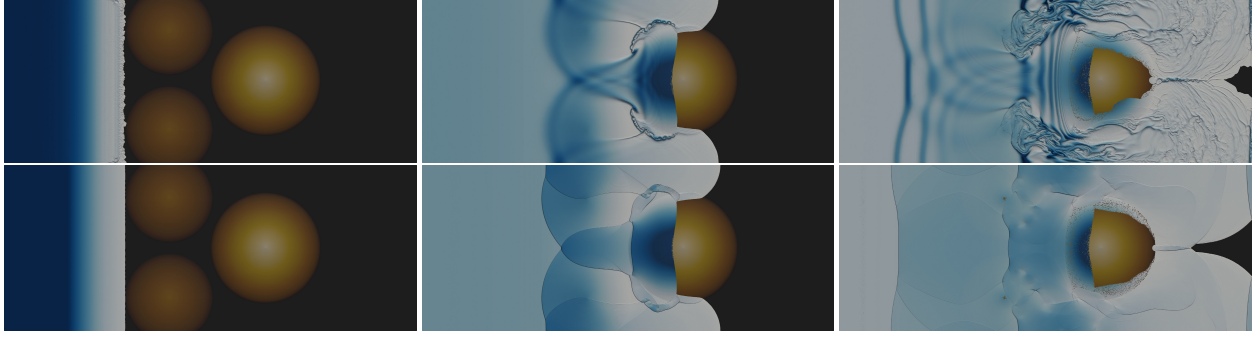


Figure 6.17: Dam break with bumps – Comparison with Hyperbolic Serre (top) and Saint-Venant (bottom) at  $t = \{1, 7.8, 15\}$ .

be:

$$h(\mathbf{x}, 0) = \begin{cases} 1.8 \text{ m}, & \sqrt{(x - 250)^2 + (y - 250)^2} < R, \\ 1 \text{ m}, & \text{otherwise,} \end{cases}, \quad \mathbf{q}(\mathbf{x}, 0) = h(\mathbf{x}, 0) \begin{pmatrix} y - y_c \\ -(x - x_x) \end{pmatrix}.$$

For this test case, we again use the *Ryujin* software. The computational domain is set to  $D = (0, 500 \text{ m}) \times (0, 500 \text{ m})$  with 4,004,001  $\mathbb{Q}_1$  nodes. The final time is set to  $T = 10 \text{ s}$  with CFL number of 0.075. In Figure 6.18, we compare the water depth profiles for both the hyperbolic Serre model (top) and the Saint-Venant model (bottom) at the times  $t = \{1, 5, 10\} \text{ s}$ . Notice that at time  $t = 1 \text{ s}$ , the Hyperbolic Serre model is producing structures reminiscent of a circular Kelvin–Helmholtz instability.

## 6.7 Laboratory experiments

We continue the numerical illustrations by reproducing several laboratory experiments that have been documented in the literature.

### 6.7.1 Shoaling of solitary waves over sloped beach

We now consider the 1994 experiments of Guibourg [35] conducted at LEGI (Laboratoire des Écoulements Géophysiques et Industriels) in Grenoble, France, to investigate the shoaling of solitary waves over a sloped beach. We consider 4 series of experiments proposed in [35] with a

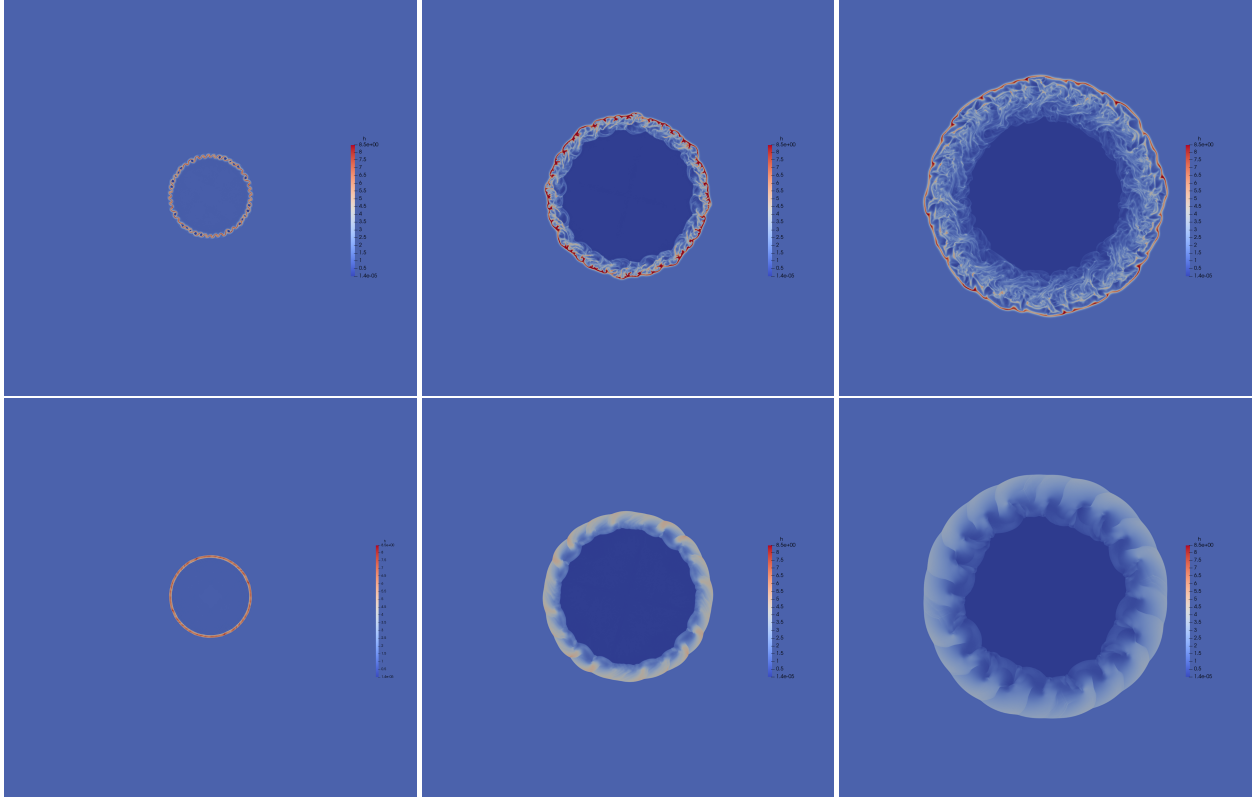


Figure 6.18: Circular dam break with divergent-free velocity field – Comparison with Hyperbolic Serre (top) and Saint-Venant (bottom) at  $t = \{1, 5, 10\}s$ .

reference water depth of  $h_0 = 0.25$  m and different solitary wave amplitudes (see: Table 6.5). We simulate the experiments in one spatial dimension and reproduce the bathymetry as follows:

$$z(x) = \begin{cases} \frac{1}{30}(x - 2.5) - h_0, & x \geq 25 \\ -h_0, & \text{otherwise.} \end{cases}$$

The computational domain is set to  $D = (-5 \text{ m}, 35 \text{ m})$ . For each experiment, we initialize the solitary wave at  $x_0 = 0$  m with the profiles defined in (6.3.1) and the amplitudes shown in Table 6.5. We run the computations to the final time  $T = 10$  s on the three difference meshes with respective mesh-size:  $h = \{0.05 \text{ m}, 0.025 \text{ m}, 0.0125 \text{ m}\}$  (corresponding to 800, 1600, 3200  $\mathbb{P}_1$  elements). The CFL number is set to 0.1. We set wall boundary conditions at both ends of the domain.

In the experiments, the wave elevation was measured with three wave gauges (WGs) which

	Case 1	Case 2	Case 3	Case 4
$\alpha/h_0$	0.096	0.2975	0.456	0.5343
WG1	7.75 m	5.75 m	4.25 m	4.25 m
WG2	8.25 m	6.25 m	5.0 m	5.0 m
WG3	8.75 m	6.75 m	5.75 m	5.75 m

Table 6.5: Solitary wave shoaling experiment [35] – configuration values

were moved for each case. We report the location of the wave gauges in Table 6.5. In Figure 6.19, we show the comparisons with the numerical computations and the experimental data for each case. We observe that the numerical results match the experimental data very well. This set of experiments reinforces that the dispersive hyperbolic Serre model captures well the shoaling phenomenon induced by topography.

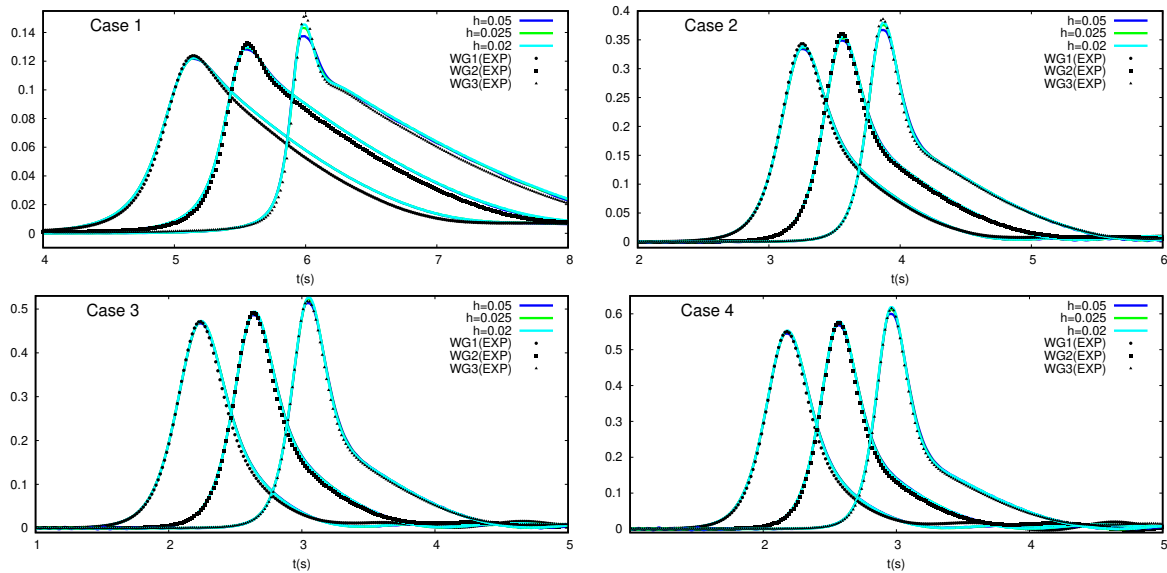


Figure 6.19: Comparison of numerical results with experimental data for solitary wave shoaling experiments of [35].

## 6.7.2 Periodic waves propagation over a submerged bar

We now consider the 1994 experiments conducted in Beji and Battjes [7] which investigate the propagation of periodic waves over a submerged trapezoidal bar. The goal of the experiments is to model the interaction of highly dispersive waves, and in particular, the release of higher-harmonics into a deeper region after the shoaling process.

We consider two of the experimental setups described in [7]: (i) sinusoidal long waves (SL) with target amplitude  $a = 1$  cm and period  $T_p = 2$  s; (ii) sinusoidal high-frequency waves (SH) with target amplitude  $a = 1$  cm and period  $T_p = 1.25$  s. We simulate these experiments in one spatial dimension and reproduce the bathymetry of the submerged bar as follows:

$$z(x) = \begin{cases} \frac{1}{20}(x - 6), & 6 \leq x \leq 12 \\ 0.3, & 12 \leq x \leq 14 \\ 0.3 - \frac{1}{10}(x - 14), & 14 \leq x \leq 17 \\ 0, & \text{otherwise.} \end{cases}$$

The computational domain is set to be  $D = (-12.3 \text{ m}, 37.7 \text{ m})$ . We impose two relaxation zones

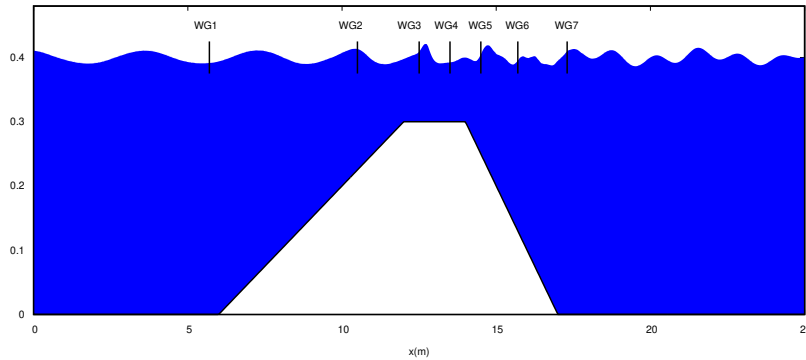


Figure 6.20: Submerged bar set up with gauge locations.

in the domain for the generation and absorption of waves (as described in Chapter 2). The length

of the generation zone for the SL case is 6 m (approximately 1.5 wavelengths) and 4 m for the SH case (approximately 2.0 wavelengths). The absorption zone for both cases is set to  $D_{\text{abs}} = (25 \text{ m}, 37.7 \text{ m})$ . We set the reference water depth to  $H_0 = 0.4 \text{ m}$  and initialize the water height profile with  $h_0(x) = H_0 - z(x)$  and discharge  $q_0(x) = 0$ . The periodic waves are introduced into the domain via the generation zone with the profiles given by:

$$h(x, t) = h_0 + a \sin(kx - \sigma t), \quad u(x, t) = \frac{a}{h_0} \frac{\sigma}{k} \sin(kx - \sigma t),$$

where  $a$  is the amplitude,  $k$  the wave number and  $\sigma$  the wave frequency. Here, we define the wave frequency by  $\sigma = \frac{2\pi}{T_p}$  and  $k$  is found by using the dispersion relation for the full Serre model:  $k^2 = 3\sigma^2 / (3gh_0 - h_0^2\sigma^2)$ . We set the final time to be  $t = 60 \text{ s}$  and run with CFL=0.175. We run the computations on three different meshes with meshsizes  $h = \{0.05 \text{ m}, 0.025 \text{ m}, 0.0125 \text{ m}\}$  (corresponding to 1000, 2000, and 4000  $\mathbb{P}_1$  cells.).

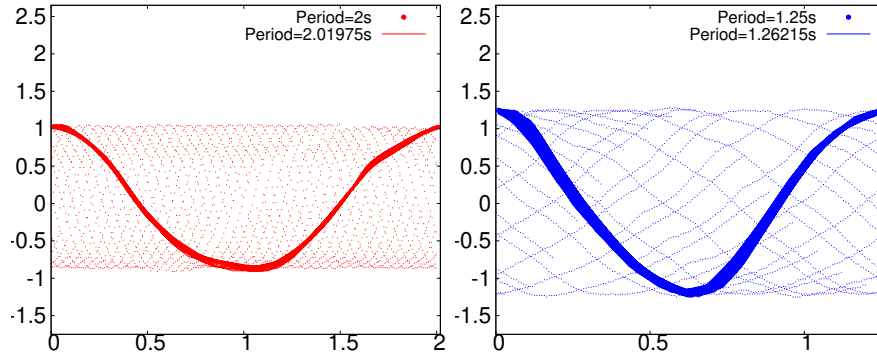


Figure 6.21: Illustration of *period-folding* with experimental wave gauge 1 with SL case (left) and SH case (right).

In the original experiments, seven wave gauges (WGs) were used to measure the water elevation: WG1( $x = 5.7 \text{ m}$ ), WG2( $x = 10.5 \text{ m}$ ), WG3( $x = 12.5 \text{ m}$ ), WG4( $x = 13.5 \text{ m}$ ), WG5( $x = 14.5 \text{ m}$ ), WG6( $x = 15.7 \text{ m}$ ), WG7( $x = 17.3 \text{ m}$ ). In Figure 6.20, we show the locations of these wave gauges with respect to the bathymetry. The experimental data used here was obtained from

the original author of the experiments, Serdar Beji. It was our experience that the experimental data did not quite match the targeted values of the period mentioned above. To illustrate this, we introduce a post-processing technique of the experimental data that we call *period-folding*. The idea is that given some experimental time series that is supposedly periodic with period  $T_p$  in the time interval  $[t_0, T_{\text{final}}]$ , the folding of the sequence obtained by the mapping  $t \mapsto t - t_0 - \lfloor \frac{t-t_0}{T_p} \rfloor T_p$  should represent the evolution of the signal during one period (here  $\lfloor \cdot \rfloor$  is the floor function). Doing this folding gives a better idea of the long time behavior of the experimental data than just looking at one specific window of length  $T_p$  as often done in the literature. In particular it reveals whether the signal is indeed periodic with period  $T_p$ . In Figure 6.21, we use period folding for the experimental data at WG1 with the targeted values of  $T_p$ . This process shows that the experimental data have not exactly the alleged period. We have been able to discover a good approximation of the actual period  $T_p^{\text{adj}}$  by doing the period folding with various values of  $T_p$ : (i) SL case,  $T_p = 2$  s,  $T_p^{\text{adj}} = 2.019,75$  s; (ii) SH case,  $T_p = 1.25$  s,  $T_p^{\text{adj}} = 1.262,15$  s. We also note that the wave amplitude value for the SH case is closer to 0.014 m and this is what we use for our computations.

Using the adjusted values above and the period-folding technique, we compare in Figures 6.22 and 6.23 the experimental data and the results of the computations at the wave gauges 2–7 (Figure 6.22 for the SL experiment and and Figure 6.23 for the SH experiment). For the experimental data, we choose  $t_0$  to be the time corresponding to the second maximum wave height of the signal at WG2. For the numerical simulations, we choose  $t_0$  to be the time corresponding to the maximum wave height around  $t \approx 40$  s at WG2. We observe that the numerical results converge as the mesh is refined. The SH experiment is relatively well reproduced at all the gauges. There are slight deviations at the last two gauges behind the bar for the SL experiment. It is possible that some wave breaking occurs between gauges 5 and 6 in this case.

### 6.7.3 Propagation of periodic waves over an elliptic shoal

We consider the 1982 experiments of Berkhoff et al. [8] conducted to study the propagation of monochromatic waves over an elliptic shoal. The goal of the experiments were to model the

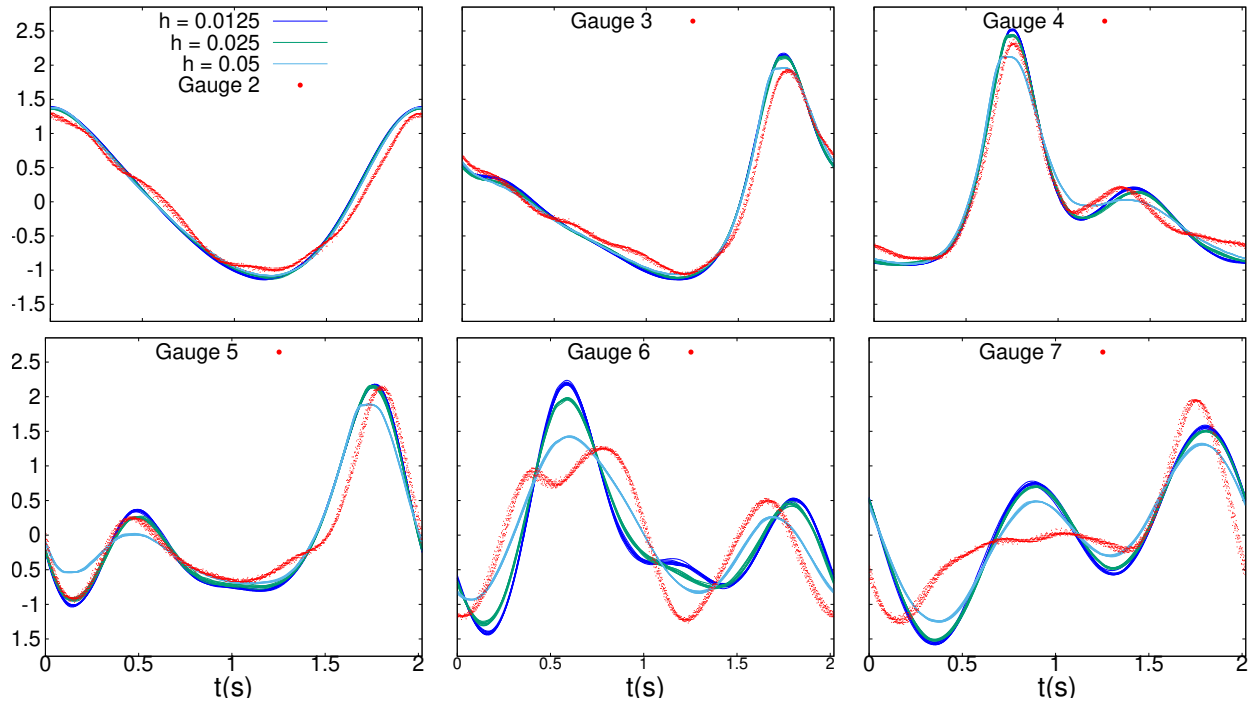


Figure 6.22: SL Case. Water elevation at seven gauges. Numerical results using three meshes,  $h = \{0.05 \text{ m}, 0.025 \text{ m}, 0.0125 \text{ m}\}$  (solid lines). Experimental data (red points).

refraction and diffraction of waves when propagating over a varying bottom. These experiment have become a benchmark for validating dispersive wave models (see: Duran and Marche [15], Ricchiuto and Filippini [52].)

The experimental basin is composed of a  $\frac{1}{50}$  sloping bottom which forms a  $20^\circ$  angle with the  $y$ -axis and an elliptic-shaped shoal built on the ramp. We reproduce this bathymetry as follows. We first define the rotated coordinates  $\boldsymbol{x} \mapsto \boldsymbol{x}_r(\boldsymbol{x})$ :

$$x_r := x \cos(20^\circ) - y \sin(20^\circ), \quad y_r := x \sin(20^\circ) + y \cos(20^\circ).$$

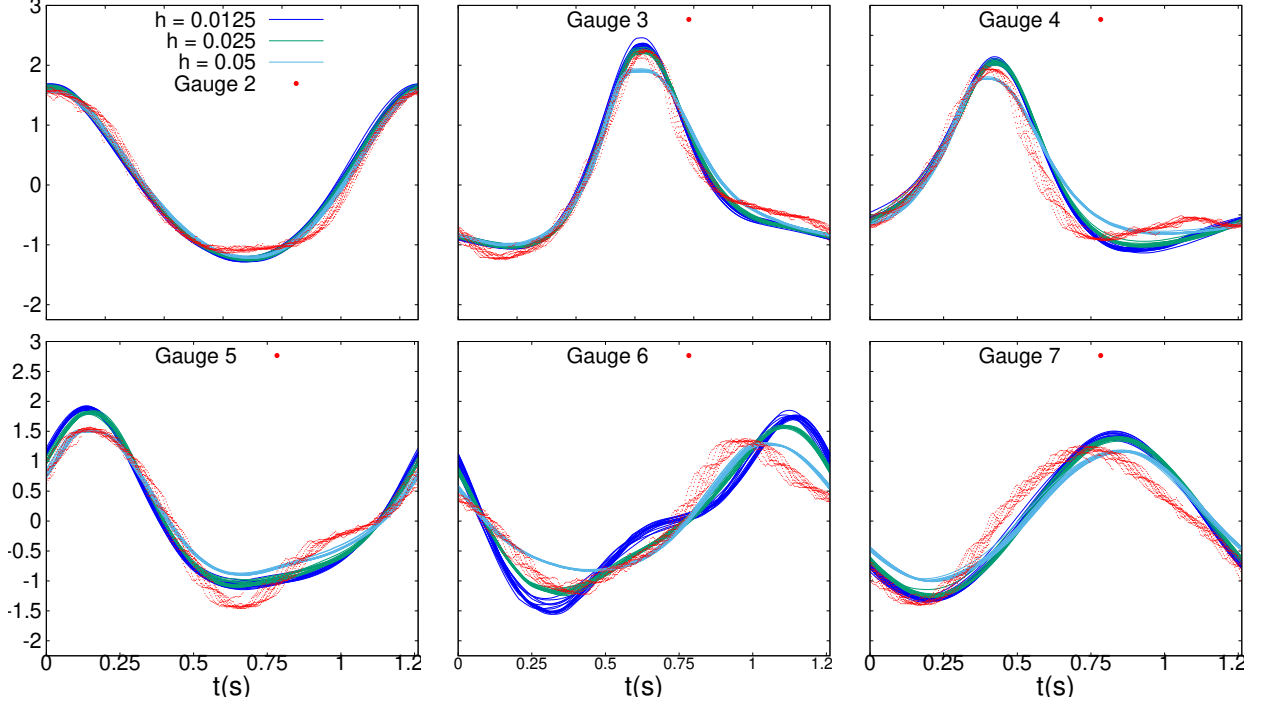


Figure 6.23: SH Case. Water elevation at seven gauges. Numerical results using three meshes,  $h = \{0.05 \text{ m}, 0.025 \text{ m}, 0.0125 \text{ m}\}$  (solid lines). experimental data (red points).

Then we define the sloping bottom and elliptic shoal profiles as:

$$z_{\text{ramp}}(\mathbf{x}) := \begin{cases} \frac{1}{50}(x_r(\mathbf{x}) + 5.82), & -5.82 \leq x_r(\mathbf{x}) \leq 14 \\ 0.3964, & 14 \leq x_r(\mathbf{x}) \\ 0, & \text{otherwise,} \end{cases}$$

$$z_{\text{shoal}}(\mathbf{x}) := \begin{cases} -0.3 + \frac{1}{2}\sqrt{1 - \left(\frac{x_r(\mathbf{x})}{3.75}\right)^2 - \left(\frac{y_r(\mathbf{x})}{5}\right)^2}, & \left(\frac{x_r(\mathbf{x})}{3.75}\right)^2 + \left(\frac{y_r(\mathbf{x})}{5}\right)^2 \leq 1 \\ 0, & \text{otherwise.} \end{cases}$$

The full bathymetry is defined as  $z(\mathbf{x}) := z_{\text{ramp}}(\mathbf{x}_r(\mathbf{x})) + z_{\text{shoal}}(\mathbf{x}_r(\mathbf{x}))$ . Note that this bathymetry is slightly modified from that proposed in Berkhoff et al. [8] to include a flat portion at the right-end of the basin.

The reference water depth is set to  $H_0 = 0.45 \text{ m}$ . We initialize the water height with  $h_0(\mathbf{x}) =$



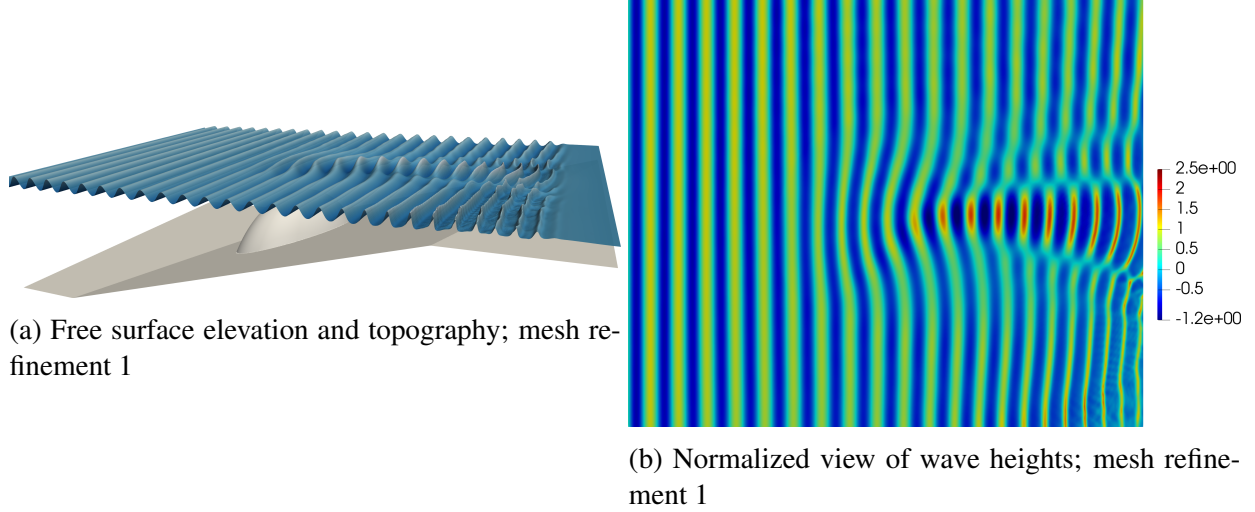


Figure 6.24: Elliptic shoal experiments

$H_0 - z(\mathbf{x})$  and discharge  $\mathbf{q}_0(\mathbf{x}) = \mathbf{0}$ . For the simulation of the experiments, we use the `Ryujin` code. The computational domain is set to be  $D = (-14, 18\text{m}) \times (-10, 10\text{m})$ . We generate the periodic waves via the generation zone methodology described in §2.4.3 with the profiles:

$$h(x, t) = h_0 + a \sin(kx - \sigma t), \quad u(x, t) = \frac{a}{h_0} \frac{\sigma}{k} \sin(kx - \sigma t),$$

The amplitude is set to  $a = 0.0232$  m and the period to  $T_p = 1$  s. The wave frequency is given by  $\sigma = \frac{2\pi}{T_p}$  and  $k$  is found by using the dispersion relation for the full Serre model:  $k^2 = 3\sigma^2 / (3gh_0 - h_0^2\sigma^2)$ . The generation zone length and absorption zone length are set to 4 m, i.e.,  $L_{\text{gen}} = L_{\text{abs}} := 4$  m, and we set  $x_{\text{min}} := -14$  and  $x_{\text{max}} := 18$ . We run the computations until the final time  $T = 60$  s to allow the waves to reach a steady state. To verify our results, we run the computations on three different meshes composed of 657,025, 2,624,769 and 10,492,417  $\mathbb{Q}_1$  nodes labeled refinement 1, 2, and 3 respectively. In Figure 6.24, we show a snapshot of the free surface elevation and topography at time  $T = 60$  s and a normalized view of the generated wave heights (i.e.,  $\frac{h+z(\mathbf{x})-H_0}{a}$ ).

In the experiments, the water elevation is measured at eight sections throughout the basin.

These sections are:

$$\begin{aligned}
\text{section 1} &: \{x = 1 \text{ m}, -5 \text{ m} \leq y \leq 5 \text{ m}\}, & \text{section 2} &: \{x = 3 \text{ m}, -5 \text{ m} \leq y \leq 5 \text{ m}\}, \\
\text{section 3} &: \{x = 5 \text{ m}, -5 \text{ m} \leq y \leq 5 \text{ m}\}, & \text{section 4} &: \{x = 7 \text{ m}, -5 \text{ m} \leq y \leq 5 \text{ m}\}, \\
\text{section 5} &: \{x = 9 \text{ m}, -5 \text{ m} \leq y \leq 5 \text{ m}\}, & \text{section 6} &: \{y = -2 \text{ m}, 0 \text{ m} \leq x \leq 11 \text{ m}\}, \\
\text{section 7} &: \{y = 0 \text{ m}, 0 \text{ m} \leq x \leq 11 \text{ m}\}, & \text{section 8} &: \{y = 2 \text{ m}, 0 \text{ m} \leq x \leq 11 \text{ m}\}.
\end{aligned}$$

To properly compare our numerical results with the experimental data, we do the following: we extract the data over the temporal window  $t \in [40 \text{ s}, T]$  of the reference water elevation  $h + z - H_0$ ; we then take the maximum of this data over every period in the temporal interval. We then normalize the wave heights with the incoming wave amplitude  $a = 0.0232 \text{ m}$ . In Figure 6.25, we show the comparison with the computational results for the three different meshes. We see that the approximate solutions converge and we observe that the computational results compare reasonably well with the experimental data.

#### 6.7.4 Propagation of periodic waves over semi-circular shoal

We now consider the 1971 experiments of Whalin [63] performed at the U.S. Army Engineer Waterways Experiment Station (now the U.S. Army Engineer Research and Development Center) in Vicksburg, Mississippi. The goal of the experiments is to study the refraction and diffraction of periodic waves propagating over a semi-circular shoal. In particular, we reproduce the experiments conducted where the wave period is  $T = 2 \text{ s}$  and amplitude  $a = 0.0075 \text{ m}$  (see [63, Fig. 68]).

The experimental basin is designed to be  $25.603 \text{ m}$  in length and  $6.096 \text{ m}$  wide and the still water elevation is set to  $0.4572 \text{ m}$ . We reproduce these experiments with the `RyuJin` code. We define the computational domain as  $(-10, 33 \text{ m}) \times (0, 6.096 \text{ m})$ . The lengths of the generation and relaxation zones are defined to be  $L_{\text{gen}} = L_{\text{abs}} = 8 \text{ m}$  (which is roughly 2 wave lengths) with  $x_{\text{min}} = -10$  and  $x_{\text{max}} = 33$ . The bathymetry is reproduced as follows: Defining  $G(y) := \sqrt{y(6.096 - y)}$ ,

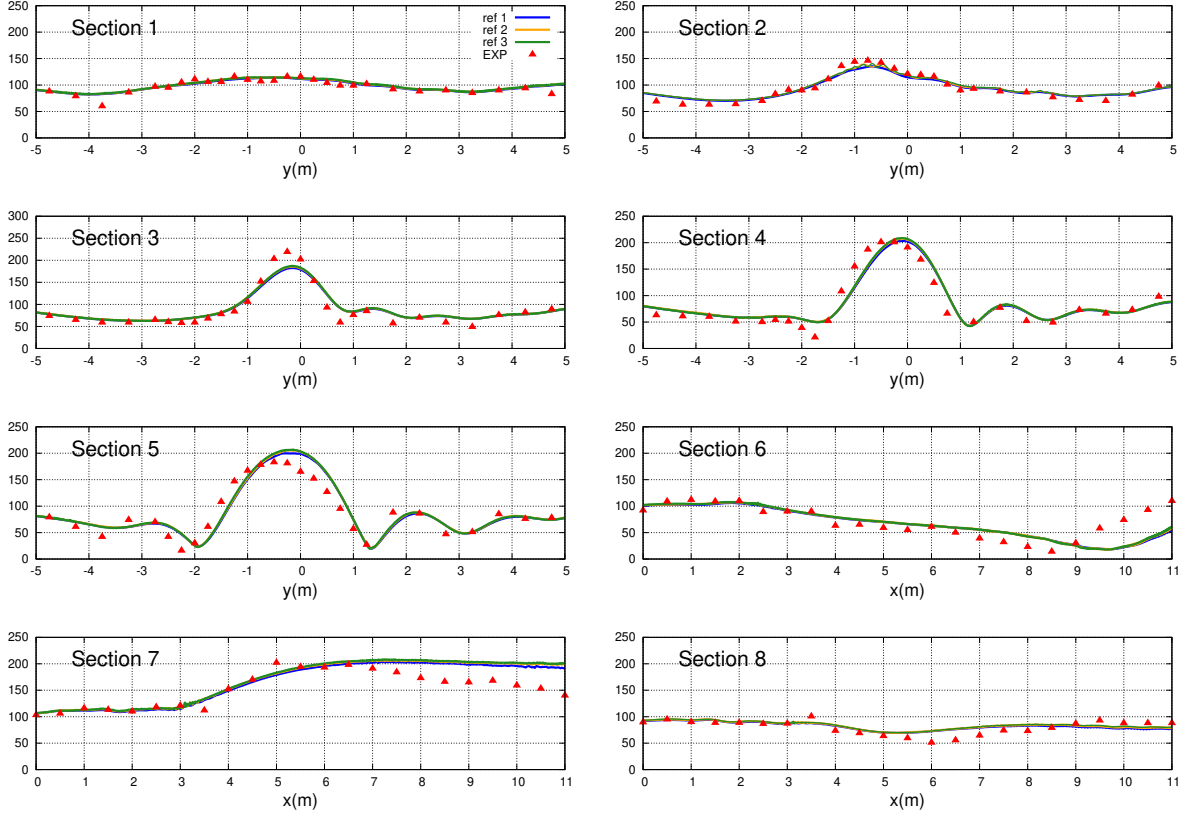


Figure 6.25: Elliptic shoal – Comparison of numerical results (3 mesh refinements) with the experimental data along the 8 sections. Experimental data: red triangles.

we set

$$z(\mathbf{x}) = \begin{cases} 0, & 0 \leq x \leq 10.67 - G(y), \\ -0.04(10.67 - G(y) - x), & 10.67 - G(y) \leq x \leq 18.297 - G(y), \\ 0.3048, & 18.297 - G(y) \leq x \end{cases}$$

The computational domain is composed of a quadrilateral mesh with 265,761  $\mathbb{Q}_1$  dofs. We run the numerical simulations until the time  $T = 60$  s to allow the waves to reach a steady state. The CFL number is set to 0.125.

In Whalin [63], the authors perform the harmonic analysis of the wave elevation data at the centerline of the basin over one period. This is done to study the non-linear transfer of energy from

lower to higher frequency components as the waves propagate and focus over the topography. We numerically reproduce this harmonic analysis as follows: We interpolate the centerline  $y = 3.048$  m with roughly 1400 points along the  $x$ -axis at every 0.001 s in the interval  $t \in [58, 60]$  s. We then perform the discrete Fourier Transform of the time-series wave elevation data at each point along centerline. In Figure 6.26, we show (a) the computational free surface elevation at  $T = 60$  s; (b) the comparison of the amplitude spectrum with the numerical first, second and third harmonics (solid lines) and the experimental data of Whalin (black geometric shapes). The amplitude spectrum of the waves in the numerical simulations is very close to the experimental one.

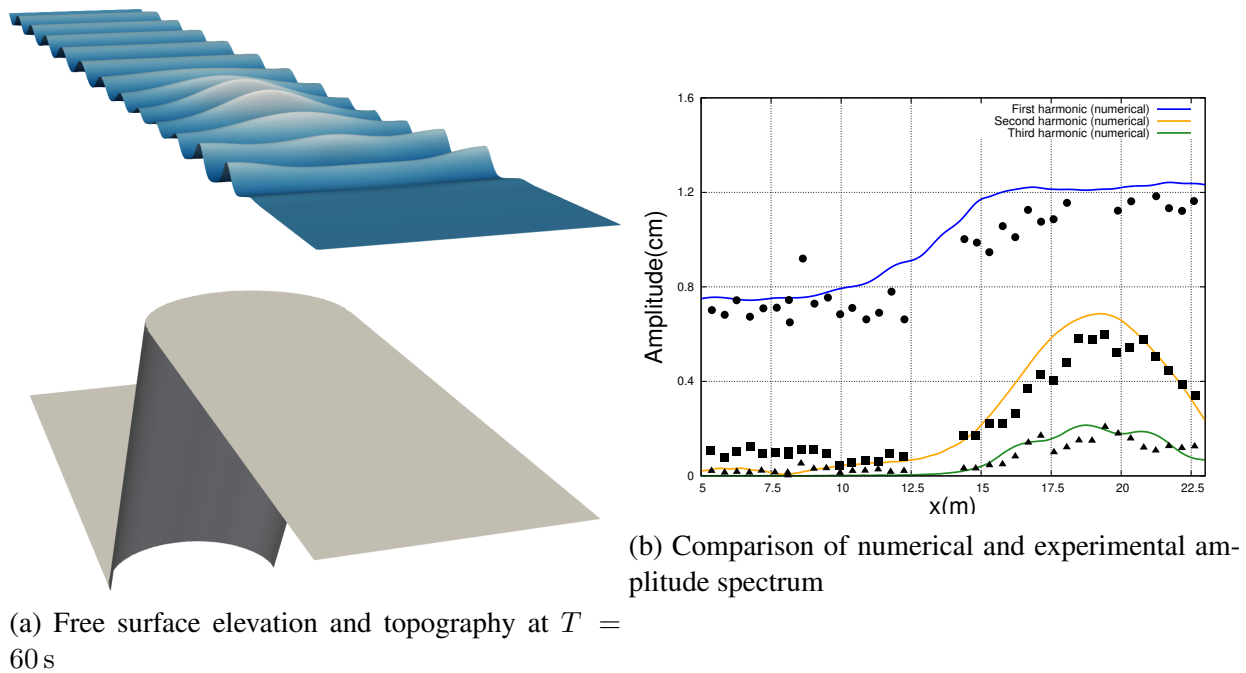


Figure 6.26: Whalin semi-circular shoal results

### 6.7.5 2D Solitary wave run-up over a conical island

We consider the 1995 laboratory experiments conducted by Briggs et al. [11] at the US Army Waterways Experiment Station in Vicksburg, Mississippi (now the US Army Engineer Research and Development Center). The laboratory experiments were motivated by several tsunami events

in the 1990s where large unexpected run-up heights were observed on the back (or lee) side of small islands. Several authors have used this experiment to study the run-up phenomena using the classical Shallow Water model and other dispersive models (see: Hou et al. [36], Lannes and Marche [45], Kazolea et al. [39]).

Let  $r(\boldsymbol{x})$  by the radius from the center of the island located at (12.96 m, 13.80 m). Then the conical island bathymetry is defined by

$$z(\boldsymbol{x}) = \begin{cases} \min(h_{\text{top}}, h_{\text{cone}} - r(\boldsymbol{x})/s_{\text{cone}}), & r(\boldsymbol{x}) < r_{\text{cone}}, \\ 0, & \text{otherwise,} \end{cases} \quad (6.7.1)$$

where  $h_{\text{top}} = 0.625\text{m}$ ,  $h_{\text{cone}} = 0.9\text{m}$  and  $r_{\text{cone}} = 3.6\text{m}$ . We reproduce two experiments, which we call Case B and Case C, with  $\alpha/h_0 = 0.091$  and  $\alpha/h_0 = 0.181$  where  $h_0 = 0.32\text{m}$  and  $\alpha$  is the amplitude of the solitary wave.

The computations are done in the domain  $(0, 25\text{ m}) \times (0, 30\text{ m})$  until the final time  $T = 12\text{ s}$  with wall boundary conditions and CFL number 0.25. We initiate the solitary wave at  $x_0 = 9.36 - \frac{L}{2}$  using (6.3.1) with  $L = \frac{2h_0}{k} \operatorname{arccosh} \sqrt{20}$  and  $k = \sqrt{\frac{3\alpha}{4h_0}}$ . Here  $x_0$  is the location of the experimental wave gauge 3 (WG3) which was used to measure the free surface elevation away from the island. In Figure 6.27, we show the surface plots of the free surface elevation  $h + z$  on a mesh composed of 52,129  $\mathbb{P}_1$  nodes at  $t = \{0, 5.8, 8\text{ s}\}$ .

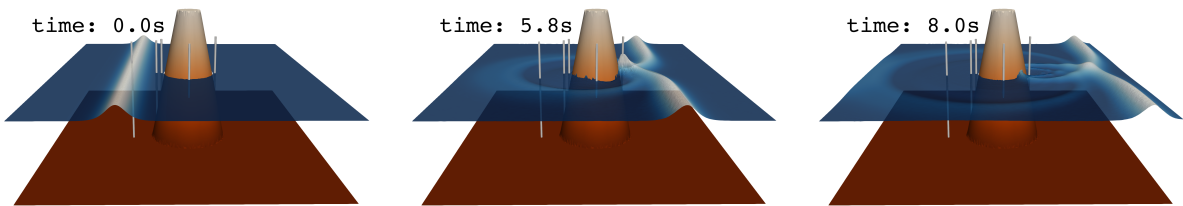


Figure 6.27: Experiment 4 – Surface plot of the water elevation  $h + z$  at several times for Case C. The thin grey cylinders represent the wave gauges WG3, WG6, WG9, WG16, WG22 (left to right).

In the experiment, several wave gauges were placed around the island to measure the free surface elevation and wave run-up. We compare the numerical results with the measurements at four of the experimental wave gauges: WG6(9.36 m, 13.80 m), WG9(10.36 m, 13.80 m), WG16(12.96 m, 11.22 m), WG22(15.56 m, 13.80 m). We show in Figure 6.28 the comparison with the experimental data and numerical simulations for both Case B (on the left) and Case C (on the right). For both cases, the numerical results show good agreement with the experimental data. We capture well the magnitudes of the run-up and draw-down at the front side of the island at WG9 with a slight overshoot in Case C. For both cases, we see very good comparison with WG16 which corresponds to the run-up and draw-down on the side of the island. We note that the experimental data shows subsequent free surface oscillations after impacting the island which is most notable in WG9, but our numerical simulation do not capture this effect. This phenomena is consistent with the literature and has been observed by others (see: Lannes and Marche [45], Kazolea et al. [39], Yamazaki et al. [66]), and is likely due to inconsistency in the original experiments.

### 6.7.6 Propagation over a solitary wave over a triangular shelf with conical island

We now reproduce the experiments of Swigler [59] and Lynett et al. [47] performed at the O.H. Hinsdale Wave Research Laboratory of Oregon State University. The experiments were conducted to study specific phenomena that are known to occur when solitary waves propagate over irregular bathymetry such as shoaling, refraction, breaking, etc. Several others (see: Duran and Marche [15], Kazolea et al. [40], Roeber and Cheung [53]) have used these experiments for validation.

We reproduce the bathymetry of the experiments as follows: Let  $r = 3$ ,  $h_{\text{cone}} = 0.45$ ,  $d(y) := 1 - \min(1, |y|/13.25)$ ,  $a_x(y) := 12.5 + 12.4999(1 - d(y))$ ,  $a_z(y) := 0.7 + 0.05(1 - d(y))$ . We define separately the cone, base and triangular shelf portions of the bathymetry:

$$\text{cone}(x, y) := \max \left( h_{\text{cone}} - \sqrt{\frac{(x - 17)^2 + y^2}{(\frac{3}{0.45})^2}}, 0 \right),$$

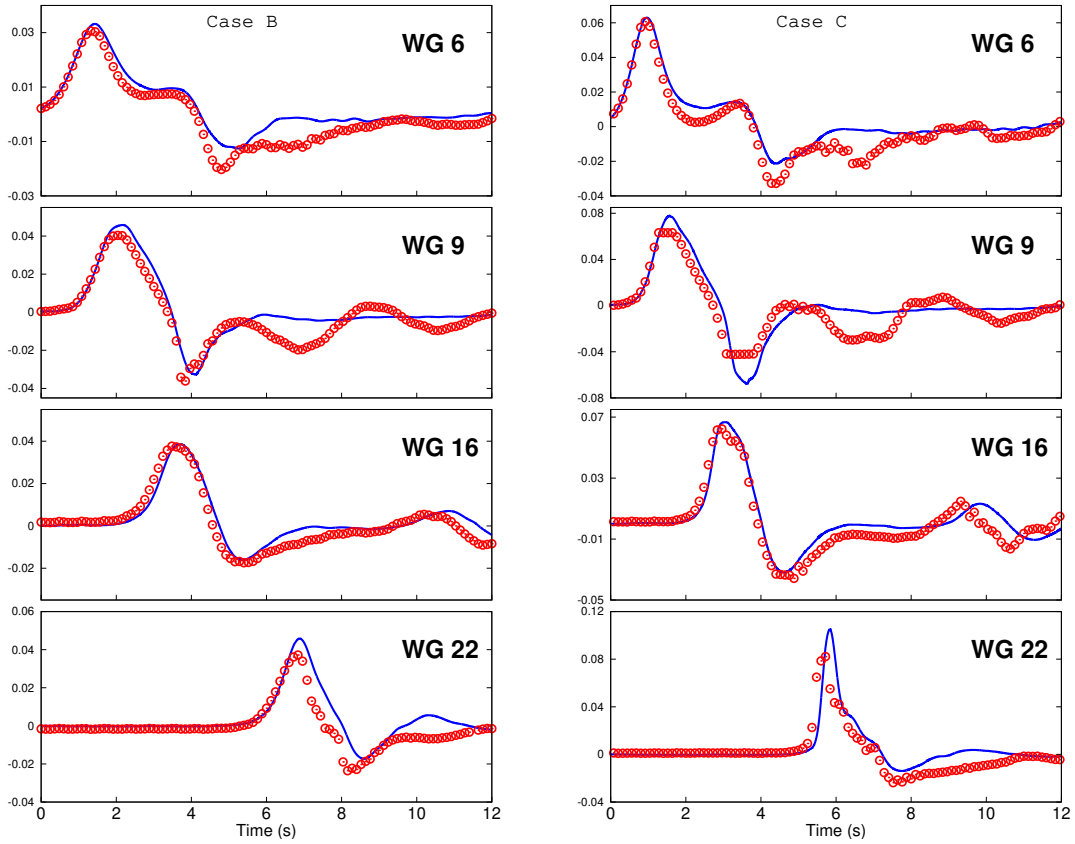


Figure 6.28: Experiment 4 – Temporal series over the period  $t \in [0, 12 \text{ s}]$  of the free surface elevation  $h + z$  in meters at the four WGs (blue solid) compared to the experimental data (red circles) for Case B (on the left) and Case C (on the right).

$$\text{base}(x, y) := \begin{cases} 0, & x < 10.2, \\ \frac{0.5-0.0}{17.5-10.2}(x - 10.2), & 10.2 \leq x \leq 17.5, \\ 1 + \frac{1-0.5}{32.5-17.5}(x - 32.5), & 17.5 \leq x \leq 32.5, \\ 1, & \text{otherwise} \end{cases},$$

$$\text{shelf}(x, y) := \begin{cases} 0, & x < 10.2, \\ \frac{a_z(y)}{a_x(y)-10.2}(x - 10.2), & 10.2 \leq x \leq a_x(y), \\ 0.75 + \frac{a_z(y)-0.75}{a_x(y)-25}(x - 25), & a_x(y) \leq x \leq 25, \\ 1 + \frac{1-0.5}{32.5-17.5}(x - 32.5), & 25 \leq x \leq 32.5, \\ 1, & \text{otherwise} \end{cases},$$

Then the full bathymetry is defined by

$$z(x, y) := \text{cone}(x, y) + \max(\text{base}(x, y), \text{shelf}(x, y)).$$

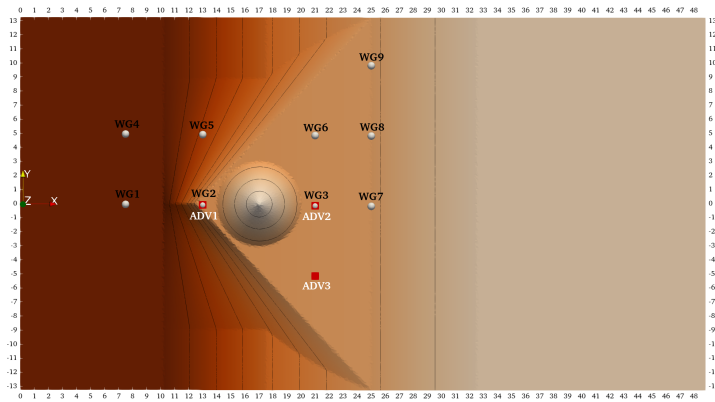
The setup of this complex bathymetry can be seen in Figure 6.29b.

The computations are done in the domain  $(0, 48.8 \text{ m}) \times (-13.25, 13.25 \text{ m})$ . The solitary wave is initiated at  $x_0 = 5 \text{ m}$  with reference water depth  $h_0 = 0.78 \text{ m}$  and amplitude  $\alpha = 0.39 \text{ m}$  using (6.3.1). We run the computations until  $T = 40 \text{ s}$  with a CFL number of 0.25. We note that for this particular problem, it is our experience that no friction is needed to reproduce correctly the experiment. In Figure 6.31, we show the surface plots of the free surface elevation  $h + z$  on a mesh composed of 57,854  $\mathbb{P}_1$  nodes at various times using the TAMU code.

In the experiments, nine wave gauges (WGs) are placed along the basin to capture the free surface elevation along with three Acoustic Doppler Velocimeters (ADV) that measure the velocity. In Figure 6.29, we show on the left panel the coordinates of the wave gauges and ADVs, and their respective locations on the bathymetry in the right panel of the figure. We show in Figure 6.30a, the comparison between the free surface elevation values of the numerical simulation and the experimental data over the temporal period  $t \in [0, 40 \text{ s}]$  using both codes. In Figure 6.30b, we show the comparison between the numerical velocities and the experimental data from the ADVs. For both the free surface and velocities, our results compare exceptionally well with the experimental data. We also see that the results of the TAMU and Proteus codes agree very closely and are almost indistinguishable. The Proteus computations were done on a mesh composed of 57,188



Gauge	x(m)	y(m)
WG1	7.5	0.0
WG2	13.0	0.0
WG3	21.0	0.0
WG4	7.5	5.0
WG5	13.0	5.0
WG6	21.0	5.0
WG7	25.0	0.0
WG8	25.0	5.0
WG9	25.0	10.0
ADV1	13.0	0.0
ADV2	21.0	0.0
ADV3	21.0	-5.0

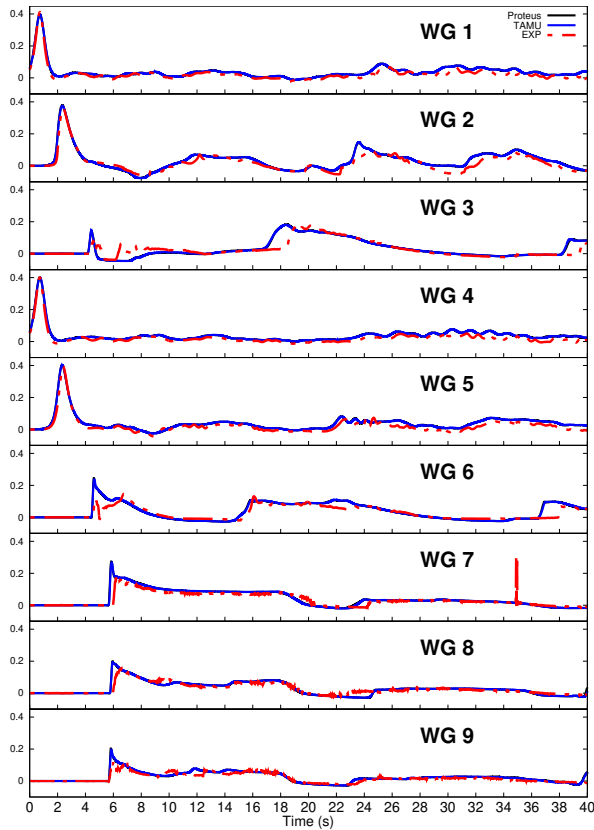


(b) Overview of bathymetry with WGs (white spheres) and ADVs locations (red boxes).

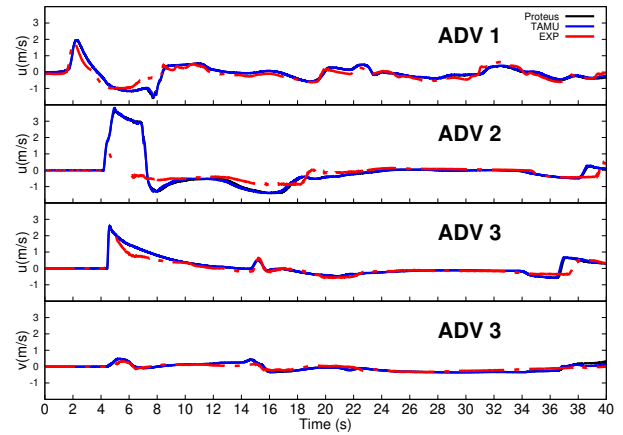
(a) WGs and ADVs coordinates.

Figure 6.29: (a) Coordinates of the wave gauges and ADVs in meters; (b) Overview of their respective locations on the bathymetry.

$\mathbb{P}_1$  nodes and CFL number of 0.25.



(a) Wave gauges



(b) Acoustic Doppler Velocimeters

Figure 6.30: (a) Temporal series over the period  $t \in [0, 40s]$  of the free surface elevation  $h + z$  compared to the experimental data (red dashed). The TAMU code results are in blue (solid) and Proteus code results in black (solid). (b) Temporal series over the period  $t \in [0, 40s]$  of velocity  $v$  (blue solid) and experimental ADVs (red dashed).

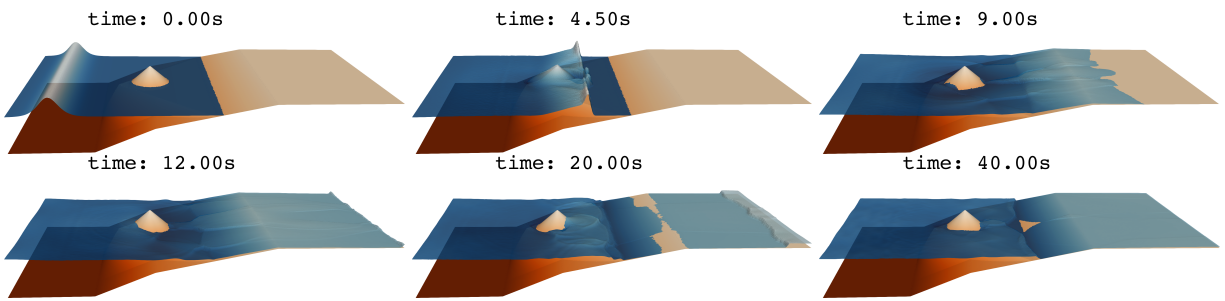


Figure 6.31: Surface plot of the water elevation  $h + z$  at several times.

## REFERENCES

- [1] Borys Alvarez-Samaniego and David Lannes. Large time existence for 3D water-waves and asymptotics. *Invent. Math.*, 171(3):485–541, 2008.
- [2] George J. Arcement and Verne R. Schneider. Guide for selecting manning’s roughness coefficients for natural channels and flood plains. Technical Report 2339, U.S. Geological Survey, 1989.
- [3] Daniel Arndt, Wolfgang Bangerth, Bruno Blais, Marc Fehling, Rene Gassmüller, Timo Heister, Luca Heltai, Uwe Köcher, Martin Kronbichler, Matthias Maier, Peter Munch, Jean-Paul Pelteret, Sebastian Proell, Konrad Simon, Bruno Turcksin, David Wells, and Jiaqi Zhang. The deal.II library, version 9.3. *Journal of Numerical Mathematics*, 2021, accepted for publication. URL <https://dealii.org/deal93-preprint.pdf>.
- [4] Emmanuel Audusse, François Bouchut, Marie-Odile Bristeau, Rupert Klein, and B. Perthame. A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows. *SIAM J. Sci. Comput.*, 25(6):2050–2065, 2004.
- [5] Pascal Azerad, Jean-Luc Guermond, and Bojan Popov. Well-balanced second-order approximation of the shallow water equation with continuous finite elements. *SIAM J. Numer. Anal.*, 55(6):3203–3224, 2017.
- [6] Eric Barthélemy. Nonlinear shallow water theories for coastal waves. *Surveys in Geophysics*, 25:315–337, 2004.
- [7] S. Beji and J.A. Battjes. Numerical simulation of nonlinear wave propagation over a bar. *Coastal Engineering*, 23(1):1 – 16, 1994.
- [8] J.C.W. Berkhoff, N. Booy, and A.C. Radder. Verification of numerical wave propagation models for simple harmonic linear water waves. *Coastal Engineering*, 6(3):255 – 279, 1982.
- [9] Alfredo Bermúdez and Ma. Elena Vázquez. Upwind methods for hyperbolic conservation laws with source terms. *Comput. Fluids*, 23(8):1049–1071, 1994.
- [10] P. Bonneton, F. Chazel, D. Lannes, F. Marche, and M. Tissier. A splitting approach for

- the fully nonlinear and weakly dispersive Green–Naghdi model. *J. Comput. Phys.*, 230(4): 1479–1498, 2011.
- [11] Michael J. Briggs, Costas E. Synolakis, Gordon S. Harkins, and Debra R. Green. Laboratory experiments of tsunami runup on a circular island. *pure and applied geophysics*, 144(3): 569–593, 1995.
- [12] Marie-Odile Bristeau, Anne Mangeney, Jacques Sainte-Marie, and Nicolas Seguin. An energy-consistent depth-averaged Euler system: derivation and properties. *Discrete Contin. Dyn. Syst. Ser. B*, 20(4):961–988, 2015.
- [13] A. Chertock, S. Cui, A. Kurganov, and T. Wu. Well-balanced positivity preserving central-upwind scheme for the shallow water system with friction terms. *Internat. J. Numer. Methods Fluids*, 78(6):355–383, 2015.
- [14] Paul J. Dellar and Rick Salmon. Shallow water equations with a complete coriolis force and topography. *Physics of Fluids*, 17(10):106601, 2005. doi: 10.1063/1.2116747.
- [15] A. Duran and F. Marche. A discontinuous Galerkin method for a new class of Green–Naghdi equations on simplicial unstructured meshes. *Appl. Math. Model.*, 45:840–864, 2017.
- [16] Denys Dutykh, Didier Clamond, Paul Milewski, and Dimitrios Mitsotakis. Finite volume and pseudo-spectral schemes for the fully nonlinear 1d serre equations. *European Journal of Applied Mathematics*, 24(5):761–787, 2013. doi: 10.1017/S0956792513000168.
- [17] G. A. El, R. H. J. Grimshaw, and N. F. Smyth. Unsteady undular bores in fully nonlinear shallow-water theory. *Physics of Fluids*, 18(2):027104, 2006.
- [18] Alexandre Ern and Jean-Luc Guermond. *Finite Elements III*. Springer International Publishing, 2021.
- [19] C. Escalante, M. Dumbser, and M. J. Castro. An efficient hyperbolic relaxation system for dispersive non-hydrostatic water waves and its solution with high order discontinuous Galerkin schemes. *J. Comput. Phys.*, 394:385–416, 2019.
- [20] N. Favrie and S. Gavriluk. A rapid numerical method for solving Serre–Green–Naghdi equations describing long free surface gravity waves. *Nonlinearity*, 30(7):2718–2736, 2017.

- [21] Enrique D. Fernandez-Nieto, Martin Parisot, Yohan Penel, and Jacques Sainte-Marie. A hierarchy of dispersive layer-averaged approximations of euler equations for free surface flows. *Communications in Mathematical Sciences*, 16(5):1169–1202, 2019.
- [22] S. L. Gavriluyk and S. M. Shugrin. Media with equations of state that depend on derivatives. *Journal of Applied Mechanics and Technical Physics*, 37(2):177–189, 1996.
- [23] Sergey L. Gavriluyk, Boniface Nkonga, Keh-Ming Shyue, and Lev Truskinovsky. Generalized riemann problem for dispersive equations. 2018.
- [24] Edwige Godlewski and Pierre-Arnaud Raviart. *Introduction*. Springer New York, New York, NvY, 1996. ISBN 978-1-4612-0713-9.
- [25] A. E. Green and P. M. Naghdi. A derivation of equations for wave propagation in water of variable depth. *Journal of Fluid Mechanics*, 78(2):237–246, 1976. doi: 10.1017/S0022112076002425.
- [26] A. E. Green, N. Laws, and P. M. Naghdi. On the theory of water waves. *Proc. Roy. Soc. (London) Ser. A*, 338:43–55, 1974.
- [27] J. M. Greenberg and A.-Y. Le Roux. A well-balanced scheme for the numerical processing of source terms in hyperbolic equations. *SIAM J. Numer. Anal.*, 33(1):1–16, 1996.
- [28] J.-L. Guermond, M. Quezada de Luna, B. Popov, C. Kees, and M. Farthing. Well-balanced second-order finite element approximation of the shallow water equations with friction. *SIAM Journal on Scientific Computing*, 40(6):A3873–A3901, 2018.
- [29] Jean-Luc Guermond and Richard Pasquetti. A correction technique for the dispersive effects of mass lumping for transport problems. *Computer Methods in Applied Mechanics and Engineering*, 253:186 – 198, 2013.
- [30] Jean-Luc Guermond, Bojan Popov, and Yong Yang. The effect of the consistent mass matrix on the maximum-principle for scalar conservation equations. *Journal of Scientific Computing*, 70(3):1358–1366, 2017.
- [31] Jean-Luc Guermond, Murtazo Nazarov, Bojan Popov, and Ignacio Tomas. Second-order invariant domain preserving approximation of the euler equations using convex limiting. *SIAM*

- J. Scientific Computing*, 40(5):A3211–A3239, 2018.
- [32] Jean-Luc Guermond, Bojan Popov, and Ignacio Tomas. Invariant domain preserving discretization-independent schemes and convex limiting for hyperbolic systems. *Computer Methods in Applied Mechanics and Engineering*, 347:143 – 175, 2019. ISSN 0045-7825.
- [33] Jean-Luc Guermond, Bojan Popov, Eric Tovar, and Chris Kees. Robust explicit relaxation technique for solving the Green–Naghdi equations. *J. Comput. Phys.*, 399:108917, 17, 2019.
- [34] Jean-Luc Guermond, Matthias Maier, Bojan Popov, and Ignacio Tomas. Second-order invariant domain preserving approximation of the compressible navier–stokes equations. *Computer Methods in Applied Mechanics and Engineering*, 375(1):113608, 2021.
- [35] Sandrine Guibourg. *Modélisations numérique et expérimentale des houles bidimensionnelles en zone cotière*. PhD thesis, 1994. URL <http://www.theses.fr/1994GRE10160>.
- [36] Jingming Hou, Qiuhua Liang, Franz Simons, and Reinhard Hinkelmann. A stable 2d unstructured shallow flow model for simulations of wetting and drying over rough terrains. *Computers & Fluids*, 82:132 – 147, 2013.
- [37] Yuxin Huang, Ningchuan Zhang, and Yuguo Pei. Well-balanced finite volume scheme for shallow water flooding and drying over arbitrary topography. *Engineering Applications of Computational Fluid Mechanics*, 7(1):40–54, 2013.
- [38] Mutsuto Kawahara and Tsuyoshi Umetsu. Finite element method for moving boundary problems in river flow. *International Journal for Numerical Methods in Fluids*, 6(6):365–386, 1986.
- [39] M. Kazolea, A.I. Delis, I.K. Nikolos, and C.E. Synolakis. An unstructured finite volume numerical scheme for extended 2d boussinesq-type equations. *Coastal Engineering*, 69:42 – 66, 2012.
- [40] M. Kazolea, A.I. Delis, and C.E. Synolakis. Numerical treatment of wave breaking on unstructured finite volume approximations for extended boussinesq-type equations. *Journal of Computational Physics*, 271:281 – 305, 2014. *Frontiers in Computational Physics*.
- [41] C. E. Kees and M. W. Farthing. Parallel computational methods and simulation for coastal

- and hydraulic applications using the proteus toolkit. *Supercomputing 11: Proceedings of the PyHPC11 Workshop.*, 2011.
- [42] B. Khobalatte and B. Perthame. Maximum principle on the entropy and second-order kinetic schemes. *Math. Comput.*, 62(205):119–131, 1994.
- [43] Alexander Kurganov and Guergana Petrova. A second-order well-balanced positivity preserving central-upwind scheme for the Saint-Venant system. *Commun. Math. Sci.*, 5(1):133–160, 2007.
- [44] D Lannes. Modeling shallow water waves. *Nonlinearity*, 33(5):R1–R57, 2020.
- [45] D. Lannes and F. Marche. A new class of fully nonlinear and weakly dispersive Green–Naghdi models for efficient 2d simulations. *Journal of Computational Physics*, 282:238 – 268, 2015.
- [46] David Lannes and Philippe Bonneton. Derivation of asymptotic two-dimensional time-dependent equations for surface water wave propagation. *Physics of Fluids*, 21(1):016601, 2009. doi: 10.1063/1.3053183.
- [47] Patrick Lynett, David Swigler, Hoda El Safty, Luis Motoya, Adam Keen, Sangoung Son, and Pablo Higuera. Study of the three-dimensional hydrodynamics associated with a solitary wave traveling over an alongshore-variable, shallow shelf. *Journal of Waterway, Port, Coastal, and Ocean Engineering (ASCE)*, 2019.
- [48] P. A. Madsen, H. B. Bingham, and H. A. Schäffer. Boussinesq-type formulations for fully nonlinear and extremely dispersive water waves: derivation and analysis. *Proc. Roy. Soc. (London) Ser. A*, 459:1075–1104, 2003.
- [49] Matthias Maier and Martin Kronbichler. Efficient parallel 3d computation of the compressible euler equations with an invariant-domain preserving second-order finite-element scheme. *ACM Transactions on Parallel Computing*, accepted, 2021. URL <https://arxiv.org/abs/2007.00094>.
- [50] Fabien Marche. Combined hybridizable discontinuous galerkin (hdg) and runge-kutta discontinuous galerkin (rk-dg) formulations for green-naghdi equations on unstructured meshes.



- Journal of Computational Physics*, 418:109637, 2020. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2020.109637>.
- [51] Dimitrios Mitsotakis, Denys Dutykh, and John D. Carter. On the nonlinear dynamics of the traveling-wave solutions of the Serre system. *Wave Motion*, 70:166–182, April 2017. doi: 10.1016/j.wavemoti.2016.09.008.
- [52] M. Ricchiuto and A. G. Filippini. Upwind residual discretization of enhanced Boussinesq equations for wave propagation over complex bathymetries. *J. Comput. Phys.*, 271:306–341, 2014.
- [53] Volker Roeber and Kwok Fai Cheung. Boussinesq-type model for energetic breaking waves in fringing reef environments. *Coastal Engineering*, 70:1 – 20, 2012. ISSN 0378-3839.
- [54] A.J.C. 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. *C. R. Acad. Sc. Paris*, (73):147–154, 1871.
- [55] Ali Samii and Clint Dawson. An explicit hybridized discontinuous galerkin method for Serre–Green–Naghdi wave model. *Computer Methods in Applied Mechanics and Engineering*, 330: 447 – 470, 2018.
- [56] Fernando J. Seabra-Santos, Dominique P. Renouard, and André M. Temperville. Numerical and experimental study of the transformation of a solitary wave over a shelf or isolated obstacle. *Journal of Fluid Mechanics*, 176:117–134, 1987.
- [57] François Serre. Contribution à l’étude des écoulements permanents et variables dans les canaux. *La Houille Blanche*, (6):830–872, 1953.
- [58] C. H. Su and C. S. Gardner. Korteweg-de Vries equation and generalizations. III. Derivation of the Korteweg-de Vries equation and Burgers equation. *J. Mathematical Phys.*, 10:536–539, 1969.
- [59] David Townley Swigler. Laboratory Study Investigating the Three-dimensional Turbulence and Kinematic Properties Associated with a Breaking Solitary Wave. Master’s thesis, Texas A&M University., College Station, Texas, 2009.

- [60] Sergey Tkachenko. *Analytical and numerical study of a dispersive shallow water model*. PhD thesis, Aix-Marseille University, 2020.
- [61] E.F. Toro. *Shock-capturing methods for free-surface shallow flows*. John Wiley, 2001.
- [62] C. B. Vreugdenhil. *Shallow-water flows*. Springer Netherlands, 1994. ISBN 978-94-015-8354-1.
- [63] R. Whalin. *The limit of applicability of linear wave refraction theory in a convergence zone*. PhD thesis, Texas A&M University, 1971.
- [64] Gerald B. Whitham. *Linear Dispersive Waves*. John Wiley & Sons, Ltd, 1999. ISBN 9781118032954.
- [65] Inc. Wolfram Research. Mathematica, Version 12.3.1. URL <https://www.wolfram.com/mathematica>. Champaign, IL, 2021.
- [66] Yoshiki Yamazaki, Zygmunt Kowalik, and Kwok Fai Cheung. Depth-integrated, non-hydrostatic model for wave breaking and run-up. *International Journal for Numerical Methods in Fluids*, 61(5):473–497, 2009.
- [67] Steven T Zalesak. Fully multidimensional flux-corrected transport algorithms for fluids. *Journal of Computational Physics*, 31(3):335–362, 1979. ISSN 0021-9991.
- [68] Yao Zhang, Andrew B. Kennedy, Nishant Panda, Clint Dawson, and Joannes J. Westerink. Generating–absorbing sponge layers for phase-resolving wave models. *Coastal Engineering*, 84:1–9, 2014. ISSN 0378-3839.