

A ROBUST SECOND ORDER INVARIANT-DOMAIN PRESERVING
APPROXIMATION OF THE COMPRESSIBLE EULER EQUATIONS WITH AN
ARBITRARY EQUATION OF STATE

A Thesis

by

BENNETT GILES CLAYTON

Submitted to the Office of 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, Bojan Popov

Committee Members, Jean-Luc Guermond

Matthias Maier

Jean Ragusa

Head of Department, Sarah Witherspoon

May 2023

Major Subject: Mathematics

Copyright 2023 Bennett Giles Clayton

ABSTRACT

For many physical problems, robust numerical methods for solving the compressible Euler equations are essential. For the Euler equations to accurately describe the fluid behavior a suitable equation of state (EOS), which describes the relationship between the thermodynamic variables, must be chosen. However, a robust numerical method which can handle an arbitrary equation of state has been unavailable.

In this thesis, we present a second order invariant-domain preserving method for the compressible Euler equations with an arbitrary equation of state. The description of the second order method first requires the development of a first order method that preserves certain thermodynamic properties of the fluid. A method which preserves this physical aspect is referred to as an invariant-domain preserving method. The fundamental methodology of the first order method relies on estimating the maximum wave speed of local Riemann problems. For an arbitrary equation of state this estimation can be impossible. We circumvent the issue by extending the system with an interpolatory EOS, and rigorously justify that the use of the max wave speed of this extended problem implies the invariant-domain preserving properties of the method.

Using a higher order graph viscosity cannot guarantee the invariant-domain preserving property. We resolve this issue through the use of quasiconcave limiting on the density and a surrogate entropy. For an arbitrary equation of state, access to the entropy may not be possible, therefore, this surrogate entropy is used to guarantee that a local approximation of the entropy will increase across shock waves. Furthermore, limiting on the surrogate entropy guarantees that the specific internal energy satisfies the invariant domain constraint.

DEDICATION

To my wife, Unkyung, and daughter, Alice Sena.

ACKNOWLEDGMENTS

I would like to thank my advisor, Dr. Bojan Popov, as well as my committee, Dr. Jean-Luc Guermond, Dr. Matthias Maier, and Dr. Jean Ragusa, for helping me throughout this arduous journey. I would not have been able to achieve this work without their support and guidance. I would also like to thank Texas A&M University for giving me the opportunity to pursue this degree.

I would also like to thank my friends and colleagues: Dr. Eric Tovar and Dr. Weston Baines for providing insightful discussion, practicing presentations, and all around general help for a variety of issues.

Lastly, I thank my wife for surviving with me and supporting me along the way. I may never have even tried to enter into the PhD program without her.

CONTRIBUTORS AND FUNDING SOURCES

Contributors

This work was supported by a dissertation committee consisting of Professor Bojan Popov (advisor), Professor Jean-Luc Guermond, and Professor Matthias Maier, of the Mathematics Department and Professor Jean Ragusa of the Department of Nuclear Engineering. In addition, part of the work was done in collaboration with Dr. Eric Tovar of Los Alamos National Lab. All other work conducted for the thesis was completed by the student independently.

Funding Sources

The graduate study was supported under the following grants: The National Science Foundation DMS-2110868, The Air Force Office of Scientific Research, USAF, under the grant/contract number FA9550-18-1-0397, The Army Research Office under grant number W911NF-19-1-0431, and the U.S. Department of Energy by Lawrence Livermore National Laboratory under the contracts, B640889 and B641173.

NOMENCLATURE

DoF(s)	Degree(s) of Freedom
EOS	Equation of State
FEM	Finite Element Method
JWL	Jones-Wilkins-Lee
NASG	Noble-Abel Stiffened Gas
TAMU	Texas A&M University
vdW	van der Waals

TABLE OF CONTENTS

	Page
ABSTRACT	ii
DEDICATION	iii
ACKNOWLEDGMENTS	iv
CONTRIBUTORS AND FUNDING SOURCES	v
NOMENCLATURE	vi
TABLE OF CONTENTS	vii
LIST OF FIGURES	x
LIST OF TABLES	xii
1. HYPERBOLIC CONSERVATION LAWS	1
1.1 Introduction	1
1.2 Review of the Theory of Conservation Laws	5
1.3 The Riemann Problem*	8
1.4 The Compressible Euler Equations	12
1.5 The Riemann Problem for the Euler Equations	14
2. THE EQUATION OF STATE.....	17
2.1 Relevant Thermodynamic Relations	20
2.2 The Noble-Abel Stiffened Gas EOS	20
2.3 The van der Waals EOS	21
2.4 The Cubic EOS	22
2.4.1 The Redlich-Kwong EOS	23
2.5 The Jones-Wilkins-Lee EOS	25
2.6 The Mie-Gruneisen EOS	26
3. FINITE ELEMENT APPROXIMATION OF THE EULER EQUATIONS	27
3.1 The Continuous Galerkin (cG) Framework	27
3.1.1 Semi-discrete Scheme	29
3.2 The Fully Discrete Scheme	29

4.	THE FIRST ORDER APPROXIMATION*	33
4.1	Invariant Domain Preserving Method	33
4.2	Extended Riemann Problem	35
4.2.1	The Wave Structure	37
4.3	The Solution to the Extended Riemann Problem	40
4.3.1	Shock Wave	42
4.3.2	Rarefaction Wave	46
4.4	Connecting the L- and R-waves	49
4.5	A Weak Solution to the Extended Riemann Problem	51
4.5.1	Weak solution with Vacuum state	54
4.6	Upper Bound on the Max Wave Speed	55
4.6.1	Case 0: Vacuum	57
4.6.2	Case 1: $p^* > 0$ and $\varphi(p_{\min}) > 0$	57
4.6.3	Case 2: $\varphi(p_{\min}) \leq 0 \leq \varphi(p_{\max})$	58
4.6.3.1	Case 2a: $\gamma_{\min} = \gamma_m$	59
4.6.4	Case 2b: $\gamma_{\min} = \gamma_M$	60
4.6.4.1	Case 3: $\varphi(p_{\max}) < 0$	60
5.	HIGH ORDER APPROXIMATION OF THE EULER EQUATIONS WITH TAB- ULATED EOS*	62
5.1	The Use of the Consistent Mass Matrix	62
5.2	A Review of the Entropy Solution	64
5.3	The Entropy Viscosity Method	66
6.	QUASICONCAVE LIMITING*	70
6.1	Local Bounds	70
6.1.1	Relaxation on the Density Bounds	71
6.2	The Flux Corrected Transport Method	73
6.3	Quasiconcave Limiting	74
6.3.1	Quasiconcave Functionals	74
6.3.2	The Abstract Scheme	75
6.3.3	The Limiter	77
6.3.4	Limiting on the Density	78
6.3.5	Limiting on the Internal Energy	79
6.4	The Entropy Surrogate	80
6.4.1	The Entropy for the NASG EOS	80
6.4.2	Limiting on the Surrogate Entropy	86
6.5	The Quadratic Newton Method	86
6.5.1	Review of Divided Differences	88
6.5.2	The Quadratic Newton Method	88
6.6	Relaxation on the Surrogate Entropy	91
7.	NUMERICAL RESULTS*	92

7.1	Convergence Tests	93
7.1.1	The Fan-Jump-Fan Composite Wave	93
7.1.2	Smooth Wave with Various EOS	95
7.1.2.1	Ideal EOS.....	96
7.1.2.2	Van der Waals EOS.....	96
7.1.2.3	Jones-Wilkins-Lee EOS.....	96
7.1.2.4	Mie-Gruneisen EOS	96
7.1.3	The Isentropic Vortex with van der Waals EOS.....	96
7.2	The Two-Expansion Wave Speed Estimate	98
7.2.1	Underestimation of Max Wave Speed: Test 1	99
7.2.2	Underestimation of Max Wave Speed: Test 2	99
7.2.3	Overestimation of Max Wave Speed: Test 3.....	99
7.3	The SESAME Database	104
7.3.1	Expansion-Contact-Shock Comparison.....	105
7.4	Benchmark Configurations	105
7.4.1	EOS Comparison in a Riemann Problem	105
7.4.2	The Woodward-Colella Blast Wave	105
7.4.3	Shock Collision with Triangular Obstacle.....	108
7.4.4	Shock Bubble Interaction	110
7.4.5	Shock Diffraction	112
	REFERENCES	116
	APPENDIX A. Derivation of Exact Solutions for the Euler Equations	125
A.1	Isentropic Vortex*	125
	APPENDIX B. Numerical Algorithms	130
B.1	Upper Bound on Maximum Wave Speed.....	130
B.2	The Quadratic Newton Method for Limiting the Surrogate Entropy	131
B.3	Approximate Godunov-type Solver	132
B.3.1	The Method.....	132
B.3.2	The Algorithm	132

LIST OF FIGURES

FIGURE	Page
1.1	An example Riemann fan consisting of 5 waves and 6 constants states 10
3.1	A typical \mathbb{P}_1 basis function. 29
4.1	An example solution for $\gamma(x, t)$ in the extended Riemann problem 41
6.1	Two example visual descriptions of limiting using the quadratic newton method. Left: quasiconcave function that is not concave. Right: concave function. 90
7.1	Comparison of the exact solution for density (left) and pressure (right) for the fan-jump-fan composite wave. Approximate solutions are computed using 400 and 1600 DoFs and the corresponding mesh sizes are $h = 0.005$ and $h = 0.00125$, respectively. 95
7.2	Plots for Test 1 of the underestimation problem using \widehat{p}^* for computing the maximum wave speed. (Top left): density, (top right): pressure, (bottom left): velocity, (bottom right): sound speed. 100
7.3	Comparison of the sound speed for our method versus the two expansion approximation. This test was run with 800 DoFs up to time $t = 1.2$. The simulation using the two-expansion estimate immediately crashes after $t = 1.2$ generating complex sound speed. 101
7.4	Plots for Test 2 of the underestimation problem using \widehat{p}^* for computing the maximum wave speed. (Top left): density, (top right): pressure, (bottom left): velocity, (bottom right): specific internal energy. Final time is $t = 0.4$ 102
7.5	Plots for Test 3 of the overestimation problem using \widehat{p}^* for computing the maximum wave speed. From left to right: density, pressure, velocity, sound speed. Final time is $t = 0.005$ and using CFL = 1.42. Figures were similar when using λ^{exp} 103
7.6	Comparison of ρ (top left), p (top right), max wave speed, $\widetilde{\lambda}^{\text{max}}$ (bottom left), and γ (bottom right), for the various materials at the final time, $t = 1.2 \times 10^{-5}$ s. 106
7.7	Comparison of the density (top left), pressure (top right), max wave speed (bottom left), and γ (bottom right), for the various different EOS at the final time, $t = 0.1$ 107

7.8	Case 1 of the Woodward-Colella blast wave with the JWL EOS. (Left) density, (right) pressure.	109
7.9	Case 2 of the Woodward-Colella blast wave with the JWL EOS. (Left) density, (right) pressure.	109
7.10	Schlieren plot of a shock wave interacting with a triangular obstacle at $t = 1$ ms, 1.6 ms, and 2.2 ms. Reprinted with permission from [1].	110
7.11	Visual description of the initial state for the shock bubble interaction.	112
7.12	Schlieren plots for the shock-bubble interaction benchmark for $t = 40 \mu\text{s}$. Reprinted with permission from [1].	113
7.13	Schlieren plots for the shock-bubble interaction benchmark for $t = 70 \mu\text{s}$. Reprinted with permission from [1].	114
7.14	Schlieren plots for the shock-bubble interaction benchmark for $t = 100 \mu\text{s}$ (bottom). [1].	114
7.15	Comparison of numerical Schlieren plots for first-order accurate (top row) and second-order accurate (bottom row) solutions at time $t_{\text{final}} = 0.02$ s. The mesh resolution increases from left to right as follows: 1,116,289, 4,460,801, and 17,834,497 \mathbb{Q}_1 -nodes. Reprinted with permission from [1].	115
B.1	Figures describing different solutions of $\mathbf{u}(0, t)$. (Top): the solution in this case is $\mathbf{u}(0, t) = \frac{1}{2}(\hat{\mathbf{u}}_{\text{L}}^* + \hat{\mathbf{u}}_{\text{R}}^*)$. (Middle): $\mathbf{u}(0, t)$ is the solution on the expansion wave which connects the two states, \mathbf{u}_{R} to $\hat{\mathbf{u}}_{\text{R}}^*$ at $x/t = 0$. (Bottom): $\mathbf{u}(0, t) = \hat{\mathbf{u}}_{\text{L}}^*$...	134

LIST OF TABLES

TABLE	Page
2.1	A quick reference for common equations of state in $p = p(\rho, e)$ form. More information on these EOS can be found in Chapter 2. 24
2.2	A quick reference for common sound speeds for several EOS $a = a(\rho, p)$ form. More information on these EOS can be found in Chapter 2. 25
7.1	Consolidated errors and convergence rates for the fan-jump-fan composite wave. Solution computed at $t = 5.0$. Reprinted with permission from [2]. 94
7.2	$\delta_\infty(t_{\text{final}})$ error defined in equation (7.4) and corresponding convergence rates with various EOS for the one-dimensional smooth traveling wave problem with exact solution (7.7) under uniform refinement of the interval $D = (0, 1)$. Reprinted with permission from [1]. 97
7.3	The consolidated error defined in equation (7.4) and convergence rates for the isentropic vortex problem with the Van der Waals EOS. The exact solution is given in (7.8). Reprinted with permission from [1]. 98
7.4	JWL parameters for Woodward-Colella interacting blast wave benchmark. Reprinted with permission from [1]. 108

1. HYPERBOLIC CONSERVATION LAWS

1.1 Introduction

The partial differential equations (PDEs) known as “conservation laws” are used to model a wide variety of physical phenomena. For example, the shallow water models can be used to predict evacuation zones for a storm surge brought about by a hurricane or due to a tsunami. One can also predict and categorize flood zones for which insurance companies and home-buyers use to assess their respective risk. Another active area of research is in compressible flow. Engineers and scientists are researching supersonic aircraft designs which mitigate the sonic boom effect that happens when the aircraft breaks the sound speed barrier. This research is carried out by numerically solving the compressible Euler or Navier-Stokes equations. Furthermore, the study and design of hypersonic objects is now of increasing importance for both national defense and atmospheric reentry vehicles.

These motivations play a major role in the development of robust numerical methods. These robust methods must maintain several important properties. One, the numerical method should preserve the physical properties of quantity or material being studied. For example, preservation of positive water height in the shallow water equations, positive specific internal energy in compressible Euler equations, and so on. Two, the numerical method should produce a “physically relevant” solution. It is known that there are infinitely many weak solutions to a scalar conservation law. In the case of the compressible Euler equations, the numerical method should be able to exclude non-physical solutions through the enforcement of the minimum principle on the specific entropy. Three, the method should not have any tunable parameters that must be adjusted depending on the initial data or computational domain. Four, the method should be scalable. Real applications often require an extremely fine grid to accurately model the object being studied, or operate over a very large spatial domain as in the case of modeling hurricanes. Therefore, it is necessary to simulate

these experiments on a supercomputer. Hence, scalability implies that the numerical method can be programmed efficiently in parallel. Five, the numerical method should be high order. That is, the method converges to the solution rapidly as the mesh size decreases.

Numerical methods for the compressible Euler equations go all the way back to von Neumann [3] and Lax [4] in the early 1950s with the use of finite difference schemes in one dimension. In 1960, Lax and Wendroff [5] introduced a second order method for systems of conservation laws with artificial viscosity to stabilize the method. Also in 1959, Godunov published a new finite difference method, (which could be regarded as a finite volume method), for hydrodynamics. It uses the solution to local Riemann problems at the interface between cells to define the flux on these interfaces. As the original paper was published in Russian in *Matematicheskii Sbornik*; an English translation provided in [6]. In 1969, a novel second order accurate method for hydrodynamics was published by MacCormack [7] (the citation provided is a reprint of the original paper). This method is often referred to as the “predictor-corrector” method. The method works by first computing an intermediate solution using a first order forward-differencing method. This solution is then used to recompute the flux at the interfaces. The final solution is then computed using a backward differencing method but with the new fluxes. This results in a second order method. In 1981, a new approach was developed by Steger and Warming [8] which involves splitting the flux vector into two components based on the positivity or negativity of the eigenvalues of the Jacobian matrix. This method is referred to as the flux-vector splitting method, see also, van Leer [9]. Around the same time, an approximate Riemann solver was used to compute the flux on the cell interface by Roe [10]. Further advances in the use of approximate Riemann solvers appear in the HLL scheme [11], and the HLLC scheme [12]. Since then many new ideas using flux-vector splitting and approximate Riemann solvers were developed like the AUSM scheme by Liou; see [13]. We would also like to point out the popular ENO [14] and WENO [15] schemes which work to minimize oscillations near discontinuities. Lastly, a new approach to numerical methods for conservation laws, referred to as, invariant-domain

preserving methods, began in 2016 with Guermond & Popov [16].

The specific focus of this thesis is on the development of a second order accurate invariant-domain preserving method for the compressible Euler equations with an arbitrary or tabulated equation of state (EOS). That is, the method will be shown to always preserve the necessary thermodynamic properties enforced by the equation of state. For example, positive density and positive specific internal energy. We do this by discretizing the compressible Euler equations with \mathbb{P}_1 (or \mathbb{Q}_1) finite elements and also introduce a “graph viscosity” to the equations which is essential for invariant-domain preserving properties. The numerical update is then written as a convex combination of some “states”. These states motivate the definition of the graph viscosity which is determined by the maximum wave speed to local Riemann problems. It is the determination of this maximum wave speed that is the fundamental difficulty when the equation of state provided is tabulated. There are several numerical methods (e.g. Godunov-type methods) which require the solution to the Riemann problem. Several papers have addressed the issue of the solution to the Riemann problem, but only if the EOS is defined analytically (not tabulated), see [17, Sec. 1], [18], and [19]. Another approach for handling an arbitrary EOS is to use an approximate Riemann solver. See the work done in, Dukowicz [20], Roe & Pike [21], Pike [22], and Lee et. al. [23]. However, these methods make no guarantee on preserving invariant-domain properties other than positivity of the density.

Once the maximum wave speeds are determined, we are able to justify that the states belong to a so-called, “invariant set”. This invariant set is defined by quasiconcave functionals which express the physical quantities we are trying to preserve. This invariant set is slightly different but inspired from the notion of an invariant region described in Chueh et. al. [24], Hoff [25], and Frid [26]. This invariant set is convex and hence a convex combination of states inside a convex set imply that the updated state also belongs to that set. Thus we guarantee that the physical properties we desire, are held.

The extension to second order is done by defining a so-called “entropy viscosity” which is

small in the regions of smooth flow and large when there are discontinuities. The motivating factor for this viscosity measures the change in entropy. I.e., no change in entropy indicates that the flow is smooth and hence we can take the graph viscosity to be small. This second order method by itself is not invariant-domain preserving and a local “limiting” process based on local bounds is required to preserve the physical properties. Limiting is performed on the density and an entropy. However, for a given EOS, we may not know the specific entropy or the specific entropy may not be a concave function of the specific volume and the specific internal energy. Therefore, we suggest the use of a surrogate entropy which behaves similar to a physical entropy; in that, the surrogate entropy increases across shocks.

In Chapter 2 we go over some of the necessary thermodynamic properties for describing the equation of state as well as, a description of several EOS which we use in our numerical demonstrations. In Chapter 3 we give a brief survey of the finite element method of which we employ in the discretization of the problem. Note, however, that the numerical method described in Chapters 4, 5, and 6 is discretization independent. That is, one could equivalently describe the problem in the finite volume or finite difference context; see Remark 4.0.1. In Chapter 4, we derive the first order invariant-domain preserving method for an arbitrary or tabulated equation of state. The estimation of the maximum wave speed for local Riemann problems is detailed here. Then in Chapter 5 we extend the first order method to second order by using the consistent mass matrix and a smaller graph viscosity. It is known that the use of a graph viscosity which is too small can lead to instabilities in multiple forms. This issue is addressed with the use of quasiconcave limiting described in Chapter 6. Lastly, numerical results are presented in Chapter 7.

1.2 Review of the Theory of Conservation Laws

We begin by reviewing some relevant facts regarding the theory of hyperbolic conservation laws. We are interested in the partial differential equation (PDE),

$$\partial_t \mathbf{u} + \nabla \cdot \mathbf{g}(\mathbf{u}) = \mathbf{0}, \quad \text{for } \mathbf{x} \in \mathbb{R}^d, t > 0, \quad (1.1a)$$

$$\mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}) \in \mathcal{B} \subset \mathbb{R}^m, \quad (1.1b)$$

where d denotes the spatial dimension and $\mathbf{u} = \mathbf{u}(\mathbf{x}, t) = (u_1(\mathbf{x}, t), \dots, u_m(\mathbf{x}, t))^\top$ is the unknown vector of conserved quantities with $^\top$ denoting the transpose, and $\mathbf{g} \in C^2(\mathbb{R}^m; \mathbb{R}^{m \times d})$ is the flux. In particular, $\mathbf{g}(\mathbf{u}) = (\mathbf{g}_1(\mathbf{u}), \dots, \mathbf{g}_d(\mathbf{u}))$, where $\mathbf{g}_k : \mathcal{B} \subset \mathbb{R}^m \rightarrow \mathbb{R}^m$ for $k \in \{1 : d\}$ and $\mathbf{g}_k(\mathbf{u}) = (g_{k1}(\mathbf{u}), \dots, g_{km}(\mathbf{u}))^\top$. The set \mathcal{B} is some subset of the phase space and will be explored more in Chapter 4.

The form (1.1) is referred to as the conservative form and can be equivalently written in the quasi-linear form,

$$\partial_t \mathbf{u} + \sum_{k=1}^d \mathbb{A}_k(\mathbf{u}) \partial_{x_k} \mathbf{u} = \mathbf{0}, \quad \text{for } \mathbf{x} \in \mathbb{R}^d, t > 0. \quad (1.2)$$

where $\mathbb{A}_k(\mathbf{u}) = \left(\frac{\partial g_{ki}}{\partial u_j}(\mathbf{u}) \right)_{1 \leq i, j \leq m}$.

Definition 1.2.1 (Hyperbolic System). Define $\mathbb{A}(\mathbf{u}; \mathbf{n}) := \sum_{k=1}^d \mathbb{A}_k(\mathbf{u}) n_k$, where $\mathbf{n} := (n_1, \dots, n_d)^\top$ and $\|\mathbf{n}\|_{\ell^d} = 1$. We say that (1.1) is *hyperbolic*, if $\mathbb{A}(\mathbf{u}; \mathbf{n})$ has m real eigenvalues $\lambda_1(\mathbf{u}; \mathbf{n}), \dots, \lambda_m(\mathbf{u}; \mathbf{n})$, for all $\mathbf{u} \in \mathcal{B}$ and $\mathbf{n} \in \mathbb{S}^{d-1}$. Furthermore, we say that (1.1) is *strictly hyperbolic* if the eigenvalues are all distinct. \square

It is a well known fact that solutions to hyperbolic conservation laws can develop discontinuities in finite time, even if the initial data is smooth. The goal instead, is to find **weak** solutions to (1.1).

Definition 1.2.2 (Weak Solutions). A function, $\mathbf{u} \in [L^\infty(\mathbb{R}^d \times (0, \infty))]^m$ is said to be a

weak solution to (1.1) if

$$\int_{\mathbb{R}^d} \int_0^\infty \mathbf{u} \partial_t \varphi + \mathbf{g}(\mathbf{u}) \nabla \varphi \, dt \, d\mathbf{x} = - \int_{\mathbb{R}^d} \mathbf{u}_0(\mathbf{x}) \varphi(\mathbf{x}, 0) \, d\mathbf{x}, \quad (1.3)$$

holds for all $\varphi \in C_c^1(\mathbb{R}^d \times [0, \infty))$, with C_c^1 denoting compactly supported C^1 functions. \square

For a system of strictly hyperbolic conservation laws in one dimension, under the assumption that all of the characteristic fields are either genuinely nonlinear or linearly degenerate (see Definition 1.2.5), Glimm [27], proved the existence of weak solutions utilizing a random choice method.

In the coming numerical method, we seek physically relevant weak solutions; that is, solutions in the vanishing viscosity sense. This vanishing viscosity solutions will satisfy the so-called, entropy inequalities, see [28, Chapter 11, Theorem 2].

Definition 1.2.3 (Entropy Solutions). A weak solution $\mathbf{u} \in [L^\infty(\mathbb{R}^d \times [0, \infty))]^{d+2}$ to (1.1) is said to be an entropy solution if,

$$\int_{\mathbb{R}^d} \int_0^\infty \eta(\mathbf{u}) \partial_t \varphi(\mathbf{x}, t) + \mathbf{F}(\mathbf{u}) \cdot \nabla \varphi(\mathbf{x}, t) \, dt \, d\mathbf{x} \geq 0 \quad (1.4)$$

for all $\varphi \in C_c^1(\mathbb{R} \times [0, \infty))$ and all entropy, entropy-flux pairs (η, \mathbf{F}) with η convex. In short, we say that $\partial_t \eta(\mathbf{u}) + \nabla \cdot \mathbf{F}(\mathbf{u}) \leq 0$ in the *weak sense*. \square

It is still an open problem on whether the entropy solution will be unique for a hyperbolic system of conservation laws. However, there are results for special cases. For scalar conservation laws, the weak solution only needs to satisfy the so-called Kruzhkov entropies, to ensure uniqueness of the solution, see Kruzhkov [29]. For the Cauchy problem with viscosity, $\partial_t \mathbf{u} + \mathbb{A}(\mathbf{u}) \partial_x \mathbf{u} = \varepsilon \partial_{xx} \mathbf{u}$, if $\mathbb{A}(\mathbf{u})$ is strictly hyperbolic and if the initial data is small enough; that is, $\|\mathbf{u}\|_{\text{BV}} < \delta$ for some $\delta > 0$, then there exists a unique solution $\mathbf{u}^\varepsilon(x, t)$. This solution, \mathbf{u}^ε is referred to as the viscosity solution. For $\varepsilon \rightarrow 0^+$, we have that \mathbf{u}^ε converges to the unique solution to $\partial_t \mathbf{u} + \mathbb{A}(\mathbf{u}) \partial_x \mathbf{u} = \mathbf{0}$. This solution is referred to as the *vanishing*

viscosity solution. For the general system $\partial_t \mathbf{u} + \partial_x \mathbf{g}(\mathbf{u}) = \mathbf{0}$, the vanishing viscosity solution is an entropy solution. This solution was shown in Bianchini & Bressan [30] to be the same solution derived by Glimm in [27]. For more on the theory of vanishing viscosity solutions, see [30]. Despite the lack of results for systems in higher dimensions, entropy solutions are still sought after since they exclude many non-physical solutions.

One of the main motivations of our numerical method, the invariant-domain preserving method, is that the solution should reside in region of the phase space for which the physical properties are satisfied. This motivates the definition of an *invariant region*.

Definition 1.2.4 (Invariant Region). A set $\mathcal{B} \subset \mathbb{R}^m$ is said to be an invariant region for the PDE (1.1) if $\mathbf{u}(\mathbf{x}, 0) \in \mathcal{B}$ for all $\mathbf{x} \in \mathbb{R}^d$ and there exists $\delta > 0$ such that a solution $\mathbf{u}(\mathbf{x}, t)$ exists for all $0 < t < \delta$, and if $\mathbf{u}(\mathbf{x}, t) \in \mathcal{B}$, then we say that \mathcal{B} is an invariant region. \square

Remark 1.2.1. This set \mathcal{B} will be defined in terms of some quasiconcave functionals. That is, let $\{\Psi_i\}_{i=1}^N$ be a collection of quasiconcave functionals, then

$$\mathcal{B} := \{\mathbf{u} \in \mathbb{R}^m : \Psi_i(\mathbf{u}) > 0, i = 1, \dots, N\}, \quad (1.5)$$

is an invariant region. This definition was originally motivated by Chueh et. al. [24] in the context of convection-diffusion problems and the functionals, $\{\Psi_i\}$ are chosen to be convex, with the constraints $\Psi_i(\mathbf{u}) \leq 0$. \square

As much of the upcoming material (see Chapter 4) requires analysis of the Riemann problem which is one-dimensional, we switch our focus to conservation laws in one-dimension.

Consider the one-dimensional problem,

$$\partial_t \mathbf{u} + \partial_x (\mathbf{g}(\mathbf{u})) = \mathbf{0}, \quad \text{for } x \in \mathbb{R}, t > 0. \quad (1.6)$$

Let $\mathbf{A}(\mathbf{u})$ be the Jacobian matrix for the flux $\mathbf{g}(\mathbf{u})$. Assume the eigenvalues of $\mathbf{A}(\mathbf{u})$ are all real and distinct, $\lambda_1(\mathbf{u}) < \dots < \lambda_m(\mathbf{u})$. Let $\mathbf{r}_1(\mathbf{u}), \dots, \mathbf{r}_m(\mathbf{u})$ be the associated right

eigenvectors. Then we have the following definitions,

Definition 1.2.5 (Wave Type Definitions). For the one-dimensional conservation law, (1.6), we have the following definitions,

- We call $(\lambda_i(\mathbf{u}), \mathbf{r}_i(\mathbf{u}))$ the *i-characteristic field*.
- The *i-characteristic field* is said to be *genuinely nonlinear* if $D\lambda_i(\mathbf{u}) \cdot \mathbf{r}_i(\mathbf{u}) \neq 0$ for all $\mathbf{u} \in \mathcal{B}$.
- The *i-characteristic field* is said to be *linearly degenerate* if, $D\lambda_i(\mathbf{u}) \cdot \mathbf{r}_i(\mathbf{u}) = 0$ for all $\mathbf{u} \in \mathcal{B}$.

□

Definition 1.2.6 (k-Riemann invariant). A smooth function $w : \mathcal{B} \rightarrow \mathbb{R}$ is said to be a *k-Riemann invariant*, if $Dw(\mathbf{u}) \cdot \mathbf{r}_k(\mathbf{u}) = 0$ for all $\mathbf{u} \in \mathcal{B}$.

□

Definition 1.2.7 (Rankine-Hugoniot Conditions). The Rankine-Hugoniot conditions are defined as,

$$S(\mathbf{u}_L - \mathbf{u}_R) = \mathbf{g}(\mathbf{u}_L) - \mathbf{g}(\mathbf{u}_R), \quad (1.7)$$

where S is the instantaneous speed of the discontinuity and \mathbf{u}_L and \mathbf{u}_R are the instantaneous states to the left and right of the discontinuity.

□

1.3 The Riemann Problem*

A large number of numerical methods for conservation laws require, in some way, the concept of a Riemann problem.

Definition 1.3.1. Let $\mathbf{g} \in C^1(\mathbb{R}^m; \mathbb{R}^m)$ be the flux, then the Riemann problem is written

* Lemma 1.3.1 and its proof are taken from [2] and are reprinted with permission from [2].

as: Find a self-similar (weak) solution $\mathbf{u} \in L^\infty(\mathbb{R} \times (0, \infty); \mathbb{R}^m) \cap C^0((0, \infty); L^1_{\text{loc}}(\mathbb{R}; \mathbb{R}^m))$ to

$$\partial_t \mathbf{u} + \partial_x \mathbf{g}(\mathbf{u}) = \mathbf{0}, \quad (1.8a)$$

$$\mathbf{u}(x, 0) = \begin{cases} \mathbf{u}_L, & x < 0, \\ \mathbf{u}_R, & x > 0, \end{cases} \quad (1.8b)$$

□

We then have a fundamental result for the existence and uniqueness of solutions to the Riemann problem,

Theorem 1.3.1 (Existence and Uniqueness). *Assume the system (1.8a) is strictly hyperbolic and the characteristic fields are either genuinely nonlinear or linearly degenerate. If $\|\mathbf{u}_L - \mathbf{u}_R\| < \delta$ for $\delta > 0$, is sufficiently small, then there exists a unique self-similar solution to the Riemann problem.*

For the proof of this theorem see, [31, Chapt. I. Theorem 6.1]. The Riemann problem was first studied by Riemann in the context of the isentropic Euler equations in the seminal paper [32]; an English translation can be found in [33, pg. 109]. In particular, the Riemann solution can be constructed using Lax's method, see [34, Sec. 9]. In particular, the solution to the Riemann problem consists of $m + 1$ constant states separated by m waves. That is, there exists $2m$ numbers,

$$\lambda_1^- \leq \lambda_1^+ \leq \lambda_2^- \leq \lambda_2^+ \leq \cdots \leq \lambda_m^- \leq \lambda_m^+. \quad (1.9)$$

These $2m$ numbers define $2m + 1$ sectors in the (x, t) -plane. In particular, they can be defined through the self-similarity parameter, $\xi := x/t$. That is, ξ belongs to one of the intervals, $(-\infty, \lambda_1^-)$, $(\lambda_1^-, \lambda_1^+)$, $(\lambda_1^+, \lambda_2^-)$, \dots , $(\lambda_m^-, \lambda_m^+)$, and (λ_m^+, ∞) . Note, some of these intervals can be empty. If the interval $(\lambda_i^-, \lambda_i^+)$ is non-empty for $i \in \{1 : m\}$, then the solution inside this interval must be a rarefaction wave. All other wave solutions are either shocks or contact

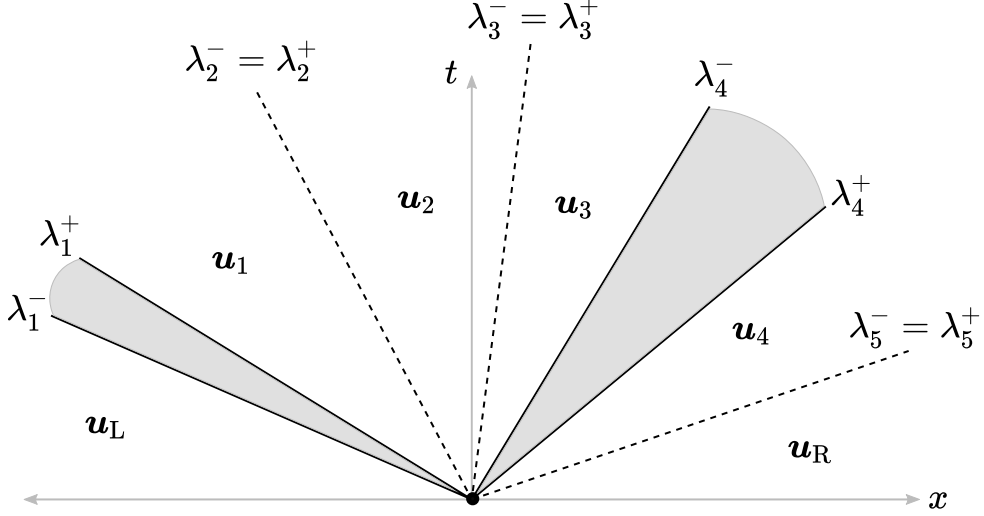


Figure 1.1: An example Riemann fan consisting of 5 waves and 6 constants states

discontinuities. Note that $\mathbf{u}(x, t) = \mathbf{u}_L$ for $\frac{x}{t} \in (-\infty, \lambda_1^-)$ and $\mathbf{u}(x, t) = \mathbf{u}_R$ for $\frac{x}{t} \in (\lambda_m^+, \infty)$. See Figure 1.1 for a visual description.

The main motivation for introducing these $2m$ numbers is that the numerical method described in Chapter 4 relies on estimation of the maximum speed of propagation. That is, we need an upper bound on

$$\lambda_{\max} := \max\{|\lambda_1^-|, |\lambda_m^+|\}. \quad (1.10)$$

For the purposes of our numerical method, we modify slightly the definition of the invariant region, see Definition 1.2.4.

Definition 1.3.2 (Invariant Set). We say that $\mathcal{B} \subset \mathbb{R}^m$ is an **invariant set** for (1.1) if for every pair $(\mathbf{U}_L, \mathbf{U}_R) \in \mathcal{B} \times \mathcal{B}$ we have that average of the entropy solution to (1.8a)–(1.8b) with $\mathbf{g}(\mathbf{u}) := g(\mathbf{u})\mathbf{n}$ for any $\mathbf{n} \in \mathbb{S}^{d-1}$ over the Riemann fan, remains in \mathcal{B} ; that is,

$$\frac{1}{t(\lambda_m^+ - \lambda_1^-)} \int_{\lambda_1^- t}^{\lambda_m^+ t} \mathbf{u}(x, t) dx \in \mathcal{B} \quad (1.11)$$

for any $t > 0$. □

Lemma 1.3.1. Let $\mathbf{u} \in L^\infty(\mathbb{R} \times (0, \infty); \mathbb{R}^m) \cap C^0((0, \infty); L^1_{\text{loc}}(\mathbb{R}; \mathbb{R}^m))$ be a solution to the

Riemann problem (1.8a)–(1.8b) with $\mathbf{u}_L, \mathbf{u}_R \in \mathcal{B}$. Then the following holds for all $t \in (0, \frac{1}{2\lambda_{\max}}]$ where λ_{\max} is the maximal wave speed to the Riemann problem defined in (1.10),

1. $\bar{\mathbf{u}}(t) := \int_{-1/2}^{1/2} \mathbf{u}(x, t) dx = \frac{1}{2}(\mathbf{u}_L + \mathbf{u}_R) - t(\mathbf{g}(\mathbf{u}_R) - \mathbf{g}(\mathbf{u}_L)).$

2. $\bar{\mathbf{u}}(t) \in \mathcal{B}.$

3. Let $\Psi \in C^1(\mathcal{B}; \mathbb{R})$ be a quasiconcave functional. Assume that $\Psi(\mathbf{u}(x, t)) \geq 0$, for a.e. $x \in \mathbb{R}$ and all $t > 0$. Then $\Psi(\bar{\mathbf{u}}) \geq 0$.

4. Let $\Psi \in C^1(\mathcal{B}, \mathbb{R})$ be a concave functional. Assume that $\Psi(\mathbf{u}(x, t)) \geq 0$ for a.e. $x \in \mathbb{R}$ and all $t > 0$. Assume that there exists $\lambda_b, \lambda_{\sharp} \in [-\lambda_{\max}, \lambda_{\max}]$, $\lambda_b < \lambda_{\sharp}$, so that $\Psi(\mathbf{w}(x, t)) > 0$ for a.e. $\frac{x}{t} \in (\lambda_b, \lambda_{\sharp})$. Then $\Psi(\bar{\mathbf{w}}(t)) > 0$.

Proof. The proof of this lemma is taken from Clayton et. al. [2, Lemma 3.2]. For the entire proof t is a fixed real number in $(0, \frac{1}{2\lambda_{\max}})$.

(i) Let u_1, \dots, u_m be the m components of \mathbf{u} , and let g_1, \dots, g_m be the m components of the flux \mathbf{g} . Let $l \in \{1 : m\}$. Since \mathbf{u} is a weak solution to (1.8a), we have

$$0 = \int_{-\infty}^{\infty} \int_0^{\infty} (-u_l \partial_{\tau} \phi - g_l(\mathbf{u}) \partial_x \phi) d\tau dx - u_{l,L} \int_{-\infty}^0 \phi(x, 0) dx - u_{l,R} \int_0^{\infty} \phi(x, 0) dx$$

for all $\phi \in W^{1,\infty}(\mathbb{R} \times [0, \infty); \mathbb{R})$ with compact support in $\mathbb{R} \times [0, \infty)$. Here, $u_{l,Z}$ is the l -th component of \mathbf{u}_Z . Now we define a sequence of smooth functions $(\phi_{\epsilon})_{\epsilon > 0}$ with $\phi_{\epsilon}(x, \tau) = \phi_{1,\epsilon}(|x|) \phi_{2,\epsilon}(\tau)$,

$$\phi_{1,\epsilon}(x) = \begin{cases} 1 & 0 \leq x \leq \frac{1}{2}, \\ \frac{1}{\epsilon}(-x + \frac{1}{2} + \epsilon) & \frac{1}{2} \leq x \leq \frac{1}{2} + \epsilon, \\ 0 & \frac{1}{2} + \epsilon \leq x, \end{cases} \quad \phi_{2,\epsilon}(\tau) = \begin{cases} 1 & 0 \leq \tau \leq t, \\ \frac{1}{\epsilon}(-\tau + t + \epsilon) & t \leq \tau \leq t + \epsilon, \\ 0 & t + \epsilon \leq \tau. \end{cases}$$

Using that $u_l \in C^0([0, \infty); L^1_{\text{loc}}(\mathbb{R}))$, we infer that $\int_{-\infty}^{\infty} \int_0^{\infty} -u_l \partial_{\tau} \phi_{\epsilon} dx d\tau \rightarrow \int_{-\frac{1}{2}}^{\frac{1}{2}} u_l(x, t) dx$ as $\epsilon \rightarrow 0$. Likewise, we have $\int_{-\infty}^{\infty} \int_0^{\infty} -g_l(\mathbf{u}) \partial_x \phi_{\epsilon} d\tau dx \rightarrow \int_0^t (g_l(\mathbf{u}_R) - g_l(\mathbf{u}_L)) d\tau = (g_l(\mathbf{u}_R) -$

$g_l(\mathbf{u}_L))t$ as $\epsilon \rightarrow 0$. Finally, $-u_{l,L} \int_{-\infty}^0 \phi_\epsilon(x, 0)dx - u_{l,R} \int_0^\infty \phi_\epsilon(x, 0)dx \rightarrow -\frac{1}{2}(u_{l,L} + u_{l,R})$ as $\epsilon \rightarrow 0$. In conclusion, we have established that

$$0 = \bar{\mathbf{u}}(t) + (\mathbf{g}(\mathbf{u}_R) - \mathbf{g}(\mathbf{u}_L))t - \frac{1}{2}(\mathbf{u}_L + \mathbf{u}_R). \quad (1.12)$$

(ii) Since \mathcal{B} is convex, $\mathbf{u}(x, t) \in \mathcal{B}$ for $\forall x \in \mathbb{R}$ and all $t > 0$, and the length of the interval $[-\frac{1}{2}, \frac{1}{2}]$ is 1, we infer that $\bar{\mathbf{u}}(t) \in \mathcal{B}$.

(iii) Let $\Psi \in C^1(\mathcal{B}; \mathbb{R})$ be a quasiconcave functional. The quasiconcavity implies that $\Psi(\bar{\mathbf{u}}(t)) \geq \text{ess inf}_{x \in (-\frac{1}{2}, \frac{1}{2})} \Psi(\mathbf{u}(x, t)) \geq 0$.

(iv) Let $\Psi \in C^1(\mathcal{B}; \mathbb{R})$ be a concave functional. Jensen's inequality implies

$$\Psi(\bar{\mathbf{u}}(t)) \geq \int_{-\frac{1}{2}}^{\frac{1}{2}} \Psi(\mathbf{u}(x, t))dx \geq \int_{\lambda_b t}^{\lambda_\sharp t} \Psi(\mathbf{u}(x, t))dx > 0, \quad (1.13)$$

where we used $-\frac{1}{2} \leq \lambda_b t < \lambda_\sharp t \leq \frac{1}{2}$. This concludes the proof. \square

Results 1. and 2. in Lemma 1.3.1 are used in Chapter 4 to prove that the numerical method is invariant-domain preserving. Results 3. and 4. are used for the purposes of quasiconcave limiting, see Chapter 6. For more information regarding the existence of self similar solutions to the Riemann problem see, Dafermos [35, Chapter IX].

1.4 The Compressible Euler Equations

The main focus of this thesis is on numerical methods for solving the compressible Euler equations. These equations are a fundamental and important model for simulating fluid flow with very small viscosity (or, equivalently, a very large Reynolds number). They also form the basis for more complicated multi-physics models.

The Euler equations represent the conservation of mass, momentum and total energy;

they are given respectively as,

$$\partial_t \rho(\mathbf{x}, t) + \nabla \cdot (\rho(\mathbf{x}, t) \mathbf{v}(\mathbf{x}, t)) = 0, \quad (\mathbf{x}, t) \in \mathbb{R}^d \times (0, \infty), \quad (1.14a)$$

$$\partial_t \mathbf{m}(\mathbf{x}, t) + \nabla \cdot \left(\frac{\mathbf{m}(\mathbf{x}, t)}{\rho(\mathbf{x}, t)} \otimes \mathbf{m}(\mathbf{x}, t) + p(\rho, \mathbf{e}(\mathbf{u})) \mathbb{I}_d \right) = \mathbf{0}, \quad (\mathbf{x}, t) \in \mathbb{R}^d \times (0, \infty), \quad (1.14b)$$

$$\partial_t E(\mathbf{x}, t) + \nabla \cdot \left(\frac{\mathbf{m}(\mathbf{x}, t)}{\rho(\mathbf{x}, t)} (E(\mathbf{x}, t) + p(\rho, \mathbf{e}(\mathbf{u}))) \right) = 0, \quad (\mathbf{x}, t) \in \mathbb{R}^d \times (0, \infty), \quad (1.14c)$$

where ρ is the density, $\mathbf{m} = (m_1, \dots, m_d)$ is the momentum, E is the total energy, \otimes is the outer product, d is the spatial dimension (where d is either 1, 2, or 3) and $\mathbf{u} = (\rho, \mathbf{m}, E)^T$. The pressure mapping $\mathbb{R}^2 \supset \mathcal{A} \ni (\rho, e) \mapsto p(\rho, e) \in \mathbb{R}$ is defined by an *equation of state* on some suitable thermodynamic region, \mathcal{A} , for which the equation of state can be inverted as, $e = e(\rho, p)$; see Chapter 2. Furthermore, we refer to the pressure, $p = p(\rho, e)$ as the **oracle** as we make no assumptions on the exact description of the equation of state (EOS). The total energy is, $E = \rho e + \frac{1}{2} \rho \|\mathbf{v}\|_{\ell^2}$, where ρe is the internal energy and $\frac{1}{2} \rho \|\mathbf{v}\|_{\ell^2}$ is the kinetic energy with $\mathbf{v} := \mathbf{m}/\rho$ being the velocity and $\|\cdot\|_{\ell^2}$ the usual Euclidean norm. Solving for e , we can also write, $e(\rho, p) = \mathbf{e}(\mathbf{u}) := \frac{E}{\rho} - \frac{\|\mathbf{m}\|_{\ell^2}^2}{2\rho^2}$. We write the Euler equations in the more general form,

$$\partial_t \mathbf{u} + \nabla \cdot \mathbf{f}(\mathbf{u}) = \mathbf{0}, \quad (1.15)$$

with,

$$\mathbf{f}(\mathbf{u}) := \begin{pmatrix} \mathbf{m}^\top \\ \frac{\mathbf{m}}{\rho} \otimes \mathbf{m} + p \mathbb{I}_d \\ \frac{\mathbf{m}^\top}{\rho} (E + p) \end{pmatrix}. \quad (1.16)$$

Additionally, for practical purposes, we must numerically solve the Euler equations on a bounded domain $D \subset \mathbb{R}^d$. This requires the implementation of boundary conditions. Implementation of boundary conditions is reserved for Chapter 7.

Remark 1.4.1 (Assumptions on the Oracle). The only assumption we make on the equation of state is that, $p + p_\infty \geq 0$ where p_∞ is a reference pressure state or is the absolute value of the global minimum pressure. If p_∞ is, a priori, not known, then we set $p_\infty = 0$ and require

$p \geq 0$. □

Remark 1.4.2 (Parameters Given by the Oracle). There are certain parameters that the oracle may possess that can be used in the approximation of the problem. We list them here. The *maximum compressibility* (or *covolume*) constant, $b > 0$ which represents the smallest volume the fluid can occupy; that is, $\tau - b > 0$ or $\rho < b^{-1}$. A reference specific internal energy, $q \in \mathbb{R}$, so that $e - q > 0$. A reference pressure p_∞ as described in Remark 1.4.1. The exact use of these parameters will be evident in Chapter 4. □

1.5 The Riemann Problem for the Euler Equations

As the compressible Euler equations are a multi-dimensional system of conservation laws, the notion of a Riemann problem is not necessarily clear. For many numerical methods, the Riemann problem for the Euler equations is defined in terms of a direction, \mathbf{n} . Let $\mathbf{n} \in \mathbb{S}^{d-1}(\mathbf{0}, 1)$ be given, then define the orthonormal basis, $\{\mathbf{n}, \mathbf{t}_1, \dots, \mathbf{t}_{d-1}\}$. The exact choice of $\{\mathbf{t}_i\}_{i=1}^{d-1}$ is not important. With respect to this basis, we write the momentum as, $\mathbf{m} = (m, \mathbf{m}^\perp)^\top$ where $m := \mathbf{m} \cdot \mathbf{n}$ and $\mathbf{m}^\perp := (\mathbf{m} \cdot \mathbf{t}_1, \dots, \mathbf{m} \cdot \mathbf{t}_{d-1})^\top$. The velocity is defined similarly, $\mathbf{v} = (v, \mathbf{v}^\perp)^\top$ where $v := m/\rho$ and $\mathbf{v}^\perp := \mathbf{m}^\perp/\rho$. Thus the Riemann problem in the direction \mathbf{n} is defined by,

$$\partial_t \mathbf{u} + \partial_x (\mathbf{f}(\mathbf{u})\mathbf{n}) = \mathbf{0}, \quad \text{for } x \in \mathbb{R}, t > 0 \tag{1.17}$$

$$\mathbf{u}(x, 0) = \begin{cases} \mathbf{u}_L, & x < 0, \\ \mathbf{u}_R, & x > 0, \end{cases} \tag{1.18}$$

This can be explicitly written as,

$$\partial_t \begin{pmatrix} \rho \\ m \\ \mathbf{m}^\perp \\ E \end{pmatrix} + \partial_x \begin{pmatrix} m \\ \frac{1}{\rho} m^2 + p(\mathbf{u}) \\ \frac{m}{\rho} \mathbf{m}^\perp \\ \frac{m}{\rho} (E + p(\mathbf{u})) \end{pmatrix} = \mathbf{0} \tag{1.19}$$

The solution to this Riemann problem is carried out by first solving,

$$\partial_t \begin{pmatrix} \rho \\ m \\ \mathcal{E} \end{pmatrix} + \partial_x \begin{pmatrix} m \\ \frac{1}{\rho}m^2 + p(\mathbf{u}) \\ \frac{m}{\rho}(\mathcal{E} + p(\mathbf{u})) \end{pmatrix} = \mathbf{0} \quad (1.20)$$

where $\mathcal{E} := E - \frac{\|\mathbf{m}^\perp\|_2^2}{2\rho}$. Then \mathbf{m}^\perp is found afterwards by solving $\partial_t \mathbf{m}^\perp + \partial_x (\frac{m}{\rho} \mathbf{m}^\perp) = \mathbf{0}$. Note also that $\rho e = \mathcal{E} - \frac{m^2}{2\rho} = E - \frac{\|\mathbf{m}\|_2^2}{2\rho}$; this says that the internal energy does not depend on the change of basis. More on the solution to the Riemann problem will be presented in Chapter 4.

Remark 1.5.1 (Invariant Sets for the Euler Equations). The compressible Euler equations depend on the EOS, therefore, each EOS may define a different invariant set. For example, if the EOS is given by the ideal gas law, $p(\rho, e) = (\gamma - 1)\rho e$, then, a possible invariant set would be,

$$\mathcal{B} := \{\mathbf{u} \in \mathbb{R}^{d+2} : \rho > 0, \mathbf{e}(\mathbf{u}) > 0, \mathbf{s}(\mathbf{u}) \geq s_{\min}\}. \quad (1.21)$$

Note $\mathbf{s}(\mathbf{u}) \geq s_{\min}$ is the minimum principle on the specific entropy and is discussed more in Chapter 2. Furthermore, the sets, $\mathcal{B} := \{\mathbf{u} \in \mathbb{R}^{d+2} : \rho > 0, \mathbf{e}(\mathbf{u}) > 0\}$. and $\mathcal{B} := \{\mathbf{u} \in \mathbb{R}^{d+2} : \rho > 0\}$ are also invariant sets.

If the EOS is given by the covolume EOS, $p(\rho, e) = (\gamma - 1)\frac{\rho e}{1 - b\rho}$ then an invariant set involves the maximum compressibility constant, b^{-1} (for $b \neq 0$). That is,

$$\mathcal{B}(b) := \{\mathbf{u} \in \mathbb{R}^{d+2} : 0 < \rho < b^{-1}, \mathbf{e}(\mathbf{u}) > 0, \mathbf{s}(\mathbf{u}) \geq s_{\min}\}. \quad (1.22)$$

As we will see through the course of this thesis, we will be working with the Nobel-Abel Stiffened Gas EOS (see Section 2.2). An invariant set for this EOS is,

$$\mathcal{B}(b, q, p_\infty) := \{\mathbf{u} \in \mathbb{R}^{d+2} : 0 < \rho < b^{-1}, \mathbf{e}(\mathbf{u}) - q > p_\infty(\rho^{-1} - b)\}. \quad (1.23)$$

□

Lastly, we would also like to present the Rankine-Hugoniot conditions for the compressible Euler equations for reference.

Definition 1.5.1. The Rankine-Hugoniot conditions for the compressible Euler equations are given by,

$$S(\rho_L - \rho_R) = \rho_L v_L - \rho_R v_R, \tag{1.24a}$$

$$S(\rho_L v_L - \rho_R v_R) = \rho_L v_L^2 + p_L - (\rho_R v_R^2 + p_R), \tag{1.24b}$$

$$S(\mathcal{E}_L - \mathcal{E}_R) = v_L(\mathcal{E}_L + p_L) - v_R(\mathcal{E}_R + p_R). \tag{1.24c}$$

□

2. THE EQUATION OF STATE

In this chapter we cover specific topics pertaining to the many different EOS.

The *equation of state* (EOS) is one of the most important aspects of the Euler equations, it is an equation which relates three thermodynamic quantities; for example, the density, temperature and pressure or volume, specific internal energy, and specific entropy. The choice of the EOS is chosen based on the physical problem of interest. In particular, the EOS models the material. The EOS are often designed to accurately model specific materials under certain conditions. Typically, the Euler equations are numerically solved using the ideal gas law as it is the simplest. In addition, most theoretical results for the Euler equations are proven for the ideal gas law. There are, in fact, hundreds of different EOS and numerically solving the Euler equations with invariant domain preserving properties becomes a highly non-trivial task for each of these equations. To further exacerbate the issue, many industrial and research laboratories rely on EOS which are tabulated; that is, there is no analytic function, but rather, thermodynamic quantities are generated from a database.

We begin by reviewing some essential principles in thermodynamics. The first law of thermodynamics states that the change in energy in a non-adiabatic system; that is, a system for which energy can be transferred through the walls of the system, is equal to the change in heat plus the work done on that system. Letting Q denote the heat, the first law of thermodynamics is written as,

$$de = dQ - pdV. \tag{2.1}$$

Since the EOS can be written in terms of many different thermodynamic variables, it is often convenient to overload the notation. For example, we may use the same symbol p to write $p = p(\tau, e)$ and $p = p(\rho, s)$. The context should be clear but comments will be made when necessary to avoid confusion. To further emphasize the variables being used, the notation for partial differentiation is written in the form $(\frac{\partial f}{\partial x})_y$, where $(\cdot)_y$ emphasizes

that the variable held constant is y . This also displays that f is a function of x and y .

A common assumption for the EOS is the existence of a physical entropy. That is, the assumption that there is a function $s = s(\tau, e)$ which satisfies the second law of thermodynamics

$$T ds = de + p d\tau. \quad (2.2)$$

We must assume that the entropy, $s(\tau, e)$, is strictly convex (as a function of (τ, e)). Otherwise, the pressure can be multivalued as a function of τ and s . Furthermore, from Clairaut's theorem we have

$$\frac{\partial}{\partial \tau} \left(\frac{\partial s}{\partial e}(\tau, e) \right) = \frac{\partial}{\partial e} \left(\frac{\partial s}{\partial \tau}(\tau, e) \right). \quad (2.3)$$

Using appropriate identities for $\left(\frac{\partial s}{\partial e}\right)_\tau$ and $\left(\frac{\partial s}{\partial \tau}\right)_e$, (2.3) gives us the so-called, *Maxwell's relation*, see Callen [36, Chapter 7]. However, the existence of an entropy is not always guaranteed.

Definition 2.0.1 (Incomplete EOS). An *incomplete* equation of state is a relation of the form, $p = p(\rho, e)$. It is called *incomplete* as the equation cannot completely describe all of the thermodynamical properties of the system. \square

Definition 2.0.2 (Complete EOS). A *complete* equation of state is one for which a concave entropy, $s = s(\tau, e)$, exists such that $p = p(\tau, s)$. Or equivalently, the mapping $(\tau, e) \mapsto -s(\tau, e)$, is convex. This equation of state completely describes all of the thermodynamic properties of the system. \square

So an *incomplete* EOS makes no assumption on the existence of a convex entropy. Any given complete EOS describes a unique incomplete EOS; however, the converse is not true, there can be more than one complete EOS corresponding to the same incomplete EOS. For more on this, see Menikoff & Plohr [37, Sec. II. F.]. Note that the first order method described in Chapter 4 does not require a complete EOS. However, the second order method in some way requires the use of an entropy (which may not exist), this is resolved in Chapter 5 and Chapter 6.

We now review some important facts regarding the entropy. First, concavity of $s(\tau, e)$ requires that

$$\left(\frac{\partial^2 s}{\partial \tau^2}\right)_e < 0, \quad \left(\frac{\partial^2 s}{\partial e^2}\right)_\tau < 0, \quad \text{and} \quad \left(\frac{\partial^2 s}{\partial \tau^2}\right)_e \left(\frac{\partial^2 s}{\partial e^2}\right)_\tau - \left(\frac{\partial^2 s}{\partial \tau \partial e}\right)^2 < 0, \quad (2.4)$$

i.e., the Hermitian matrix of $-s(\tau, e)$ is positive definite.

Remark 2.0.1 (Convex Entropy in the Literature). In the literature, concavity of the entropy is often framed as $(\tau, e) \mapsto -s(\tau, e)$ being convex. \square

Remark 2.0.2 (Assumptions on the Oracle). Throughout the course of this Thesis, the only assumption we make on the oracle, is that there exists a minimum bound on the pressure, $-p_\infty$, for $p_\infty \geq 0$. Furthermore, we make no assumption on the existence of an entropy for the oracle. Other remarks will be made regarding the oracle when relevant. \square

Note that a convex entropy From Equation (2.2), we have the following identities,

$$\left(\frac{\partial s}{\partial e}\right)_\tau = T^{-1}, \quad \text{and} \quad \left(\frac{\partial s}{\partial \tau}\right)_e = pT^{-1}. \quad (2.5)$$

The pressure can then be defined by,

$$p(\tau, e) = \frac{\left(\frac{\partial s}{\partial \tau}\right)_e}{\left(\frac{\partial s}{\partial e}\right)_\tau} \quad (2.6)$$

or by,

$$\frac{\partial s}{\partial \tau} - p(\tau, e) \frac{\partial s}{\partial e} = 0. \quad (2.7)$$

We also have $\left(\frac{\partial e}{\partial \tau}\right)_s = -p(\tau, s)$ and $\left(\frac{\partial e}{\partial s}\right)_\tau = T$ from Equation (2.2). It is also assumed that the EOS must satisfy $\left(\frac{\partial e}{\partial s}\right)_\tau = T > 0$.

2.1 Relevant Thermodynamic Relations

The material speed of sound is defined by,

$$a := \sqrt{\left(\frac{\partial p}{\partial \rho}\right)_s} = \sqrt{-\tau^2 \left(\frac{\partial p}{\partial \tau}\right)_s}. \quad (2.8)$$

As with the pressure, we also leave the independent variables of the speed of sound unspecified; the appropriate variables should be clear from context. Finding the specific entropy, s , is not always an easy task. For an incomplete EOS the sound speed can more easily be computed by applying the chain rule and (2.2), to (2.8) to find,

$$a = \sqrt{-\tau^2 \left[\left(\frac{\partial p}{\partial \tau}\right)_e + \left(\frac{\partial p}{\partial e}\right)_\tau \left(\frac{\partial e}{\partial \tau}\right)_s \right]} = \sqrt{-\tau^2 \left[\left(\frac{\partial p}{\partial \tau}\right)_e - p \left(\frac{\partial p}{\partial e}\right)_\tau \right]}. \quad (2.9)$$

The specific heat of a fluid is defined by $\frac{dQ}{dT}$. From the first law of thermodynamics, we have, $dQ = de + p d\tau = T ds$. We then are able to define the specific heat at constant volume and constant pressure using the first law of thermodynamics by,

$$c_v := \left(\frac{\partial e}{\partial T}\right)_\tau = T \left(\frac{\partial s}{\partial T}\right)_\tau \quad c_p := T \left(\frac{\partial s}{\partial T}\right)_p = \left(\frac{\partial(e + p\tau)}{\partial T}\right)_p = \left(\frac{\partial h}{\partial T}\right)_p. \quad (2.10)$$

2.2 The Noble-Abel Stiffened Gas EOS

A large portion of the material in this thesis relies on the use of the Noble-Abel Stiffened Gas EOS (NASG EOS). This equation of state was first introduced by Le M'etayer & Saurel in [38] and has been extended to more general form in [39]. The NASG EOS can also be viewed as an extension of the Noble-Abel (or covolume) EOS (see [40] for a more recent discussion and use of the Noble-Abel EOS).

The NASG EOS is defined by the incomplete EOS,

$$p(\tau, e) = (\gamma - 1) \frac{e - q}{\tau - b} - \gamma p_\infty \quad (2.11)$$

where q is some reference internal energy, p_∞ is also a reference pressure state, and b represents the maximum compressibility of the fluid; that is, the fluid cannot be compressed to an infinitely small volume. It can be shown that the specific entropy is given by,

$$s(\tau, e) = \log \left((e - q - p_\infty(\tau - b))^{\frac{1}{\gamma-1}} (\tau - b) \right). \quad (2.12)$$

That is, $s(\tau, e)$ satisfies equation (2.7). It is important to note that this equation of state is convex, see [38, Appendix B], which is necessary in the solution to the extended Riemann problem in Section 4.2.

2.3 The van der Waals EOS

The van der Waals thermal EOS is given by,

$$p(\tau, T) := \frac{RT}{\tau - b} - \frac{a}{\tau^2}, \quad (2.13)$$

where b represents the maximum compression of the fluid; that is, $\tau > b$, a is a material dependent constant describing the attractive forces of fluid and should not be confused with the speed of sound, a , and R is the universal gas constant. This EOS was first derived by van der Waals in [41]. In order to find the complete EOS, more information must be provided. In particular, we suggest that the temperature be defined by

$$T = \frac{\gamma - 1}{R} \left(e + \frac{a}{\tau} \right). \quad (2.14)$$

The reasoning behind this follows from Maxwell's relation (2.3). The details can be found in [36, Sec. 3.5]. From this definition, we find the caloric EOS to be,

$$p(\tau, e) = (\gamma - 1) \frac{e + a/\tau}{\tau - b} - \frac{a}{\tau^2} \quad (2.15)$$

One also sees that, under the above definitions of T and p , the specific entropy,

$$s(\tau, e) = R \log \left(\left(e + \frac{a}{\tau} \right)^{\frac{1}{\gamma-1}} (\tau - b) \right) - s_0, \quad s_0 \in \mathbb{R}, \quad (2.16)$$

satisfies the second law of thermodynamics, (2.2). Using (2.15) and (2.16) we can eliminate e and find the complete EOS for p ,

$$p(\tau, s) = \frac{\gamma - 1}{(\tau - b)^\gamma} \exp \left(\frac{\gamma - 1}{R} (s - s_0) \right) - \frac{a}{\tau^2}. \quad (2.17)$$

Recall the sound speed is,

$$a^2 = \frac{(\gamma - 1)\gamma\tau^2}{(\tau - b)^{\gamma+1}} \exp \left(\frac{\gamma - 1}{R} (s - s_0) \right) - \frac{2a}{\tau}. \quad (2.18)$$

Using (2.17) we can compute the sound speed in terms of p and ρ .

$$a^2 = \gamma \frac{p + a\rho^2}{\rho(1 - b\rho)} - 2a\rho. \quad (2.19)$$

2.4 The Cubic EOS

There is a general class of EOS referred to as the cubic EOS as the pressure is a cubic function of \sqrt{T} . A general form of this EOS is can be defined as,

$$p(\tau, T) = \frac{RT}{\tau - b} - \frac{\alpha(T)}{(\tau - br_1)(\tau - br_2)}, \quad (2.20)$$

where $\alpha(T)$ is referred to as the attractive term. We refer to [42] for a survey of this general class of EOS. Note that the van der Waals EOS for the pressure is recovered if one sets $\alpha(T) = a$ and $r_1 = r_2 = 0$. Next, we suppose that the specific internal energy is given by,

$$e(\tau, T) = c_v T + \frac{\alpha(T) - T\alpha'(T)}{b(r_1 - r_2)} \log \left(\frac{\tau - br_1}{\tau - br_2} \right) + e_0, \quad (2.21)$$

where e_0 is a reference specific internal energy and $r_1 \neq r_2$. This assumption follows from the relation, $de = c_v dT + (T(\frac{\partial p}{\partial T})_\tau - p)d\tau$.

2.4.1 The Redlich-Kwong EOS

The Redlich-Kwong EOS was first introduced in by Redlich and Kwong in [43]. It can be obtained from the general cubic EOS by setting $\alpha(T) = a/\sqrt{T}$, $r_1 = 0$, and $r_2 = -1$. The parameters a and b are the same parameters from the van der Waals EOS. The equation of state is thus,

$$p(\tau, T) = \frac{RT}{\tau - b} - \frac{a}{\sqrt{T}\tau(\tau + b)}, \quad (2.22a)$$

$$e(\tau, T) = c_v T + \frac{3a}{2b\sqrt{T}} \log\left(\frac{\tau}{\tau + b}\right) + e_0. \quad (2.22b)$$

Note that an explicit expression for $p = p(\tau, e)$ cannot be written. Instead, one must first solve a cubic equation and then use that solution in $p(\tau, T)$ or $e(\tau, T)$ depending on the desired quantity. The two cubic equations to solve for temperature, given p or e , respectively, are,

$$\tilde{T}^3 - \frac{p(\tau - b)}{R}\tilde{T} - \frac{a(\tau - b)}{\tau(\tau + b)} = 0, \quad (2.23a)$$

$$\tilde{T}^3 - \frac{e - e_0}{c_v}\tilde{T} + \frac{3a}{2bc_v} \log\left(\frac{\tau}{\tau + b}\right) = 0, \quad (2.23b)$$

for $\tilde{T} := \sqrt{T}$. We can also see that the specific entropy defined by,

$$s(\tau, T) := c_v \log(T) + \frac{a}{2bT^{3/2}} \log\left(\frac{\tau}{\tau + b}\right) + R \log(\tau - b), \quad (2.24)$$

satisfies the 2nd law of thermodynamics, (2.2).

Remark 2.4.1 (Solving the Cubic). For completion, we briefly describe how to compute the roots for a cubic equation of the form, $x^3 + c_1x + c_0 = 0$. Define $\Delta := \frac{c_1^3}{27} + \frac{c_0^2}{4}$. If $\Delta \geq 0$,

then this cubic equation has one real root,

$$x_1 = \sqrt[3]{-\frac{c_0}{2} + \sqrt{\Delta}} + \sqrt[3]{-\frac{c_0}{2} - \sqrt{\Delta}}. \quad (2.25)$$

If $\Delta < 0$, then there are three real roots and they are defined by,

$$x_k := 2\sqrt{-\frac{c_1}{3}} \cos \left[\frac{1}{3} \arccos \left(\frac{3c_0}{2c_1} \sqrt{\frac{-3}{c_1}} \right) - \frac{2\pi(k-1)}{3} \right], \quad \text{for } k = 1, 2, 3. \quad (2.26)$$

Note, in both cases, (2.23a)–(2.23b), for $T = 0$, the cubic equation is negative, and since the coefficient of \tilde{T}^3 is positive, we must have at least one positive root. In the event that there are multiple positive roots, then we simply take the largest root. \square

Equations of State in $p = p(\rho, e)$ Form	
Ideal	$p = (\gamma - 1)\rho e \quad (2.27)$
Covolume	$p = (\gamma - 1) \frac{\rho e}{1 - b\rho} \quad (2.28)$
van der Waals	$p = (\gamma - 1) \frac{\rho e + a\rho^2}{1 - b\rho} - a\rho^2 \quad (2.29)$
Redlich-Kwong	No explicit formula
Jones-Wilkins-Lee	$p = A \left(1 - \frac{\omega}{R_1} \frac{\rho}{\rho_0} \right) e^{-R_1 \frac{\rho_0}{\rho}} + B \left(1 - \frac{\omega}{R_2} \frac{\rho}{\rho_0} \right) e^{-R_2 \frac{\rho_0}{\rho}} + \omega \rho e \quad (2.30)$

Table 2.1: A quick reference for common equations of state in $p = p(\rho, e)$ form. More information on these EOS can be found in Chapter 2.

Sound Speeds for Different EOS in $a = a(\rho, p)$ Form	
Ideal	$a^2 = \frac{\gamma p}{\rho} \quad (2.31)$
Covolume	$a^2 = \frac{\gamma p}{\rho(1 - b\rho)} \quad (2.32)$
van der Waals	$a^2 = \gamma \frac{p + a\rho^2}{\rho(1 - b\rho)} - 2a\rho \quad (2.33)$
Redlich-Kwong	No explicit formula
Jones-Wilkins-Lee	$a^2 = \frac{(\omega + 1)p}{\rho} + A \left(\frac{R_1 \rho_0}{\rho^2} + \frac{\omega(\omega + 1 + \rho)}{R_1 \rho_0} - 1 \right) \exp \left(- R_1 \frac{\rho_0}{\rho} \right) + B \left(\frac{R_2 \rho_0}{\rho^2} + \frac{\omega(\omega + 1 + \rho)}{R_2 \rho_0} - 1 \right) \exp \left(- R_2 \frac{\rho_0}{\rho} \right) \quad (2.34)$

Table 2.2: A quick reference for common sound speeds for several EOS $a = a(\rho, p)$ form. More information on these EOS can be found in Chapter 2.

2.5 The Jones-Wilkins-Lee EOS

The Jones-Wilkins-Lee (JWL) EOS is an empirical EOS used to model detonation products in multi-material reactive flow. The ideal gas law is often used to model the surrounding ambient fluid. As the purpose of this thesis is only for single material compressible Euler equations, we use the JWL EOS as a demonstration of the numerical method. The original paper for the development of this EOS can be found in [44].

The pressure is defined by,

$$p(\rho, e) := A \left(1 - \frac{\omega}{R_1} \frac{\rho}{\rho_0} \right) \exp \left(- R_1 \frac{\rho_0}{\rho} \right) + B \left(1 - \frac{\omega}{R_2} \frac{\rho}{\rho_0} \right) \exp \left(- R_2 \frac{\rho_0}{\rho} \right) + \omega \rho (e - e_0), \quad (2.35)$$

where A , B , R_1 and R_2 are parameters specific to the model or experiment, ρ_0 and e_0 are

some reference density and specific internal energy, and ω is also a chosen constant depending on the material, and is related to the Grüneisen coefficient.

Remark 2.5.1 (JWL Misinterpretation). There is a lot of confusion surrounding the exact definition of the JWL EOS. In some cases in the literature, the JWL EOS is taken as an isentrope and not an EOS by replacing $e - e_0$ with e_0 . The technical report, [45], provides an overview of this confusion. \square

2.6 The Mie-Gruneisen EOS

The Mie-Gruneisen EOS is an incomplete EOS defined by,

$$p(\rho, e) := p_{\text{ref}}(\rho) + \rho\Gamma(\rho)(e - e_{\text{ref}}(\rho)) \quad (2.36)$$

where p_{ref} and e_{ref} are some reference pressure and specific internal energy curves, respectively. For example the reference curve could be shock locus, an isotherm, an isentrope and so on. A quick overview of this incomplete EOS as well as a means to *complete* it, is given in [46].

For our numerical demonstrations we use the linear Hugoniot locus, defined by,

$$p_{\text{ref}}(\rho) := P_0 + \rho_0 c_0^2 \frac{1 - \frac{\rho_0}{\rho}}{\left(1 - s\left(1 - \frac{\rho_0}{\rho}\right)\right)^2}, \quad (2.37a)$$

$$e_{\text{ref}}(\rho) := e_0 + \frac{P_0 + p_{\text{ref}}(\rho)}{2\rho_0} \left(1 - \frac{\rho_0}{\rho}\right), \quad (2.37b)$$

where $s > 1$ and ρ_0 , c_0 , e_0 , and P_0 are reference density, sound speed, specific internal energy, and pressure, respectively. This particular equation of state is used for modeling solids under high pressures. More details on this particular equation of state can be found in [47, Sec. 4.4].

3. FINITE ELEMENT APPROXIMATION OF THE EULER EQUATIONS

The numerical method that we outline in this thesis is actually discretization independent. However, our numerical simulations are computed using the finite element method; in particular, \mathbb{P}_1 and \mathbb{Q}_1 continuous finite elements. There is a wealth of literature on the finite element method. Some recommended resources for this topic are: Grossman et. al. [48], Larsson & Thom e [49], Ciarlet [50], and the three volume series on the finite element method by Ern & Guermond [51], [52], and [53].

Therefore, in this chapter, we present some of the fundamentals regarding the continuous finite element method.

3.1 The Continuous Galerkin (cG) Framework

We use a continuous Galerkin method for solving the compressible Euler equations. We start by defining the geometric finite element (using the notation of Ciarlet, [50]) $(\widehat{K}_{\text{geo}}, \widehat{P}_{\text{geo}}, \widehat{\Sigma}_{\text{geo}})$. Here, \widehat{K}_{geo} is the reference element composed with the vertices $\{\widehat{\mathbf{a}}_i\}_{i \in \widehat{\mathcal{N}}_{\text{geo}}}$, \widehat{P}_{geo} is the geometric polynomial space used to construct the geometric mapping, and $\widehat{\Sigma}_{\text{geo}}$ are the nodal Lagrange degrees of freedom (dofs). Let $\{\widehat{\theta}_i\}_{i \in \widehat{\mathcal{N}}_{\text{geo}}}$ denote the collection of reference shape functions. That is, $\widehat{\sigma}_i(\widehat{\theta}_j) = \widehat{\theta}_j(\widehat{\mathbf{a}}_i) = \delta_{ij}$ for $i, j \in \widehat{\mathcal{N}}_{\text{geo}}$ and $\widehat{\sigma}_i \in \widehat{\Sigma}_{\text{geo}}$.

We use this reference geometric finite element to construct a collection of mappings, of which, defines our mesh. Let $\mathcal{T}_h = \{K_i\}_{i \in \mathcal{N}_{\text{shape}}}$ denote a sequence of shape regular non-overlapping elements which exactly covers our domain D . Then $\bigcup_{i \in \mathcal{N}_{\text{shape}}} K_i = D$. Let (K, P, Σ) be a local finite element where K is polytope with vertices $\{\mathbf{a}_i\}_{i \in \mathcal{N}_{\text{geo}}}$. Then the affine geometric mapping, $\mathbf{T}_K : \widehat{K}_{\text{geo}} \rightarrow K$ is defined by,

$$\mathbf{T}_K(\widehat{\mathbf{x}}) = \sum_{i \in \mathcal{N}} \widehat{\theta}_i(\widehat{\mathbf{x}}) \mathbf{a}_i, \quad \text{for } \widehat{\mathbf{x}} \in \widehat{K}_{\text{geo}}. \quad (3.1)$$

Therefore, if our domain, Ω , is polygonal, we can exactly triangulate the domain. That is, $\Omega = \bigcup_{K \in \mathcal{T}_h} K$.

However, since we will be using the first order Lagrange finite element for the approximation space, this means that the geometric and the approximation finite elements will coincide. Nevertheless, we will still make the distinction between the two elements. For the numerical method described in Chapters 4 and 5 we use the following finite element approximation space for the discretization of the PDE,

$$P(\mathcal{T}_h) := \{v \in C^0(\Omega) : v \circ \mathbf{T}_K \in \widehat{P}, \forall K \in \mathcal{T}_h\}, \quad (3.2)$$

where \mathcal{T}_h consists of either simplices or quadrangles, $\widehat{P} = \mathbb{P}_1(\widehat{K})$ or $\widehat{P} = \mathbb{Q}_1(\widehat{K})$ (that is, first order multivariate polynomials defined on the reference element \widehat{K}), and \mathbf{T}_K is defined in (3.1). Let $\{\mathbf{x}_i\}_{i \in \mathcal{V}}$ be the collection of nodes of our mesh \mathcal{T}_h , where \mathcal{V} is the index set for the nodes. A basis for $P(\mathcal{T}_h)$ is given by $\text{span}\{\varphi_i(\mathbf{x})\}_{i \in \mathcal{V}}$, where

$$\varphi_i(\mathbf{x}) = \begin{cases} \widehat{\theta}_{\text{j_dof}^{-1}(i)}(\mathbf{T}_K^{-1}(\mathbf{x})), & \text{if } \mathbf{x} \in K \in \mathcal{T}_h(i) \\ 0, & \text{otherwise} \end{cases} \quad (3.3)$$

where $\mathcal{T}_h(i) \subset \mathcal{T}_h$ such that $K \in \mathcal{T}_h$ contains the vertex $\{\mathbf{x}_i\}$ and $\text{j_dof} : \mathcal{T}_h \times \widehat{\mathcal{N}}_{\text{geo}} \rightarrow \mathcal{V}$ and is defined so that j_dof identifies the corresponding global degree of freedom from a cell $K \in \mathcal{T}_h$ and local node $j \in \widehat{\mathcal{N}}_{\text{geo}}$ from the reference element. Note that $\text{j_dof}^{-1}(i)$ is well defined since the element K is fixed. The φ_i are often referred to as *tent functions* based on their shape in two dimensions, see Figure 3.1. These tent functions share the nice property that $\varphi_i(\mathbf{x}) \geq 0$ for all $i \in \mathcal{V}$ and that they form a *partition of unity*; that is, $\sum_{i \in \mathcal{V}} \varphi_i(\mathbf{x}) = 1$.

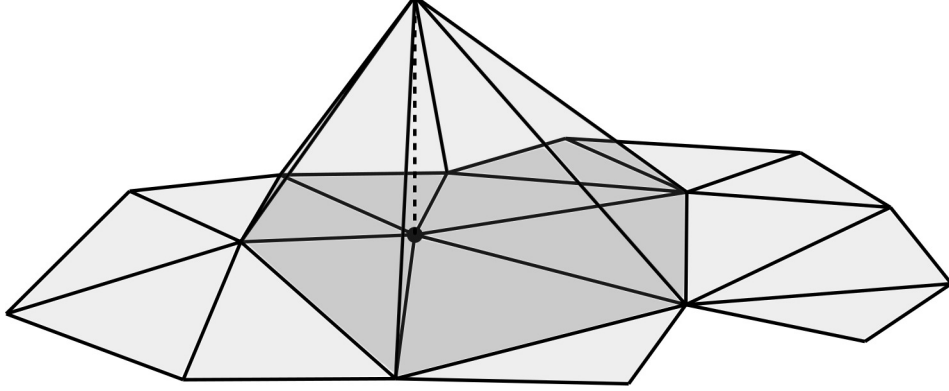


Figure 3.1: A typical \mathbb{P}_1 basis function.

3.1.1 Semi-discrete Scheme

Numerically solving the Euler equations is done through the conserved variables: ρ , \mathbf{m} , and E . The semi-discretization of these variables is given as,

$$\rho_h(\mathbf{x}, t) := \sum_{i \in \mathcal{V}} \rho_i(t) \varphi_i(\mathbf{x}), \quad \mathbf{m}_h(\mathbf{x}, t) = \sum_{i \in \mathcal{V}} \mathbf{M}_i(t) \varphi_i(\mathbf{x}), \quad E_h(\mathbf{x}, t) = \sum_{i \in \mathcal{V}} E_i(t) \varphi_i(\mathbf{x}). \quad (3.4)$$

The unknowns are the $\rho_i(t)$, $\mathbf{M}_i(t)$, and $E_i(t)$, and the global shape functions $\varphi_i(\mathbf{x})$ are the usual tent functions which form a basis for our finite dimensional space $P(\mathcal{J}_h)$. This approximation can be written in the more compact form,

$$\mathbf{u}_h(\mathbf{x}, t) := \sum_{i \in \mathcal{V}} \mathbf{u}_i(t) \varphi_i(\mathbf{x}), \quad (3.5)$$

where $\mathbf{u}_i(t) = (\rho_i(t), \mathbf{M}_i(t), E_i(t))^T$. Thus, our problem can be written as,

$$\partial_t \mathbf{u}_h + \nabla \cdot \mathbf{f}(\mathbf{u}_h) = \mathbf{0}. \quad (3.6)$$

3.2 The Fully Discrete Scheme

For computer implementation, it is not feasible to solve the semi-discrete problem. So we further discretize in time. We approximate the time derivative with the forward Euler

method.

Remark 3.2.1 (Time Stepping Methods). In the numerical results the actual time stepping is performed using the strong stability preserving 3rd order Runge-Kutta method (SSP RK3). However, a much more efficient time stepping method has been introduced in Ern & Guermond [54]. It uses an explicit Runge-Kutta method which applies a nonlinear limiting process on the high order update at each stage of the method. This allows for a less restrictive time step while still being conservative and invariant-domain preserving (see Definition 3.2.1).

□

Let $t^{n+1} := t^n + \Delta t$ where t^n is the time at the n th time step and Δt is the time step. Let $\mathbf{u}_h^n(\mathbf{x}) := \mathbf{u}_h(\mathbf{x}, t^n)$ be the approximate solution at time t^n and $\mathbf{U}_i^n := \mathbf{u}_i(t^n)$ be the coefficient of the i th basis function $\varphi_i(\mathbf{x})$ at time t^n . That is, $\mathbf{U}_i^n = (\rho_i(t^n), \mathbf{M}_i(t^n), \mathbf{E}_i(t^n))^\top =: (\rho_i^n, \mathbf{M}_i^n, \mathbf{E}_i^n)^\top$, hence

$$\mathbf{u}_h^n(\mathbf{x}) = \sum_{i \in \mathcal{V}} \mathbf{U}_i^n \varphi_i(\mathbf{x}). \quad (3.7)$$

As will be necessary later on, we approximate the flux, $\mathbf{f}(\mathbf{u}_h^n)$ by projecting it onto the discrete finite element space. We define this projection $\Pi_h : C^0(\mathbb{R}^{d+2}; \mathbb{R}^{(d+2) \times d}) \rightarrow [P(\mathcal{T}_h)]^{(d+2) \times d}$ by,

$$\mathbf{f}(\mathbf{u}_h^n) \approx \Pi_h \mathbf{f}(\mathbf{u}_h^n) = \sum_{i \in \mathcal{V}} \mathbf{f}(\mathbf{U}_i^n) \varphi_i(\mathbf{x}). \quad (3.8)$$

Putting this altogether with the forward Euler method, we arrive at the following numerical method: find \mathbf{u}_h^{n+1} such that,

$$\frac{\mathbf{u}_h^{n+1} - \mathbf{u}_h^n}{\Delta t} + \nabla \cdot (\Pi_h \mathbf{f}(\mathbf{u}_h^n)) = \mathbf{0}. \quad (3.9)$$

Solving this equation is done in the weak sense by testing the equation with φ_i for all $i \in \mathcal{V}$. That is, multiply equation (3.9) by φ_i and integrate over D . Doing so gives the

following system of card(\mathcal{V}) equations,

$$\frac{1}{\Delta t_n} \sum_{j \in \mathcal{G}(i)} (\mathbf{U}_j^{n+1} - \mathbf{U}_j^n) m_{ij} + \sum_{j \in \mathcal{G}(i)} \mathbf{f}(\mathbf{U}_j^n) \mathbf{c}_{ij} = \mathbf{0}, \quad \text{for all } i \in \mathcal{V}, \quad (3.10)$$

where

$$m_{ij} := \int_D \varphi_i(\mathbf{x}) \varphi_j(\mathbf{x}) d\mathbf{x}, \quad (3.11)$$

$$\mathbf{c}_{ij} := \int_D \varphi_i(\mathbf{x}) \nabla \varphi_j(\mathbf{x}) d\mathbf{x}. \quad (3.12)$$

Equation (3.10) is the so-called **Galerkin method**, see Ern & Guermond [52, Sec. 26.1]. Unfortunately, this method can be unstable if any discontinuities develop in the solution, a method for resolving this issue is to introduce an artificial viscosity, see von Neumann & Richtmyer [3].

Remark 3.2.2 (Partition of Unity). Recall that the basis functions $\{\varphi_i\}_{i \in \mathcal{V}}$ form a partition of unity hence $\sum_{j \in \mathcal{V}} m_{ij} = \sum_{j \in \mathcal{G}(i)} m_{ij} = \int_\Omega \varphi_i(\mathbf{x}) d\mathbf{x}$. We refer to $m_i := \int_\Omega \varphi_i(\mathbf{x}) d\mathbf{x}$ as the *lumped mass*. Furthermore, from the partition of unity, we have that $\sum_{j \in \mathcal{V}} \mathbf{c}_{ij} = \sum_{j \in \mathcal{G}(i)} \mathbf{c}_{ij} = \mathbf{0}$. \square

Remark 3.2.3 (Stability of the Galerkin Approximation). It is well known that the Galerkin approximation is stable as long as $\mathbf{u}_0(\mathbf{x})$ is smooth and the solution $\mathbf{u}(\mathbf{x}, t)$ remains smooth up to some final time, t_{final} . Furthermore, for \mathbb{P}_1 or \mathbb{Q}_1 continuous finite elements, we can achieve second order convergence. However, if the solution develops a discontinuity, then the approximate solution produces wild non-physical oscillations. \square

The method by which we resolve this issue is to supply a *graph viscosity*. That is, we modify the scheme, (3.10), to be,

$$\frac{1}{\Delta t_n} \sum_{j \in \mathcal{G}(i)} m_{ij} (\mathbf{U}_i^{n+1} - \mathbf{U}_i^n) + \sum_{j \in \mathcal{G}(i)} \mathbf{f}(\mathbf{U}_j^n) \mathbf{c}_{ij} - \sum_{j \in \mathcal{G}(i)} d_{ij}^n (\mathbf{U}_j^n - \mathbf{U}_i^n) = \mathbf{0}, \quad (3.13)$$

for all $i \in \mathcal{V}$. We also assume that d_{ij}^n satisfies the following properties,

$$d_{ij}^n \geq 0 \quad \text{for } i \neq j, \quad d_{ij}^n = d_{ji}^n, \quad \text{and } d_{ii} := - \sum_{j \in \mathcal{G}(i)} d_{ij}^n. \quad (3.14)$$

This specific discrete scheme with the graph viscosity is first introduced by Guermond & Popov in [16, Sec. 3.2]. We now introduce the definition of an invariant-domain preserving method.

Definition 3.2.1 (Invariant-Domain Preserving Method). Let \mathcal{B} be a convex invariant set. For $\mathbf{U}_i^0 \in \mathcal{B}$ for all $i \in \mathcal{V}$, if the updated states satisfy $\mathbf{U}_i^n \in \mathcal{B}$ for all $i \in \mathcal{V}$ and $n \in \mathbb{N}$, then we say the numerical method which provides the states $\{\mathbf{U}_i^n\}_{i \in \mathcal{V}}$ is said to be *invariant-domain preserving*. □

4. THE FIRST ORDER APPROXIMATION*

In this chapter, we outline the first order method and prove that it is invariant domain preserving for an arbitrary equation of state. We first approximate the mass matrix, $\{m_{ij}\}_{i,j \in \mathcal{V}}$, by the lumped mass matrix, $\{m_i\}_{i \in \mathcal{V}}$. We can rewrite, (3.13), as an explicit equation,

$$\mathbf{U}_i^{\text{L},n+1} = \mathbf{U}_i^n - \frac{\Delta t}{m_i} \left(\sum_{j \in \mathcal{G}(i)} \mathbb{f}(\mathbf{U}_j^n) \mathbf{c}_{ij} - \sum_{j \in \mathcal{G}(i)} d_{ij}^m (\mathbf{U}_j^n - \mathbf{U}_i^n) \right) \quad \text{for all } i \in \mathcal{V}. \quad (4.1)$$

Note we also use ^L superscript to signify that this is the low order update. This will be necessary for distinguishing between the high order update in the chapters to come.

Remark 4.0.1 (Discretization Independent Method). Note that the numerical method described in (4.1) can be made discretization independent. The quantities m_i and \mathbf{c}_{ij} will vary depending on the specific discretization. More on this can be found in Guermond et al. [55]. □

4.1 Invariant Domain Preserving Method

With a bit of rearrangement and using the fact that $\sum_{j \in \mathcal{G}(i)} \mathbf{c}_{ij} = 0$, we rewrite (4.1) as a convex combination of states under the CFL condition, (4.7),

$$\mathbf{U}_i^{\text{L},n+1} = \left(1 - \sum_{j \in \mathcal{G}(i) \setminus \{i\}} \frac{2\Delta t_n d_{ij}^{\text{L},n}}{m_i} \right) \mathbf{U}_i^n + \sum_{j \in \mathcal{G}(i) \setminus \{i\}} \frac{2\Delta t_n d_{ij}^{\text{L},n}}{m_i} \bar{\mathbf{U}}_{ij}^n \left(\frac{\|\mathbf{c}_{ij}\|_{\ell^2}}{2d_{ij}^{\text{L},n}} \right), \quad (4.2)$$

where

$$\bar{\mathbf{U}}_{ij}^n(t) = \frac{1}{2}(\mathbf{U}_i^n + \mathbf{U}_j^n) - t(\mathbb{f}(\mathbf{U}_j^n) - \mathbb{f}(\mathbf{U}_i^n)) \frac{\mathbf{c}_{ij}}{\|\mathbf{c}_{ij}\|_{\ell^2}}. \quad (4.3)$$

Throughout this thesis, we refer to $\bar{\mathbf{U}}_{ij}^n$ as the **bar states**.

* A majority of this chapter is a modification of the work done in [1] and is reprinted with permission from [1].

The bar states are a fundamental ingredient in proving that the numerical method is invariant domain preserving (see Definition 3.2.1). Since for a large enough $d_{ij}^{L,n}$, we claim that $\bar{\mathbf{U}}_{ij}^n \in \mathcal{B}(b)$. We define $d_{ij}^{L,n}$ to be,

$$d_{ij}^{L,n} := \max\{\lambda_{\max}(\mathbf{U}_i^n, \mathbf{U}_j^n, \mathbf{n}_{ij})\|\mathbf{c}_{ij}\|_{\ell^2}, \lambda_{\max}(\mathbf{U}_j^n, \mathbf{U}_i^n, \mathbf{n}_{ji})\|\mathbf{c}_{ji}\|_{\ell^2}\}, \quad (4.4)$$

for $i \neq j$, where $\mathbf{n}_{ij} := \mathbf{c}_{ij}/\|\mathbf{c}_{ij}\|_{\ell^2}$ and $\lambda_{\max}(\mathbf{U}_i^n, \mathbf{U}_j^n, \mathbf{n}_{ij})$ is the max wave speed to the local Riemann problem,

$$\partial_t \mathbf{u} + \partial_x(\mathbb{f}(\mathbf{u})\mathbf{n}_{ij}) = \mathbf{0}, \quad \mathbf{u}_0(x) = \begin{cases} \mathbf{U}_i^n, & \text{if } x < 0, \\ \mathbf{U}_j^n, & \text{if } x > 0. \end{cases} \quad (4.5)$$

Notice that the bar states, (4.3), are the average of the solution to the Riemann problem, (4.5), at the ‘‘artificial’’ time $t = \frac{\|\mathbf{c}_{ij}\|_{\ell^2}}{2d_{ij}^{L,n}} \leq \frac{1}{2\lambda_{\max}}$; see Theorem 1.3.1.

Remark 4.1.1 (Choice of Larger Artificial Viscosity). Notice that, taking larger values of $d_{ij}^{L,n}$; that is, finding an upper estimate on the maximum wave speed, still preserves the desired structure of the bar states. Since if $\widehat{\lambda}_{\max}$ is an upper bound on λ_{\max} and $\widehat{d}_{ij}^{L,n}$ is simply $d_{ij}^{L,n}$ with λ_{\max} replaced with $\widehat{\lambda}_{\max}$, then

$$t = \frac{\|\mathbf{c}_{ij}\|_{\ell^2}}{2\widehat{d}_{ij}^{L,n}} \leq \frac{1}{2\widehat{\lambda}_{\max}} \leq \frac{1}{2\lambda_{\max}} \quad (4.6)$$

Hence the bar states, (4.3), are still the averages of the Riemann solution for (4.5) at a different ‘‘fake’’ time. \square

Remark 4.1.2 (Justification for the Flux Approximation). Recall that the approximation of $\mathbb{f}(\mathbf{u}_h^n)$ is $\sum_{i \in \mathcal{V}} \mathbb{f}(\mathbf{U}_i^n)\varphi_i(\mathbf{x})$ which was defined in (3.8). This approximation was necessary to justify the bar states as being the average of the solution to the Riemann problem (4.5). \square

The challenge is now to determine λ_{\max} for an arbitrary EOS or a sufficiently close upper bound. This is not a simple task and will be investigated in Chapter 4. Note that, once $d_{ij}^{L,n}$

has been determined, this gives us the following CFL constraint which guarantees a convex combination in of states in (4.2),

$$\Delta t \leq \frac{m_i}{2 \sum_{j \in \mathcal{G}(i) \setminus \{i\}} d_{ij}^{\mathbf{L},n}} = -\frac{m_i}{2d_{ii}^{\mathbf{L},n}}, \quad \text{for all } i \in \mathcal{V}. \quad (4.7)$$

Theorem 4.1.1 (Invariant-Domain Preservation). *[16]] If $d_{ij}^{\mathbf{L},n}$ is defined by (4.4) and $\mathbf{U}_i^n \in \mathcal{B}(b)$ for all $i \in \mathcal{V}$, then the update provided by (4.1) is invariant domain preserving under the CFL condition, $\Delta t \leq \min_{i \in \mathcal{V}} (-\frac{m_i}{2d_{ii}^{\mathbf{L},n}})$.*

Proof. From the choice of $d_{ij}^{\mathbf{L},n}$ we have from Theorem 1.3.1 that $\bar{\mathbf{U}}_{ij}^n \in \mathcal{B}(b)$ for each $j \in \mathcal{G}(i) \setminus \{i\}$ and from the assumption, $\mathbf{U}_i^n \in \mathcal{B}(b)$. Thus the update provided in (4.2) is a convex combination of states in the convex set $\mathcal{B}(b)$ (under the CFL condition (4.7)). Hence $\mathbf{U}_i^{n+1} \in \mathcal{B}(b)$. □

4.2 Extended Riemann Problem

If the pressure is given by the ideal gas law, then the solution to the Riemann problem an exact self-similar weak solution can be determined. This was originally done by Lax in [34]; for other useful resources regarding the exact solution to this Riemann problem, see [56, Chapter 4], [31, Chapter II, Section 3], and [57, Section 5.6]. However, for an arbitrary or tabulated equations of state, computing the max wave speed, λ_{\max} , is extremely difficult if not impossible. In order to solve this issue, we propose extending the Riemann problem with an auxiliary equation in terms of a new variable, Γ . This idea is motivated by the paper by Abgrall & Karni [58] in the context of multi-fluids. We also change the EOS given by the oracle to an equation of state based on the Noble-Abel Stiffened-Gas (NASG) EOS which interpolates the left and right pressure, p_L and p_R , respectively. Recall the NASG EOS is defined by $p(\rho, e) := (\gamma - 1) \frac{\rho(e-q)}{1-b\rho} + \gamma p_\infty$. Alternatively, we write this equation of state as, $p + p_\infty = (\gamma - 1) (\frac{\rho(e-q)}{1-b\rho} - p_\infty)$. This method is first introduced in [2] but uses the covolume EOS instead of the NASG EOS.

Definition 4.2.1. The extended Riemann problem is, $\partial_t \tilde{\mathbf{u}} + \partial_x(\tilde{\mathbf{f}}(\tilde{\mathbf{u}})\mathbf{n}) = \mathbf{0}$, where

$$\tilde{\mathbf{u}} := (\mathbf{u}, \Gamma)^\top = \begin{pmatrix} \rho \\ m \\ \mathbf{m}^\perp \\ \mathcal{E} \\ \Gamma \end{pmatrix}, \quad \tilde{\mathbf{f}}(\tilde{\mathbf{u}})\mathbf{n} := \begin{pmatrix} m \\ \frac{m^2}{\rho} + \tilde{p}_{\text{nasg}}(\tilde{\mathbf{u}}) \\ \frac{m}{\rho}\mathbf{m}^\perp \\ \frac{m}{\rho}(\mathcal{E} + \tilde{p}_{\text{nasg}}(\tilde{\mathbf{u}})) \\ \frac{m}{\rho}\Gamma \end{pmatrix}, \quad (4.8)$$

with left and right data $\tilde{\mathbf{U}}_Z := (\mathbf{U}_Z, \Gamma_Z) = (\rho_Z, \mathbf{m}_Z \cdot \mathbf{n}, \mathbf{m}_Z^\perp, \mathcal{E}_Z, \Gamma_Z)^\top$ with

$$\Gamma_Z := \rho_Z \left(\frac{p_Z + p_\infty}{\frac{\rho_Z(e(\mathbf{U}_Z) - q)}{1 - b\rho_Z} - p_\infty} + 1 \right) \quad (4.9)$$

where $Z \in \{L, R\}$ and the pressure is defined by

$$\tilde{p}_{\text{nasg}}(\tilde{\mathbf{u}}) := \left(\frac{\Gamma}{\rho} - 1 \right) \left(\frac{\mathcal{E} - \frac{1}{2}m^2/\rho - \rho q}{1 - b\rho} - p_\infty \right) + \frac{\Gamma}{\rho} p_\infty. \quad (4.10)$$

□

It is often easier to work with the primitive variables, so let $\gamma := \Gamma/\rho$ then \tilde{p}_{nasg} can be rewritten as $p_{\text{nasg}}(\rho, e, \gamma) := \tilde{p}_{\text{nasg}}(\tilde{\mathbf{u}}) = (\gamma - 1) \left(\frac{\rho(e - q)}{1 - b\rho} - p_\infty \right) - p_\infty$ which is simply the NASG EOS with a variable γ . The left and right states are then $(\rho_Z, \mathbf{v}_Z \cdot \mathbf{n}, \mathbf{v}_Z^\perp, p_Z, \gamma_Z)$ where $\gamma_Z := \frac{(p_Z + p_\infty)(1 - b\rho_Z)}{\rho_Z(e(\rho_Z, p_Z) - q) - p_\infty(1 - b\rho_Z)} + 1$ and $\mathbf{v}_Z^\perp := \rho_Z^{-1} \mathbf{m}_Z^\perp$ for $Z \in \{L, R\}$.

Remark 4.2.1 (Interpolating with p_{nasg}). The reason for this choice of γ_Z is because $p_{\text{nasg}}(\rho_Z, e(\rho_Z, p_Z), \gamma_Z) = p_Z$ for $Z \in \{L, R\}$. That is, p_{nasg} interpolates the left and right pressures. □

Definition 4.2.2 (Extended Invariant Domain). We define the extended invariant domain as,

$$\tilde{\mathcal{B}}(b, q, p_\infty) := \{\tilde{\mathbf{u}} \in \mathbb{R}^{d+3} : \mathbf{u} \in \mathcal{B}(b, q, p_\infty), \Gamma > \rho\}. \quad (4.11)$$

to account for the new variable Γ . □

Remark 4.2.2 (Extended Bar States). Notice that $\tilde{\mathbf{f}}(\tilde{\mathbf{u}}_Z) = (\mathbf{f}(\mathbf{u}_Z), \mathbf{v}_Z \Gamma_Z)^\top$ since $\tilde{\rho}_{\text{nasg}}(\tilde{\mathbf{u}}_Z) = p_Z = p(\mathbf{u}_Z)$. Let $\bar{\mathbf{u}}_{\text{LR}} := \bar{\mathbf{U}}_{ij}^n$. Then, the bar state for the extended Riemann problem is,

$$\bar{\tilde{\mathbf{u}}}_{\text{LR}} = \begin{pmatrix} \bar{\mathbf{u}}_{\text{LR}} \\ \frac{1}{2}(\Gamma_L + \Gamma_R) - \frac{1}{2\lambda}(\mathbf{v}_R \Gamma_R - \mathbf{v}_L \Gamma_L) \cdot \mathbf{n} \end{pmatrix}, \quad (4.12)$$

where $\lambda = -\frac{m_i}{2d_{ii}^{\perp n}}$. Its important to note that the density, momentum, and total energy of the state $\bar{\tilde{\mathbf{u}}}_{\text{LR}}$ is the same as the state $\bar{\mathbf{u}}_{\text{LR}}$. Thus, if we can prove positivity of the density and specific internal energy of the extended Riemann problem, then this immediately carries over to the original Riemann problem. This is remark is critical in the justification of the invariant-domain preserving method. □

4.2.1 The Wave Structure

We begin by first deriving the wave structure of this Riemann problem. This is done by computing the Jacobian matrix of $\tilde{\mathbf{f}}(\tilde{\mathbf{u}})\mathbf{n}$. However, the computation is simpler if we make a change of variables. Let $\theta : \tilde{\mathcal{B}} \subset \mathbb{R}^{d+3} \rightarrow \tilde{\mathcal{B}}$ be a smooth diffeomorphism, defined by

$$\theta(\tilde{\mathbf{u}}) = \left(\rho, \frac{m}{\rho}, \frac{\mathbf{m}^\perp}{\rho}, \mathbf{e}(\mathbf{u}), \frac{\Gamma}{\rho} \right)^\top =: \tilde{\mathbf{w}}. \quad (4.13)$$

That is, θ maps to the primitive variables, $\theta(\tilde{\mathbf{u}}) = \tilde{\mathbf{w}} = (\rho, v, \mathbf{v}^\perp, e, \gamma)^\top$. Then the extended Riemann problem is formulated as,

$$\partial_t \rho + v \partial_x \rho + \rho \partial_x v = 0, \quad (4.14a)$$

$$\partial_t v + v \partial_x v + \rho^{-1} \partial_x p_{\text{nasg}} = 0, \quad (4.14b)$$

$$\partial_t \mathbf{v}^\perp + v \partial_x \mathbf{v}^\perp = 0, \quad (4.14c)$$

$$\partial_t e + \rho^{-1} p_{\text{nasg}} \partial_x v + v \partial_x e = 0, \quad (4.14d)$$

$$\partial_t \gamma + v \partial_x \gamma = 0, \quad (4.14e)$$

Thus, the Jacobian matrix of the system (4.14a)–(4.14e), is,

$$\mathbb{B}(\tilde{\mathbf{w}}) = \begin{pmatrix} v & \rho & \mathbf{0}^\top & 0 & 0 \\ \frac{1}{\rho} \frac{\partial p_{\text{nasg}}}{\partial \rho} & v & \mathbf{0}^\top & \frac{1}{\rho} \frac{\partial p_{\text{nasg}}}{\partial e} & \frac{1}{\rho} \frac{\partial p_{\text{nasg}}}{\partial \gamma} \\ \mathbf{0} & \mathbf{0} & v \mathbb{I}_{d-1} & \mathbf{0} & \mathbf{0} \\ 0 & \frac{p_{\text{nasg}}}{\rho} & \mathbf{0}^\top & v & 0 \\ 0 & 0 & \mathbf{0}^\top & 0 & v \end{pmatrix}. \quad (4.15)$$

The eigenvalues of this matrix are, $\mu_1(\tilde{\mathbf{w}}) = v - \tilde{\alpha}(\tilde{\mathbf{w}})$, $\mu_2(\tilde{\mathbf{w}}) = v$ (with multiplicity $d + 1$), and $\mu_3(\tilde{\mathbf{w}}) = v + \tilde{\alpha}(\tilde{\mathbf{w}})$, where

$$\tilde{\alpha}(\tilde{\mathbf{w}})^2 = \frac{\gamma(\tilde{p}_{\text{nasg}}(\theta^{-1}(\tilde{\mathbf{w}})) + p_\infty)}{\rho(1 - b\rho)} = \gamma(\gamma - 1) \left(\frac{e - q}{(1 - b\rho)^2} - \frac{p_\infty}{\rho(1 - b\rho)} \right). \quad (4.16)$$

Thus, the solution to the Riemann problem is composed of 3 waves. We denote them by the L-wave, C-wave, and R-wave.

Remark 4.2.3 (No Entropy for p_{nasg}). The interpolatory pressure, p_{nasg} , is not a real equation of state since it is also a function of the extra variable, γ . Therefore, we have not made a change of variable to the specific entropy which is usually done when computing the eigenvalues of the Jacobian matrix. However, in Section 4.3, the restriction of p_{nasg} to each wave, will define an equation of state. \square

Note that the change of variables does not affect the eigenvalues of the Jacobian matrix of the original system. That is, if $\mathbb{A}(\tilde{\mathbf{u}})$ is the Jacobian matrix for the flux, $\tilde{\mathbf{f}}(\tilde{\mathbf{u}})\mathbf{n}$, given in (4.8), then the eigenvalues of $\mathbb{A}(\tilde{\mathbf{u}})$ and $\mathbb{B}(\tilde{\mathbf{w}})$ are the same. Furthermore, let $\mathbf{r}_i(\tilde{\mathbf{u}})$ denote an eigenvector of $\mathbb{A}(\tilde{\mathbf{u}})$ with corresponding eigenvalue, $\lambda_i(\tilde{\mathbf{u}})$ and $\mathbf{s}_i(\tilde{\mathbf{w}})$ an eigenvector of $\mathbb{B}(\tilde{\mathbf{w}})$, with its corresponding eigenvalue, $\mu_i(\tilde{\mathbf{w}})$. Note also, eigenvectors are related by the following identity, $\mathbf{r}_i(\tilde{\mathbf{u}}) = (D_{\tilde{\mathbf{u}}}\theta(\tilde{\mathbf{u}}))^{-1}\mathbf{s}_i(\theta(\tilde{\mathbf{u}}))$. Then it can be shown that, $D\lambda_i(\tilde{\mathbf{u}}) \cdot \mathbf{r}_i(\tilde{\mathbf{u}}) = D\mu_i(\tilde{\mathbf{w}}) \cdot \mathbf{s}_i(\tilde{\mathbf{w}})$. The details on the change of variables regarding the Jacobian matrix can be found in [31, Section 2.1.1]

We now present the left and right eigenvectors for completeness. Note, that the associated eigenvectors are $\mathbf{s}_1(\tilde{\mathbf{w}}) = (1, -a/\rho, \mathbf{0}_{d-1}, 0, 0)^\top$, $\mathbf{s}_3(\tilde{\mathbf{w}}) = (1, a/\rho, \mathbf{0}_{d-1}, 0, 0)^\top$, and

$$\mathbf{s}_2(\tilde{\mathbf{w}}) \in \text{span} \left\{ \begin{pmatrix} 1 \\ 0 \\ \mathbf{0}_{d-1} \\ -\frac{e-q}{\rho(1-b\rho)} \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ \mathbf{e}_1 \\ 0 \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ 0 \\ \mathbf{e}_{d-1} \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \frac{\rho(e-q)}{1-b\rho} - p_\infty \\ 0 \\ \mathbf{0}_{d-1} \\ 0 \\ -\frac{(\gamma-1)(e-q)}{(1-b\rho)^2} \end{pmatrix} \right\}, \quad (4.17)$$

where \mathbf{e}_i are the standard basis vectors in \mathbb{R}^{d-1} . We have chosen specific eigenvectors for \mathbf{s}_1 and \mathbf{s}_3 since the dimension of their respective eigenspaces is only one.

Lemma 4.2.1 (Wave Structure of the Extended Riemann Problem). *The L- and R-waves are genuinely nonlinear and the C-wave is linearly degenerate. That is, $D\lambda_i(\tilde{\mathbf{u}}) \cdot \mathbf{r}_i(\tilde{\mathbf{u}}) \neq 0$ for $i = 1, 3$ and $D\lambda_2(\tilde{\mathbf{u}}) \cdot \mathbf{r}_2(\tilde{\mathbf{u}}) = 0$ for all $\tilde{\mathbf{u}} \in \tilde{\mathcal{B}}$.*

Proof. As mentioned above regarding the change of variables, we can equivalently work with eigenpairs (μ_i, \mathbf{s}_i) . Then we see that $D\mu_1(\tilde{\mathbf{w}}) \cdot \mathbf{s}_1(\tilde{\mathbf{w}}) = -\frac{\partial a}{\partial \rho}(\tilde{\mathbf{w}}) - \frac{a(\tilde{\mathbf{w}})}{\rho} \neq 0$ and $D\mu_3(\tilde{\mathbf{w}}) \cdot \mathbf{s}_3(\tilde{\mathbf{w}}) = \frac{\partial a}{\partial \rho}(\tilde{\mathbf{w}}) + \frac{a(\tilde{\mathbf{w}})}{\rho} \neq 0$ for all $\tilde{\mathbf{w}} \in \tilde{\mathcal{B}}$. Similarly, $D\mu_2(\tilde{\mathbf{w}}) = (0, 1, \mathbf{0}_{d-1}, 0, 0)^\top$ and is orthogonal to the space of eigenvectors defined in (4.17). \square

Lemma 4.2.2 (Continuity on the Contact). *The pressure and velocity are continuous across the contact.*

Proof. To prove this we show that p and v are 2-Riemann invariants and hence are continuous across the contact. Note the following identity holds for a i -characteristic wave,

$$\begin{aligned} D_{\tilde{\mathbf{u}}}\rho(\tilde{\mathbf{w}}) \cdot \mathbf{r}_i(\tilde{\mathbf{u}}) &= D_{\tilde{\mathbf{u}}}\rho(\theta(\tilde{\mathbf{u}})) \cdot (D_{\tilde{\mathbf{u}}}\theta(\tilde{\mathbf{u}}))^{-1} \mathbf{s}_i(\theta(\tilde{\mathbf{u}})) \\ &= ((D_{\tilde{\mathbf{w}}}\rho(\tilde{\mathbf{w}}))^\top D_{\tilde{\mathbf{u}}}\theta(\tilde{\mathbf{u}})) \cdot ((D_{\tilde{\mathbf{u}}}\theta(\tilde{\mathbf{u}}))^{-1} \mathbf{s}_i(\tilde{\mathbf{w}})) \\ &= D_{\tilde{\mathbf{w}}}\rho(\tilde{\mathbf{w}}) \cdot \mathbf{s}_i(\tilde{\mathbf{w}}) \end{aligned}$$

Therefore, we only need to show that $D_{\tilde{\mathbf{w}}}\rho(\tilde{\mathbf{w}}) \cdot \mathbf{s}_2(\tilde{\mathbf{w}}) = 0$ for ρ to be a 2-Riemann invariant. Computing the derivative of ρ , we have,

$$D_{\tilde{\mathbf{w}}}\rho(\tilde{\mathbf{w}}) = \left(\frac{(\gamma - 1)(e - q)}{(1 - b\rho)^2}, 0, \mathbf{0}^\top, \frac{(\gamma - 1)\rho}{1 - b\rho}, \frac{\rho(e - q)}{1 - b\rho} - p_\infty \right). \quad (4.18)$$

Then it is a quick check to see that $D_{\tilde{\mathbf{w}}}\rho(\tilde{\mathbf{w}}) \cdot \mathbf{s}_2(\tilde{\mathbf{w}}) = 0$ for any $\mathbf{s}_2(\tilde{\mathbf{w}})$ defined in (4.17). Therefore, ρ is a 2-Riemann invariant and hence is constant across the contact.

For the velocity, v is automatically a 2-Riemann invariant by Lemma 4.2.1, since $\lambda_2(\tilde{\mathbf{u}}) = v$. Thus the velocity is constant across the contact. □

4.3 The Solution to the Extended Riemann Problem

We suppose that the solution to the extended Riemann problem must be a self similar solution composed of three waves where the Z-wave ($Z \in \{L, R\}$) is either a shock or an expansion and the C-wave is a contact. Since γ is being transported ($\partial_t \gamma + v \partial_x \gamma = 0$), we propose that the $\gamma = \gamma_L$ left of the contact and $\gamma = \gamma_R$ right of the contact, Hence, the solution for γ is,

$$\gamma(x, t) = \begin{cases} \gamma_L, & \frac{x}{t} < v^*, \\ \gamma_R, & \frac{x}{t} > v^*, \end{cases} \quad (4.19)$$

where v^* is the speed of the contact. See Figure 4.1 for a visual description of this solution. The solution on each wave can be constructed with the NASG EOS for $\gamma = \gamma_Z$ for each respective wave. For the sake of completion, we derive the solution to this extended Riemann problem.

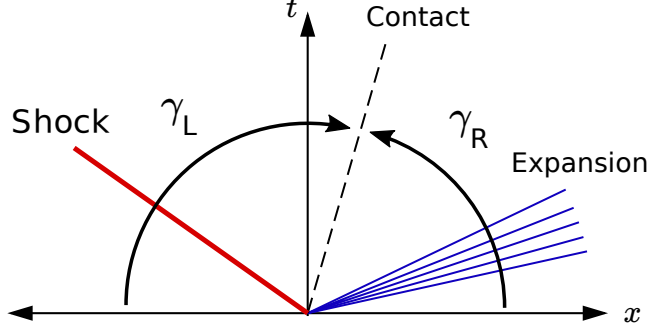


Figure 4.1: An example solution for $\gamma(x, t)$ in the extended Riemann problem

The solution to the Riemann is first constructed for the problem,

$$\partial_t \begin{pmatrix} \rho \\ m \\ \mathcal{E} \\ \Gamma \end{pmatrix} + \partial_x \begin{pmatrix} m \\ \frac{1}{\rho} m^2 + p_{\text{nasg}} \\ \frac{m}{\rho} (\mathcal{E} + p_{\text{nasg}}) \\ \frac{m}{\rho} \Gamma \end{pmatrix}, \quad (4.20)$$

with left and right data, $(\rho_Z, \mathbf{m}_Z \cdot \mathbf{n}, \mathcal{E}_Z, \Gamma_Z)^\top$ for $Z \in \{\text{L}, \text{R}\}$ where $\mathcal{E} = E - \frac{\|\mathbf{m}^\perp\|_{\mathcal{E}^2}^2}{2\rho}$. Then the complete solution is found for \mathbf{m}^\perp by solving $\partial_t \mathbf{m}^\perp + \partial_x (v \mathbf{m}^\perp) = \mathbf{0}$, see [16, Sec. 2.5]. However, the solution for \mathbf{m}^\perp is never needed, as we are only interested in the maximum wave speed to the problem, (4.20). This extended Riemann problem (4.20), is the focus of the remainder of this Chapter.

Remark 4.3.1 (Internal Energy Change of Basis). The internal energy can be written as $\rho e = E - \frac{\|\mathbf{m}^\perp\|_{\mathcal{E}^2}^2}{2\rho} = \mathcal{E} - \frac{m^2}{2\rho}$. Hence the internal energy does not depend on the basis, which is what we expect. \square

Now let $\mathbf{c} := (\rho, v, p, \gamma)^\top$ be the primitive state and set $\mathbf{c}_Z := (\rho_Z, v_Z, p_Z, \gamma_Z)^\top$. Recall that $\gamma_Z = \frac{(p_Z + p_\infty)(1 - b\rho_Z)}{\rho_Z(e_Z - q)} + 1$ and based on the assumption that the oracle provides positive pressure, see Remark 1.4.1, we have that $\min(\gamma_L, \gamma_R) > 1$. Furthermore, notice that the oracle is only invoked to compute the left and right Riemann data p_L and p_R .

We now define an important function that will appear in the solution to the Riemann problem,

$$f_Z(\rho) := \begin{cases} f_Z^{\text{exp}}(\rho) := \frac{2a_Z(1-b\rho_Z)}{\gamma_Z-1} \left(\left(\frac{\rho+p_\infty}{p_Z+p_\infty} \right)^{\frac{\gamma_Z-1}{2\gamma_Z}} - 1 \right), & \text{if } -p_\infty \leq \rho < p_Z, \\ f_Z^{\text{shock}}(\rho) := (\rho - p_Z) \sqrt{\frac{A_Z}{\rho+p_\infty+B_Z}}, & \text{if } \rho \geq p_Z. \end{cases} \quad (4.21)$$

where $A_Z := \frac{2(1-b\rho_Z)}{(\gamma_Z+1)\rho_Z}$ and $B_Z := \frac{\gamma_Z-1}{\gamma_Z+1}(p_Z + p_\infty)$.

We also introduce the wave speeds for the Riemann problem which are essential for the estimation of the maximum wave speed. They are as follows:

$$\lambda_L^-(\rho^*) := v_L - a_L \left(1 + \frac{\gamma_L + 1}{2\gamma_L} \left(\frac{\rho^* - p_L}{p_L + p_\infty} \right)_+ \right)^{\frac{1}{2}}, \quad (4.22a)$$

$$\lambda_L^+(\rho^*) := \begin{cases} v_L - f_L(\rho^*) - a_L \frac{1-b\rho_L}{1-b\rho_L^*} \left(\frac{\rho^*+p_\infty}{p_L+p_\infty} \right)^{\frac{\gamma_L-1}{2\gamma_L}}, & \text{if } \rho^* < p_L, \\ \lambda_L^-(\rho^*), & \text{if } p_L \leq \rho^*, \end{cases} \quad (4.22b)$$

$$\lambda_R^+(\rho^*) := v_R + a_R \left(1 + \frac{\gamma_R + 1}{2\gamma_R} \left(\frac{\rho^* - p_R}{p_R + p_\infty} \right)_+ \right)^{\frac{1}{2}}, \quad (4.22c)$$

$$\lambda_R^-(\rho^*) := \begin{cases} v_R + f_R(\rho^*) + a_R \frac{1-b\rho_R}{1-b\rho_R^*} \left(\frac{\rho^*+p_\infty}{p_R+p_\infty} \right)^{\frac{\gamma_R-1}{2\gamma_R}}, & \text{if } \rho^* < p_R, \\ \lambda_R^+(\rho^*), & \text{if } p_R \leq \rho^*, \end{cases} \quad (4.22d)$$

The derivation of these waves speeds is shown in the following sections.

4.3.1 Shock Wave

Assume the L-Wave is a shock. Then the solution with a shock wave must satisfy the Rankine-Hugoniot conditions. Let \mathcal{S}_L denote the speed of the shock and we denote the state across the shock by ‘*L’ subscript. Define, $\widehat{v}_L := v_L - \mathcal{S}_L$ and $\widehat{v}_{*L} := v_{*L} - \mathcal{S}_L$. Then the

Rankine Hugoniot conditions become,

$$\rho_L \widehat{v}_L = \rho_{*L} \widehat{v}_{*L}, \quad (4.23)$$

$$\rho_L \widehat{v}_L^2 + p_L = \rho_{*L} \widehat{v}_{*L}^2 + p^*, \quad (4.24)$$

$$\widehat{v}_L (\widehat{E}_L + p_L) = \widehat{v}_{*L} (\widehat{E}_{*L} + p^*), \quad (4.25)$$

where $\widehat{E}_L = \rho_L e_L + \frac{1}{2} \rho_L \widehat{v}_L^2$ and $\widehat{E}_{*L} = \rho_{*L} e_{*L} + \frac{1}{2} \rho_{*L} \widehat{v}_{*L}^2$. From the Rankine-Hugoniot conditions, one can also derive,

$$e_{*L} - e_L = \frac{1}{2} (p^* - p_L) \left(\frac{\rho_{*L} - \rho_L}{\rho_{*L} \rho_L} \right), \quad (4.26)$$

see [56, Section 3.1.3]. Note also that this identity is independent of the EOS. The goal is to determine an equation which relates the two unknowns, p^* and ρ_{*L} . So, applying the NASG EOS, we have,

$$\frac{p^* - \gamma_L p_\infty}{\gamma_L - 1} \frac{1 - b \rho_{*L}}{\rho_{*L}} - \frac{p_L - \gamma_L p_\infty}{\gamma_L - 1} \frac{1 - b \rho_L}{\rho_L} = \frac{1}{2} (p^* - p_L) \left(\frac{\rho_{*L} - \rho_L}{\rho_{*L} \rho_L} \right). \quad (4.27)$$

Now multiply the equation by $(\gamma_L - 1) \rho_{*L}$ and expand,

$$\begin{aligned} (p^* + \gamma_L p_\infty) - \rho_{*L} b (p^* + \gamma_L p_\infty) - (p_L + \gamma_L p_\infty) \frac{\rho_{*L}}{\rho_L} + \rho_{*L} b (p_L + \gamma_L p_\infty) \\ = \frac{\gamma_L - 1}{2} (p^* + p_L) \left(\frac{\rho_{*L}}{\rho_L} - 1 \right) \end{aligned} \quad (4.28)$$

We rearrange the equation so that all of the ρ_{*L} are on one side,

$$\begin{aligned} \rho_{*L} \left(b (p_L + \gamma_L p_\infty) - \frac{1}{\rho_L} (p_L + \gamma_L p_\infty) - b (p^* + \gamma_L p_\infty) - \frac{\gamma_L - 1}{2 \rho_L} (p^* + p_L) \right) \\ = - (p^* + \gamma_L p_\infty) - \frac{\gamma_L - 1}{2} (p^* + p_L). \end{aligned} \quad (4.29)$$

More rewriting yields,

$$\begin{aligned} -\frac{\rho_{*L}p_L}{2\rho_L} \left((\gamma_L - 1 + 2b\rho_L) \frac{p^*}{p_L} + \frac{2\gamma_L p_\infty}{p_L} + (\gamma_L + 1 - 2b\rho_L) \right) \\ = -\frac{1}{2}(\gamma_L + 1)p^* - \gamma_L p_\infty - \frac{1}{2}(\gamma_L - 1)p_L. \end{aligned} \quad (4.30)$$

Solving for ρ_{*L} ,

$$\rho_{*L} = \frac{\rho_L((\gamma_L + 1)p^* + 2\gamma_L p_\infty + (\gamma_L - 1)p_L)}{p_L((\gamma_L - 1 + 2b\rho_L) \frac{p^*}{p_L} + \frac{2\gamma_L p_\infty}{p_L} + (\gamma_L + 1 - 2b\rho_L))} \quad (4.31)$$

Dividing by ρ_L and rewriting, we have,

$$\frac{\rho_{*L}}{\rho_L} = \left[\frac{\frac{p^* + p_\infty}{p_L + p_\infty} + \frac{\gamma_L - 1}{\gamma_L + 1}}{\left(\frac{\gamma_L - 1 + 2b\rho_L}{\gamma_L + 1} \right) \frac{p^* + p_\infty}{p_L + p_\infty} + \frac{\gamma_L + 1 - 2b\rho_L}{\gamma_L + 1}} \right]. \quad (4.32)$$

From (4.24) we have,

$$\rho_L^2 \widehat{v}_L^2 = -\frac{p^* - p_L}{\frac{1}{\rho_{*L}} - \frac{1}{\rho_L}} = \rho_L \frac{p^* - p_L}{1 - \frac{\rho_L}{\rho_{*L}}}. \quad (4.33)$$

In addition, using (4.23) and (4.24), we can also derive the identity, $\rho_L \widehat{v}_L = \frac{p^* - p_L}{v_L - v_{*L}}$ which we can equate with (4.33), to find,

$$v_* = v_L - \sqrt{p^* - p_L} \sqrt{\frac{1}{\rho_L} \left(1 - \frac{\rho_L}{\rho_{*L}} \right)}. \quad (4.34)$$

Finally, we substitute (4.32) into (4.34) to find,

$$v^* = v_L - f_L^{\text{shock}}(p^*) \quad (4.35)$$

where f_L^{shock} is defined in (4.21). Thus the solution across the shock wave is,

$$\mathbf{u}_{*L} := \left(\frac{\rho_L((\gamma_L + 1)p^* + 2\gamma_L p_\infty + (\gamma_L - 1)p_L)}{p_L((\gamma_L - 1 + 2b\rho_L) \frac{p^*}{p_L} + \frac{2\gamma_L p_\infty}{p_L} + (\gamma_L + 1 - 2b\rho_L))}, v^* - f_L^{\text{shock}}(p^*), p^*, \gamma_L \right)^\top. \quad (4.36)$$

The density across the shock can be written in the following alternative form,

$$\mathbf{u}_{*L} := \left(\frac{\rho_L \left(\frac{p^*}{p_L} + \frac{2\gamma_L p_\infty}{(\gamma_L+1)p_L} + \frac{\gamma_L-1}{\gamma_L+1} \right)}{\frac{\gamma_L-1+2b\rho_L}{\gamma_L+1} \frac{p^*}{p_L} + \frac{2\gamma_L p_\infty}{(\gamma_L+1)p_L} + \frac{\gamma_L+1-2b\rho_L}{\gamma_L+1}}, v^* - f_L^{\text{shock}}(p^*), p^*, \gamma_L \right)^\top. \quad (4.37)$$

Next, we derive the shock speed. From (4.33), we can solve for shock speed, \mathcal{S}_L ,

$$\mathcal{S}_L(p^*) = v_L - \sqrt{\frac{p^* - p_L}{\rho_L \left(1 - \frac{\rho_L}{\rho_{*L}}\right)}}. \quad (4.38)$$

Using in (4.32) again and with a lot of simplification, we find,

$$\mathcal{S}_L = v_L - a_L \sqrt{\frac{\gamma_L + 1}{2\gamma_L} \left(\frac{p^* - p_L}{p_L + p_\infty} \right) + 1}, \quad (4.39)$$

where $a_L = \sqrt{\frac{\gamma_L(p_L + p_\infty)}{\rho_L(1-b\rho_L)}}$ is the sound speed for the NASG EOS. The shock speed is the wave speed for the L-wave. That is, we define $\lambda_1^-(p^*) = \lambda_1^+(p^*) := \mathcal{S}_L$.

The R-wave can be computed similarly. Doing so, we find,

$$v^* = v_R + f_R^{\text{shock}}(p) \quad (4.40)$$

and

$$\mathcal{S}_R(p^*) = v_R + a_R \sqrt{\frac{\gamma_R + 1}{2\gamma_R} \left(\frac{p^* - p_R}{p_R + p_\infty} \right) + 1}. \quad (4.41)$$

The wave speed on the R-wave is then $\lambda_R^-(p^*) = \lambda_R^+(p^*) = \mathcal{S}_R$ and the solution across the shock is given by,

$$\mathbf{u}_{*R} := \left(\frac{\rho_R \left((\gamma_R + 1)p^* + 2\gamma_R p_\infty + (\gamma_R - 1)p_R \right)}{p_R \left((\gamma_R - 1 + 2b\rho_R) \frac{p^*}{p_R} + \frac{2\gamma_R p_\infty}{p_R} + (\gamma_R + 1 - 2b\rho_R) \right)}, v^* + f_R^{\text{shock}}(p^*), p^*, \gamma_R \right)^\top. \quad (4.42)$$

or in the alternative form,

$$\mathbf{u}_{*R} := \left(\frac{\rho_R \left(\frac{p^*}{p_R} + \frac{2\gamma_R p_\infty}{(\gamma_R+1)p_R} + \frac{\gamma_R-1}{\gamma_R+1} \right)}{\frac{\gamma_R-1+2b\rho_R}{\gamma_R+1} \frac{p^*}{p_R} + \frac{2\gamma_R p_\infty}{(\gamma_R+1)p_R} + \frac{\gamma_R+1-2b\rho_R}{\gamma_R+1}}, v^* + f_R^{\text{shock}}(p^*), p^*, \gamma_R \right)^\top. \quad (4.43)$$

4.3.2 Rarefaction Wave

Now assume that the left wave is a rarefaction. Since the left wave uses the NASG EOS for $\gamma = \gamma_L$, we can use the existence of the specific entropy, s , for the NASG EOS. This is due to the NASG EOS being a convex equation of state; see Section 2.2. Since the specific entropy is constant along expansions, we use (2.12) to write the isentropic pressure law for the NASG EOS (with a slight abuse of the notation) as,

$$p(\rho) = C \left(\frac{\rho}{1-b\rho} \right)^{\gamma_L} - p_\infty, \quad (4.44)$$

where C is some constant depending on the specific entropy. Thus across the wave, for any state on the expansion connected to the left state, we have the following parametrization of the density,

$$\frac{1}{\rho(p)} - b = \left(\frac{1}{\rho_L} - b \right) \left(\frac{p_L + p_\infty}{p + p_\infty} \right)^{\frac{1}{\gamma_L}} \quad (4.45)$$

Using this parametrization, we can also parametrize the sound speed on the expansion, call it a_L^{exp} . Using (4.45) in the definition of the sound speed, (4.16), we see that,

$$a_L^{\text{exp}}(p) = a_L \frac{\rho_L}{\rho(p)} \left(\frac{p + p_\infty}{p_L + p_\infty} \right)^{\frac{\gamma_L+1}{2\gamma_L}} = a_L \frac{1 - b\rho_L}{1 - b\rho(p)} \left(\frac{p + p_\infty}{p_L + p_\infty} \right)^{\frac{\gamma_L-1}{2\gamma_L}}. \quad (4.46)$$

In order to determine the solution on the expansion we make use of the generalized 1-Riemann invariant on this wave. Recall the definition of the generalized Riemann invariant given in Definition 1.2.6. The corresponding right eigenvector for the 1-characteristic is $\mathbf{s}_1 = (\rho, -a_L^{\text{exp}}(p), \mathbf{0}_{d-1}, 0, 0)^\top$. The generalized 1-Riemann invariant, W_1 , can be found by solving, $DW_1(\mathbf{z}) \cdot \mathbf{s}_1 = 0$ where $\mathbf{z} = (\rho, v, p)$ is the vector of primitive variables. This gives

us the differential equation,

$$\frac{\partial W_1}{\partial \rho}(\mathbf{z}) + \frac{a_L^{\text{exp}}(\rho)}{\rho} \frac{\partial W_1}{\partial v}(\mathbf{z}) = 0. \quad (4.47)$$

A solution to this differential equation is,

$$W_1(\mathbf{z}) = v + \int_{\rho_0}^{\rho} \frac{a_L^{\text{exp}}(\mathbf{z})}{\varrho} d\varrho \quad (4.48)$$

Using the sound speed definition, (4.16), and the isentropic pressure law, (4.44), the Riemann invariant becomes,

$$W_1(\mathbf{z}) = v - \sqrt{C\gamma_L} \int_{\rho_0}^{\rho} \frac{\varrho^{\frac{1}{2}(\gamma_L-3)}}{(1-b\varrho)^{\frac{1}{2}(\gamma_L+1)}} d\varrho = v - \frac{2\sqrt{C\gamma_L}}{\gamma_L-1} \left(\frac{\varrho}{1-b\varrho} \right)^{\frac{1}{2}(\gamma_L-1)} \Big|_{\rho_0}^{\rho} \quad (4.49)$$

Furthermore, the sound speed on the expansion can be written as, $a_L^{\text{exp}}(\rho(\rho)) = \sqrt{\frac{C\gamma_L\rho^{\gamma_L-1}}{(1-b\rho)^{\gamma_L+1}}}$.

Hence the Riemann invariant takes the form,

$$W_1(\mathbf{z}) = v + \frac{2a_L^{\text{exp}}(\mathbf{z})}{\gamma_L-1} (1-b\rho) + \text{const.} \quad (4.50)$$

Note that the 1-Riemann invariant is constant along the 1-wave if the 1-wave is an expansion.

Hence, $W_1(\mathbf{z}_L) = W_1(\mathbf{z}_{*L})$. Therefore, we have the relationship for any state, \mathbf{z} on the expansion,

$$v + \frac{2a_L^{\text{exp}}(\mathbf{z})}{\gamma_L-1} (1-b\rho) = v_L + \frac{2a_L}{\gamma_L-1} (1-b\rho_L) \quad (4.51)$$

We now introduce the self similarity parameter, $\xi := x/t$. Hence, along the rarefaction, we have $\xi = v - a = v - a_L^{\text{exp}}(\rho)$. Using (4.51) we have the equation,

$$\xi(\rho) = v_L - \frac{2a_L}{\gamma_L-1} (1-b\rho_L) - \frac{2a_L^{\text{exp}}(\rho)}{\gamma_L-1} (1-b\rho(\rho)) - a_L^{\text{exp}}(\rho). \quad (4.52)$$

Using the definition of a_L^{exp} , (4.46), in (4.52) and with a lot of algebra, we arrive at the

following,

$$\xi_L(\boldsymbol{p}) = v_L - \frac{2a_L(1 - b\rho_L)}{\gamma_L - 1} \left(\left(\frac{\boldsymbol{p} + p_\infty}{p_L + p_\infty} \right)^{\frac{\gamma_L - 1}{2\gamma_L}} - 1 \right) - \frac{a_L(1 - b\rho_L)}{1 - b\rho(\boldsymbol{p})} \left(\frac{\boldsymbol{p} + p_\infty}{p_L + p_\infty} \right)^{\frac{\gamma_L - 1}{2\gamma_L}}. \quad (4.53)$$

Using the notation from (4.21), we write,

$$\xi_L(\boldsymbol{p}) = v_L - f_L^{\text{exp}}(\boldsymbol{p}) - \frac{a_L(1 - b\rho_L)}{1 - b\rho(\boldsymbol{p})} \left(\frac{\boldsymbol{p} + p_\infty}{p_L + p_\infty} \right)^{\frac{\gamma_L - 1}{2\gamma_L}}. \quad (4.54)$$

Notice that $\xi_L(\boldsymbol{p})$ is a strictly decreasing function for $\boldsymbol{p} \in [\boldsymbol{p}^*, p_L]$. Hence by the inverse function theorem, $\xi(\boldsymbol{p})$ is invertible, and thus \boldsymbol{p} can be expressed as a decreasing function of ξ for $\xi \in [\xi_L(p_L), \xi_L(\boldsymbol{p}^*)] = [\lambda_L^-(p_L), \lambda_L^+(\boldsymbol{p}^*)]$ where λ_L^- and λ_L^+ are defined in (4.22). Similarly, if an expansion on the R-wave occurs, then the parametrization of ξ is given by,

$$\xi_R(\boldsymbol{p}) = v_R + f_R^{\text{exp}}(\boldsymbol{p}) + \frac{a_R(1 - b\rho_R)}{1 - b\rho(\boldsymbol{p})} \left(\frac{\boldsymbol{p} + p_\infty}{p_R + p_\infty} \right)^{\frac{\gamma_R - 1}{2\gamma_R}}. \quad (4.55)$$

Then it can be seen that $\xi_R(\boldsymbol{p})$ is a strictly increasing function for $\boldsymbol{p} \in [\boldsymbol{p}^*, p_R]$. And again, $\xi_R(\boldsymbol{p})$ is invertible, and so we can conclude that $\boldsymbol{p} = \boldsymbol{p}(\xi)$ is a strictly increasing function for $\xi \in [\xi_R(\boldsymbol{p}^*), \xi_R(p_R)] = [\lambda_R^-(\boldsymbol{p}^*), \lambda_R^+(p_R)]$. Note however, we cannot find an explicit equation for $\boldsymbol{p}(\xi)$.

In short, we can say that the pressure, \boldsymbol{p} , decreases across expansions. Thus the solution for an expansion on the L-wave is,

$$\boldsymbol{u}_{LL}(\xi) := \begin{pmatrix} \left(b + \left(\frac{1}{\rho_L} - b \right) \left(\frac{p_L + p_\infty}{\boldsymbol{p}(\xi) + p_\infty} \right)^{\frac{1}{\gamma_L}} \right)^{-1} \\ v_L - f_L^{\text{exp}}(\boldsymbol{p}(\xi)) \\ p_{\text{nasg}}(\xi) \\ \gamma_L \end{pmatrix}, \quad (4.56)$$

and similarly, the solution on the R-wave is,

$$\mathbf{u}_{\text{RR}}(\xi) := \begin{pmatrix} \left(b + \left(\frac{1}{\rho_{\text{R}}} - b \right) \left(\frac{p_{\text{R}} + p_{\infty}}{\rho_{\text{nasg}}(\xi) + p_{\infty}} \right)^{\frac{1}{\gamma_{\text{R}}}} \right)^{-1} \\ v_{\text{R}} + f_{\text{R}}^{\text{exp}}(\rho(\xi)) \\ \rho_{\text{nasg}}(\xi) \\ \gamma_{\text{R}} \end{pmatrix}, \quad (4.57)$$

4.4 Connecting the L- and R-waves

Using the fact that the velocity and pressure are constant across the contact, we have that,

$$v^* = v_{\text{L}} - f_{\text{L}}(\rho^*) = v_{\text{R}} + f_{\text{R}}(\rho^*), \quad (4.58)$$

which defines the following nonlinear equation to solve for ρ^* ,

$$\varphi(\rho) := f_{\text{R}}(\rho) + f_{\text{L}}(\rho) + v_{\text{R}} - v_{\text{L}} = 0, \quad \text{for } \rho \in (-p_{\infty}, \infty), \quad (4.59)$$

where f_Z is defined in (4.21) for $Z \in \{\text{L}, \text{R}\}$.

Lemma 4.4.1 (Properties of φ). *The function φ defined in (4.59) is $C^2((-\infty, \infty))$, strictly increasing, concave, and $\varphi'''(\rho) > 0$ a.e.*

Proof. To see the proof, we must compute the derivative of the functions, f_Z^{exp} and f_Z^{shock} . Recall, $f_Z^{\text{exp}}(\rho) = \frac{2a_Z(1-b\rho_Z)}{\gamma_Z-1} \left(\left(\frac{\rho+p_{\infty}}{p_Z+p_{\infty}} \right)^{\frac{\gamma_Z-1}{2\gamma_Z}} - 1 \right)$ and $f_Z^{\text{shock}} = (\rho - p_Z) \sqrt{\frac{A_Z}{\rho+p_{\infty}+B_Z}}$. Then we have,

$$\frac{df_Z^{\text{exp}}}{d\rho}(\rho) = \frac{a_Z(1-b\rho_Z)}{\gamma_Z(p_Z+p_{\infty})} \left(\frac{\rho+p_{\infty}}{p_Z+p_{\infty}} \right)^{-\frac{\gamma_Z+1}{2\gamma_Z}} > 0, \quad (4.60a)$$

$$\frac{d^2f_Z^{\text{exp}}}{d\rho^2}(\rho) = -\frac{a_Z(\gamma_Z+1)(1-b\rho_Z)}{2\gamma_Z^2(p_Z+p_{\infty})^2} \left(\frac{\rho+p_{\infty}}{p_Z+p_{\infty}} \right)^{-\frac{3\gamma_Z+1}{2\gamma_Z}} < 0, \quad (4.60b)$$

$$\frac{d^3f_Z^{\text{exp}}}{d\rho^3}(\rho) = \frac{a_Z(3\gamma_Z+1)(\gamma_Z+1)(1-b\rho_Z)}{4\gamma_Z^3(p_Z+p_{\infty})^3} \left(\frac{\rho+p_{\infty}}{p_Z+p_{\infty}} \right)^{-\frac{5\gamma_Z+1}{2\gamma_Z}} > 0, \quad (4.60c)$$

and

$$\frac{df_Z^{\text{shock}}}{d\rho}(\rho) = \sqrt{\frac{A_Z}{\rho + p_\infty + B_Z}} \left(1 - \frac{\rho - p_Z}{2(\rho + p_\infty + B_Z)}\right) > 0, \quad (4.61a)$$

$$\frac{d^2 f_Z^{\text{shock}}}{d\rho^2}(\rho) = -\frac{\sqrt{A_Z}(\rho + p_\infty + 3(p_Z + p_\infty) + 4B_Z)}{4(\rho + p_\infty + B_Z)^{5/2}} < 0, \quad (4.61b)$$

$$\frac{d^3 f_Z^{\text{shock}}}{d\rho^3}(\rho) = \frac{3\sqrt{A_Z}(\rho + p_\infty + 5(p_Z + p_\infty) + 6B_Z)}{8(\rho + p_\infty + B_Z)^{7/2}} > 0. \quad (4.61c)$$

Based on the definition of $\varphi(\rho)$, (4.59), we see that φ is monotonically increasing, concave down and $\varphi'''(\rho) > 0$. Lastly, we need to show that $\varphi \in C^2((-p_\infty, \infty))$. As f_Z^{shock} and f_Z^{exp} are both C^3 functions, we just need to show that continuity of the derivatives holds at p_Z . Using the fact that $\frac{a_Z(1-b\rho_Z)}{\gamma_Z(p_Z+p_\infty)} = \sqrt{\frac{A_Z}{p_Z+p_\infty+B_Z}}$, we see that $\frac{df_Z^{\text{exp}}}{d\rho}(p_Z) = \frac{df_Z^{\text{shock}}}{d\rho}(p_Z)$, $\frac{d^2 f_Z^{\text{exp}}}{d\rho^2}(p_Z) = \frac{d^2 f_Z^{\text{shock}}}{d\rho^2}(p_Z)$, but

$$\frac{d^3 f_Z^{\text{exp}}}{d\rho^3}(p_Z) = \frac{a_Z(3\gamma_Z + 1)(\gamma_Z + 1)(1 - b\rho_Z)}{4\gamma_Z^3(p_Z + p_\infty)^3} \neq \frac{9\sqrt{A_Z}}{4(p_Z + p_\infty + B_Z)^{5/2}} = \frac{d^3 f_Z^{\text{shock}}}{d\rho^3}(p_Z). \quad (4.62)$$

This completes the proof. \square

Remark 4.4.1 (The Nonvacuum Condition). Since the function $\varphi(\rho)$ is strictly increasing, a root to the equation (4.59) only exists if $\varphi(-p_\infty) < 0$. This gives us the so-called *nonvacuum condition*

$$v_R - v_L < \frac{2a_L(1 - b\rho_L)}{\gamma_L - 1} + \frac{2a_R(1 - b\rho_R)}{\gamma_R - 1}. \quad (4.63)$$

If $\varphi(-p_\infty) \geq 0$, then we set $p^* := -p_\infty$ and the solution of the Riemann problem consists of two expansions into vacuum. \square

Lemma 4.4.2 (Unique Root for $\varphi(\rho) = 0$). *If (4.63) holds, then φ has a unique root, $p^* \in (-p_\infty, \infty)$.*

Proof. Note that condition (4.63) holds, then that implies $\varphi(-p_\infty) < 0$ and since φ is $C^3((-p_\infty, \infty))$ and strictly increasing, this implies the existence of a unique root. \square

Remark 4.4.2 (Fast Root Finding Method). Since φ has nice properties from Lemma 4.4.1, we can use the quadratic Newton method to quickly find the root of the equation. This fast estimation method was first introduced by Guermond & Popov in [59]. To avoid the nonlinear solver, we also propose an upper bound on the root in Section 4.6. \square

4.5 A Weak Solution to the Extended Riemann Problem

For $Z \in \{L, R\}$, define,

$$\mathbf{u}_Z^* := \begin{cases} \mathbf{u}_{ZZ}(\xi(\rho^*)), & \text{if } \rho^* < p_Z, \\ \mathbf{u}_{*Z}, & \text{if } p_Z \leq \rho^*. \end{cases} \quad (4.64)$$

For reference, \mathbf{u}_{ZZ} is defined in (4.56) and (4.57) and \mathbf{u}_{*Z} is defined in (4.37) and (4.43). We claim that,

$$\tilde{\mathbf{u}}(x, t) := \begin{cases} \mathbf{u}_L & \text{if } \frac{x}{t} < \lambda_L^-, \\ \mathbf{u}_{LL}(\frac{x}{t}) & \text{if } \lambda_L^- \leq \frac{x}{t} < \lambda_L^+, \\ \mathbf{u}_L^* & \text{if } \lambda_L^+ \leq \frac{x}{t} < v^* \\ \mathbf{u}_R^* & \text{if } v^* < \frac{x}{t} < \lambda_R^-, \\ \mathbf{u}_{RR}(\frac{x}{t}) & \text{if } \lambda_R^- \leq \frac{x}{t} < \lambda_R^+, \\ \mathbf{u}_R & \text{if } \lambda_R^+ \leq \frac{x}{t} \end{cases} \quad (4.65)$$

is a weak solution to the Riemann problem.

Lemma 4.5.1 (Upper Bound on the Density). *For $b > 0$, if $\tilde{\mathbf{u}}_L, \tilde{\mathbf{u}}_R \in \tilde{\mathcal{B}}(b, q, p_\infty)$, then the solution to the extended Riemann problem, $\tilde{\mathbf{u}}(x, t)$ given in (4.65) satisfies, $\rho < b^{-1}$.*

Proof. If the Z-wave is an expansion, then the parametrization of the density as a function of ρ , seen in (4.56) or (4.57) is an increasing function of ρ . As seen in Section 4.3.2, we know that the pressure decreases across an expansion hence, the density decreases across an expansion. Thus, $\rho(x, t) \leq \rho_Z < b^{-1}$ on an expansion wave.

Now assume that the Z-wave is a shock. From (4.37) or (4.43) we see that

$$\rho_{*Z} = \frac{\rho_Z \left(\frac{p^*}{p_Z} + \frac{2\gamma_Z p_\infty}{(\gamma_Z + 1)p_Z} + \frac{\gamma_Z - 1}{\gamma_Z + 1} \right)}{\frac{\gamma_Z - 1 + 2b\rho_Z}{\gamma_Z + 1} \frac{p^*}{p_Z} + \frac{2\gamma_Z p_\infty}{(\gamma_Z + 1)p_Z} + \frac{\gamma_Z + 1 - 2b\rho_Z}{\gamma_Z + 1}}. \quad (4.66)$$

Using that $\gamma_Z > 1$ and $\rho_Z < b^{-1}$, we can see through elementary calculus that ρ_{*Z} is an increasing function of p^* . Taking $p^* \rightarrow \infty$, we see that

$$\lim_{p^* \rightarrow \infty} \rho_{*Z}(p^*) = \frac{\gamma_Z + 1}{\gamma_Z - 1 + 2b\rho_Z} \rho_Z. \quad (4.67)$$

Thus, the upper bound is now seen as a decreasing function of γ_Z . Taking $\gamma_Z \rightarrow 1^+$, we find that,

$$\lim_{\gamma_Z \rightarrow 1^+} \frac{\gamma_Z + 1}{\gamma_Z - 1 + 2b\rho_Z} \rho_Z = b^{-1}. \quad (4.68)$$

Since $\gamma_Z > 1$, we have that $\rho_{*Z} < b^{-1}$ and therefore, $\rho < b^{-1}$ in the solution to the extended Riemann problem for a.e. $x \in \mathbb{R}$ and $t > 0$. \square

Lemma 4.5.2 (Weak Solution to the Extended Riemann Problem). *Assume the nonvacuum condition holds, (4.63). Then $\tilde{\mathbf{u}}(x, t)$ defined by (4.65), is a weak solution to the extended Riemann problem, (4.20). Moreover, $\tilde{\mathbf{u}}(x, t) \in \tilde{\mathcal{B}}(b, q, p_\infty)$.*

Proof. In the domain $\{x < v^*t\}$, we have that $\gamma = \gamma_L$ and hence $\Gamma = \gamma_L \rho$. Therefore, the last equation in (4.20), $\partial_t \Gamma + \partial_x(v\Gamma)$ is equivalent to the conservation of mass. Then by construction, the first three equations in (4.20) hold are satisfied in the weak sense since with the EOS defined by $p = (\gamma_L - 1) \frac{\rho(e-q)}{1-b\rho} - \gamma_L p_\infty$.

The result is the same in the domain $\{x > v^*t\}$ but with $\gamma = \gamma_R$. In order to complete the proof, we need the two states, $\tilde{\mathbf{u}}_L^* = (\rho_{*L}, v_L^*, p_L^*, \gamma_R^*)$ and $\tilde{\mathbf{u}}_R^* = (\rho_{*R}, v_R^*, p_R^*, \gamma_R^*)$ on the left and right of the contact, $\{x = v^*t\}$ to satisfy the Rankine-Hugoniot condition. Since the nonvacuum condition holds, we have that $v_L^* = v_R^* = v^*$ and $p_L^* = p_R^* = p^*$. The Rankine-Hugoniot conditions are immediately satisfied.

Since there is no vacuum state, the density remains strictly positive and by Lemma 4.5.1.

By the definition of ρ^* as the root of (4.59), we have that $\rho^* > -p_\infty$ in the nonvacuum case. In particular, in the case of a double expansion, we always have that $\rho > -p_\infty$. From the NASG EOS, we have that $(\gamma_Z - 1)\left(\frac{\rho(e-q)}{1-b\rho} - p_\infty\right) = \rho + p_\infty > 0$ on either the left or right of the contact. This implies the invariant-domain constraint that $e - q > p_\infty(\tau - b)$, since $\gamma_Z - 1 > 0$. Lastly, $\Gamma > \rho$ is equivalent to $\gamma_L > 1$ left of the contact and $\gamma_R > 1$ right of the contact. These again both hold since the Riemann data satisfy $\mathbf{u}_Z \in \mathcal{B}(b, q, p_\infty)$ for $Z \in \{L, R\}$. Therefore, we conclude that $\tilde{\mathbf{u}}(\mathbf{x}, t) \in \tilde{\mathcal{B}}(b, q, p_\infty)$. \square

This brings us to the main result from Clayton et. al. [2], which has been modified to hold for the NASG interpolatory EOS.

Theorem 4.5.1. (i) Let $\mathbf{U}_i^n, \mathbf{U}_j^n \in \mathcal{B}(b, q, p_\infty)$ where $\mathcal{B}(b, q, p_\infty)$ is defined in (1.23). Let ρ^* be the root of the equation $\varphi(\rho) = 0$ defined in (4.59) and let $\hat{\rho}^*$ be any upper bound on ρ^* ; i.e., $\hat{\rho}^* \geq \rho^*$. Let,

$$\hat{\lambda}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n) := \max(-\lambda_L^-(\hat{\rho}^*), \lambda_R^+(\hat{\rho}^*)) \quad (4.69a)$$

$$d_{ij}^{L,n} := \max(\hat{\lambda}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n) \|\mathbf{c}_{ij}\|_{\ell^2}, \hat{\lambda}(\mathbf{n}_{ji}, \mathbf{U}_j^n, \mathbf{U}_i^n) \|\mathbf{c}_{ji}\|_{\ell^2}). \quad (4.69b)$$

Let $\bar{\mathbf{U}}_{ij}^n(t_0)$ be defined in (4.3) with $t_0 = \frac{\|\mathbf{c}_{ij}\|}{2d_{ij}^{L,n}}$. Then $\bar{\mathbf{U}}_{ij}^n(t_0) \in \mathcal{B}(b, q, p_\infty)$.

(ii) Let $i \in \mathcal{V}$ and $\mathbf{U}_j^n \in \mathcal{B}(b, q, p_\infty)$ for all $j \in \mathcal{G}(i)$. Let $d_{ij}^{L,n}$ be defined by (4.69b) and assume Δt is small enough so that $\Delta t \sum_{j \in \mathcal{G}(i) \setminus \{i\}} \frac{2d_{ij}^{L,n}}{m_i} \leq 1$. If $\mathbf{U}_i^{L,n+1}$ is defined by (4.1), then $\mathbf{U}_i^{L,n+1} \in \text{conv}\{\bar{\mathbf{U}}_{ij}^n : j \in \mathcal{G}(i)\} \subset \mathcal{B}(b, q, p_\infty)$.

Proof. First notice that $-\lambda_L^-(\rho)$ and $\lambda_R^+(\rho)$ are both increasing functions of ρ . Therefore, for $\hat{\rho}^* \geq \rho^*$, we have that $\hat{\lambda}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n) \geq \max(-\lambda_L^-(\rho^*), \lambda_R^+(\rho^*))$. This implies that,

$$\left(0, \frac{1}{2\hat{\lambda}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n)}\right] \subset \left(0, \frac{1}{2\max(-\lambda_L^-(\rho^*), \lambda_R^+(\rho^*))}\right]. \quad (4.70)$$

Thus by definition of $d_{ij}^{L,n}$ and t_0 , we have that $t_0 \in (0, 1/(2\hat{\lambda}(\mathbf{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n))]$. In order to apply Lemma 1.3.1, we use $\mathbf{g}(\mathbf{u}) := \tilde{\mathbf{f}}(\tilde{\mathbf{u}})\mathbf{n}_{ij}$ with the left and right Riemann data $\tilde{\mathbf{U}}_i^n$ and

$\tilde{\mathbf{U}}_j^n$, respectively. Then we must show that the solution $\tilde{\mathbf{u}}(x, t) \in \tilde{\mathcal{B}}(b, q, p_\infty)$ for a.e. $x \in \mathbb{R}$ and all $t > 0$.

Note, the solution given in (4.65) satisfies $\rho > 0$ in every case and by Lemma 4.5.1, $\rho < b^{-1}$ (if $b > 0$). For the constraint, $e - q > p_\infty(\tau - b)$ recall that $e - q = \frac{p + \gamma p_\infty}{\gamma - 1}(\tau - b)$. By construction of the solution, $p^* \in (-p_\infty, \infty)$ and hence $e - q > p_\infty(\tau - b)$. Lastly, since $\gamma(x, t) = \gamma_L$ if $x < v^*t$, $\gamma(x, t) = \gamma_R$ if $x > v^*t$ and $\gamma_L, \gamma_R > 1$, this implies that $\Gamma(x, t) > \rho(x, t)$ for a.e. $x \in \mathbb{R}$ and all $t > 0$. Therefore, $\tilde{\mathbf{u}} \in \tilde{\mathcal{B}}(b, q, p_\infty)$.

From Lemma 4.5.1, 2., we conclude that $\tilde{\mathbf{U}}_{ij}^n(t_0) \in \tilde{\mathcal{B}}(b, q, p_\infty)$. But from (4.12), the density, momentum and total energy are the same for $\mathbf{U}_{ij}^n(t_0)$ and $\tilde{\mathbf{U}}_{ij}^n(t_0)$. So, defining, $\tilde{\Psi}_1(\tilde{\mathbf{u}}) := \rho$, $\tilde{\Psi}_2(\tilde{\mathbf{u}}) := 1 - b\rho$, and $\tilde{\Psi}_3(\tilde{\mathbf{u}}) := \mathbf{e}(\mathbf{u}) - q - p_\infty(\rho^{-1} - b)$. We have that $\Psi_l(\mathbf{U}_{ij}^n(t_0)) = \tilde{\Psi}_l(\tilde{\mathbf{U}}_{ij}^n(t_0)) > 0$ for all $l = 1, 2, 3$. Therefore, we conclude that $\mathbf{U}_{ij}^n(t_0) \in \mathcal{B}(b, q, p_\infty)$. \square

4.5.1 Weak solution with Vacuum state

In the case that the nonvacuum condition, (4.63), fails, then $p^* := -p_\infty$ and the velocity in Riemann fan is no longer continuous. That is,

$$v_L^* := v_L - f_L^{\text{exp}}(-p_\infty) = v_L + \frac{2a_L(1 - b\rho_L)}{\gamma_L - 1}, \quad (4.71a)$$

$$v_R^* := v_R + f_R^{\text{exp}}(-p_\infty) = v_R - \frac{2a_R(1 - b\rho_R)}{\gamma_R - 1}. \quad (4.71b)$$

Define $\tilde{\mathbf{u}}_Z^* := \tilde{\mathbf{u}}_{ZZ}(v_Z^*) = (0, v_Z^*, -p_\infty, \gamma_Z)^\top$ for $Z \in \{L, R\}$. Then we claim that the solution to the extended Riemann problem in the presence of vacuum is,

$$\tilde{\mathbf{u}}(x, t) := \begin{cases} \tilde{\mathbf{u}}_L & \text{if } \frac{x}{t} < v_L - a_L, \\ \tilde{\mathbf{u}}_{LL}(\frac{x}{t}) & \text{if } \lambda_L^- \leq \frac{x}{t} < v_L^*, \\ \frac{v_R^* - \frac{x}{t}}{v_R^* - v_L^*} \tilde{\mathbf{u}}_L^* + \frac{\frac{x}{t} - v_L^*}{v_R^* - v_L^*} \tilde{\mathbf{u}}_R^* & \text{if } v_L^* \leq \frac{x}{t} < v_R^* \\ \tilde{\mathbf{u}}_{RR}(\frac{x}{t}) & \text{if } v_R^* \leq \frac{x}{t} < v_R + a_R, \\ \tilde{\mathbf{u}}_R & \text{if } v_R + a_R \leq \frac{x}{t} \end{cases} \quad (4.72)$$

Lemma 4.5.3 (Weak Solution with Vacuum). *Assume that the nonvacuum condition (4.63) fails, i.e., $p^* = 0$. Then the following is true*

1. *The solution $\tilde{\mathbf{u}}$ given by (4.72) is a weak solution to (4.20).*
2. *If $\frac{x}{t} \in (v_L - a_L, v_L^*) \cup (v_R^*, v_R + a_R)$, then $\tilde{\mathbf{u}} \in \tilde{\mathcal{B}}(b, q, p_\infty)$.*
3. *$\tilde{\mathbf{u}} \in \overline{\tilde{\mathcal{B}}}(b, q, p_\infty)$; that is, $\tilde{\mathbf{u}}$ is in the closure of $\tilde{\mathcal{B}}(b, q, p_\infty)$.*

Proof. 1. We have already established that $\tilde{\mathbf{u}}$ is a weak solution to (4.20) from Lemma 4.5.2 in the region $\{x < v_L^*t\} \cup \{v_R^*t < x\}$. In the region $\{v_L^*t < x < v_R^*t\}$, the solution defined by (4.72) is a weak solution since all of the conservative variables are all zero. We only need to show to show that $\tilde{\mathbf{u}}$ is continuous on the line $\{x = v_L^*t\}$ and $\{x = v_R^*t\}$. From the definition of $\xi_L(p)$ given in (4.54), we see that $\xi_L(-p_\infty) = v_L - f_L^{\text{exp}}(-p_\infty) = v_L^*$ therefore, $p(v_L^*) = -p_\infty$. Hence $\lim_{\xi \uparrow v_L^*} \tilde{\mathbf{u}}_{LL}(\xi) = (0, v_L^*, -p_\infty, \gamma_L)^\top$. Similarly, taking $\lim_{\xi \downarrow v_L^*} \frac{v_R^* - \frac{x}{t}}{v_R^* - v_L^*} \tilde{\mathbf{u}}_L^* + \frac{\frac{x}{t} - v_L^*}{v_R^* - v_L^*} \tilde{\mathbf{u}}_R^* = (0, v_L^*, -p_\infty, \gamma_L)^\top$. Hence the solution is continuous on the line $\{x = v_L^*t\}$. A similar proof holds for the line $\{x = v_R^*t\}$.

2. This proof follows identically to what is given in the proof of Lemma 4.5.2.

3. We only need to show that the state in the region $\{v_L^*t < x < v_R^*t\}$ lies in $\overline{\tilde{\mathcal{B}}}(b, q, p_\infty)$. First note that, $\frac{v_R^* - \frac{x}{t}}{v_R^* - v_L^*} \Gamma_L^* + \frac{\frac{x}{t} - v_L^*}{v_R^* - v_L^*} \Gamma_R^* = 0$ for $\frac{x}{t} \in (v_L^*, v_R^*)$ since $\rho = 0$. Hence $\Gamma(x, t) \geq \rho(x, t)$ for a.e. $x \in \mathbb{R}$ and all $t > 0$.

In the case that $p_\infty > 0$, then $p = -p_\infty$ implies that $\frac{\rho(e-q)}{1-b\rho} = p_\infty$. In particular, as $\rho \downarrow 0$, we have that $e \rightarrow \infty$. Nevertheless, we have that, $\tilde{\mathbf{u}} \in \overline{\tilde{\mathcal{B}}}(b, q, p_\infty)$.

For $p_\infty = 0$, we see that $e \downarrow q$ as $\rho \downarrow 0$. Hence $\tilde{\mathbf{u}} \in \{\tilde{\mathbf{w}} \in \mathbb{R}^{d+3} : \rho \geq 0, 1 - b\rho > 0, e - q \geq 0, \Gamma \geq \rho\} \subset \overline{\tilde{\mathcal{B}}}(b, q, p_\infty)$. \square

4.6 Upper Bound on the Max Wave Speed

We now have the following important result regarding the max wave speed.

Theorem 4.6.1 (Upper Bound on the Max Wave Speed). *Let $\hat{\lambda}_{\max} := \lambda_{\max}(\hat{\rho}^*)$, where $\hat{\rho}^*$ is any upper bound on p^* , where p^* is the root of equation (4.59). Then $\hat{\lambda}_{\max} \geq \lambda_{\max}(p^*)$*

and the update (4.1) is invariant domain preserving with $\widehat{\lambda}$.

Recall that the equation to find ρ^* requires the solution of a nonlinear equation. This requires the use of a root finding method such as the Newton-Raphson method or the Newton-Secant method. One can also use a quadratic Newton method since $\varphi'''(p) > 0$ for all $p \geq -p_\infty$, details on this can be found in [59]. Unfortunately, despite this, it is still expensive to solve this nonlinear equation. We can alternatively use an upper estimate, $\widehat{\rho}^*$ on ρ^* to compute the max wave speed since the invariant-domain preserving properties remain unchanged. In this section, we outline the details for the computation of a close upper bound.

We begin by first introducing the *double-rarefaction* approximation of the function $\varphi(p)$,

$$\begin{aligned} \varphi_{RR}(p) := & \frac{2a_L(1 - b\rho_L)}{\gamma_L - 1} \left(\left(\frac{p + p_\infty}{p_L + p_\infty} \right)^{\frac{\gamma_L - 1}{2\gamma_L}} - 1 \right) \\ & + \frac{2a_R(1 - b\rho_R)}{\gamma_R - 1} \left(\left(\frac{p + p_\infty}{p_R + p_\infty} \right)^{\frac{\gamma_R - 1}{2\gamma_R}} - 1 \right) + v_R - v_L, \end{aligned} \quad (4.73)$$

Lemma 4.6.1 (Approximation of φ). *If $\max(\gamma_L, \gamma_R) \leq \frac{5}{3}$, then $\varphi_{RR}(p) \leq \varphi(p)$ for all $p \geq -p_\infty$. In particular, $\varphi_{RR}(p) = \varphi(p)$ for $p \in [-p_\infty, \min(p_L, p_R)]$.*

Proof. The proof for the case when $p_\infty = 0$ can be seen in [59, Theorem 4.1]. For $p_\infty \neq 0$, the proof is just a translation of the original. \square

Note that the current approximation with φ_{RR} fails for $\gamma_Z > \frac{5}{3}$. To address this we have the following lemma,

Lemma 4.6.2 (Approximation of $f_Z^{\text{shock}}(p)$). *For all $p > p_Z$, we have that $f^{\text{shock}}(p) \geq c(\gamma_Z)f^{\text{exp}}(p)$, where*

$$c(\gamma_Z) := \begin{cases} 1, & \text{if } 1 < \gamma_Z \leq \frac{5}{3} \\ \left(\frac{1}{2} + \frac{4}{3(\gamma_Z + 1)} \right)^{\frac{1}{2}}, & \text{if } \frac{5}{3} \leq \gamma_Z \leq 3, \\ \left(\frac{1}{2} + \frac{2}{\gamma_Z - 1} 3^{\frac{4 - 2\gamma_Z}{\gamma_Z - 1}} \right)^{\frac{1}{2}}, & \text{if } 3 \leq \gamma_Z. \end{cases} \quad (4.74)$$

Proof. The proof can be found in the supplementary material to [2] at [local/web 319KB] \square

Notice in passing that $\gamma_Z \mapsto c(\gamma_Z)$ is continuous and $c(\gamma_Z) \in (\frac{1}{\sqrt{2}}, 1]$. Lemma 4.6.2 enables us to get an alternative double-rarefaction approximation of the $\varphi(\boldsymbol{p})$ as we will see shortly. To simplify upcoming notation we define the following **subscripts**. Let

$$\min := \begin{cases} \text{L}, & \text{if } p_L \leq p_R, \\ \text{R}, & \text{if } p_L > p_R, \end{cases} \quad \max := \begin{cases} \text{R}, & \text{if } p_L \leq p_R, \\ \text{L}, & \text{if } p_L > p_R. \end{cases} \quad (4.75)$$

With this notation we have that $p_{\min} = \min(p_L, p_R)$ and $p_{\max} = \max(p_L, p_R)$. Or in other words, $\min = \arg \min_{Z \in \{\text{L}, \text{R}\}}(p_Z)$ and $\max = \arg \max_{Z \in \{\text{L}, \text{R}\}}(p_Z)$. Similarly, we define the following subscripts in relation to γ_Z ,

$$m := \begin{cases} \text{L}, & \text{if } \gamma_L \leq \gamma_R, \\ \text{R}, & \text{if } \gamma_L > \gamma_R, \end{cases} \quad M := \begin{cases} \text{R}, & \text{if } \gamma_L \leq \gamma_R, \\ \text{L}, & \text{if } \gamma_L > \gamma_R. \end{cases} \quad (4.76)$$

So $\gamma_m = \min(\gamma_L, \gamma_R)$ and $\gamma_M = \max(\gamma_L, \gamma_R)$. It is important to notice that the subscripts m and \min do not have to coincide (similarly for M and \max). It is worth emphasizing that “min” and “max” are **subscripts** representing either L or R.

We are now ready to define the approximation \widehat{p}^* for p^* .

4.6.1 Case 0: Vacuum

In the case of vacuum, i.e., $v_R - v_L \geq \frac{2a_L(1-b\rho_L)}{\gamma_L-1} + \frac{2a_R(1-b\rho_R)}{\gamma_R-1}$, there is no root $\varphi(\boldsymbol{p}) = 0$ and hence $p^* := 0$. This implies that $\lambda_1^-(0)v_L - a_L$ and $\lambda_3^+(0) = v_R + a_R$, which is what we expect as vacuum only occurs for a double expansion. Hence, we simply set $\widehat{p}^* = 0$.

4.6.2 Case 1: $p^* > 0$ and $\varphi(p_{\min}) > 0$

In this case the solution to the Riemann problem is composed of two expansions and therefore, the equation to solve is $\varphi_{RR}(\boldsymbol{p}) = 0$. However, for $\gamma_L \neq \gamma_R$, the equation is nonlinear. Fortunately for us, we only need the maximum wave speed, therefore an estimate

for \boldsymbol{p}^* is not necessary since, $\lambda_1^-(\boldsymbol{p}^*) = v_L - a_L$ and $\lambda_3^+(\boldsymbol{p}^*) = v_R + a_R$. If one does need an approximation to \boldsymbol{p}^* , then one can solve the equation $\widehat{\varphi}(\boldsymbol{p}) = 0$ where,

$$\widehat{\varphi}_{RR}(\boldsymbol{p}) := \alpha_L \left(\left(\frac{\boldsymbol{p} + p_\infty}{p_L + p_\infty} \right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1 \right) + \alpha_R \left(\left(\frac{\boldsymbol{p} + p_\infty}{p_R + p_\infty} \right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1 \right) + v_R - v_L. \quad (4.77)$$

Note that $\widehat{\varphi}_{RR}(\boldsymbol{p}) \leq \varphi_{RR}(\boldsymbol{p}) = \varphi(\boldsymbol{p})$ for all $\boldsymbol{p} \in [-p_\infty, p_{\min}]$. The associated root of $\widehat{\varphi}(\boldsymbol{p}) = 0$ is,

$$\widetilde{\boldsymbol{p}}^* = \left(\frac{\max(\alpha_R + \alpha_L - (v_R - v_L), 0)}{\alpha_R (p_R + p_\infty)^{-\frac{\gamma_M - 1}{2\gamma_M}} + \alpha_L (p_L + p_\infty)^{-\frac{\gamma_M - 1}{2\gamma_M}}} \right)^{\frac{2\gamma_M}{\gamma_M - 1}} - p_\infty \quad (4.78)$$

Note that $\widetilde{\boldsymbol{p}}^* \geq \boldsymbol{p}^*$, but by assumption, we also have that $p_{\min} \geq \boldsymbol{p}^*$. Therefore, we define the upper estimate by $\widehat{\boldsymbol{p}}^* := \min(\widetilde{\boldsymbol{p}}^*, p_{\min})$.

4.6.3 Case 2: $\varphi(p_{\min}) \leq 0 \leq \varphi(p_{\max})$

In this case, we have that $\boldsymbol{p}^* \in [p_{\min}, p_{\max}]$ and therefore, the “min” wave is a shock; that is $f_{\min}(\boldsymbol{p}) = f_{\min}^{\text{shock}}(\boldsymbol{p})$ and the “max” wave is an expansion; i.e., $f_{\max}(\boldsymbol{p}) = f_{\max}^{\text{exp}}(\boldsymbol{p})$. Therefore, using Lemma 4.6.2 we approximate f_{\min}^{shock} by, $c(\gamma_{\min})f_{\min}^{\text{exp}}(\boldsymbol{p}) \leq f_{\min}^{\text{shock}}(\boldsymbol{p})$ for $\boldsymbol{p} \in [p_{\min}, p_{\max}]$ and leave f_{\max}^{exp} unchanged. That is, the lower bounding function is,

$$\widehat{\varphi}_{RR}(\boldsymbol{p}) := \widehat{\alpha}_{\min} \left(\left(\frac{\boldsymbol{p} + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_{\min} - 1}{2\gamma_{\min}}} - 1 \right) + \alpha_{\max} \left(\left(\frac{\boldsymbol{p} + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_{\max} - 1}{2\gamma_{\max}}} - 1 \right) + v_R - v_L. \quad (4.79)$$

Unfortunately, this is still a nonlinear equation. So we need to coarsen our approximation a bit more. To do this we declare two sub cases.

4.6.3.1 Case 2a: $\gamma_{\min} = \gamma_m$

Define,

$$\widehat{\varphi}_1(p) := \widehat{\alpha}_{\min} \left(\delta \left(\frac{p + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1 \right) + \alpha_{\max} \left(\left(\frac{p + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1 \right) + v_R - v_L, \quad (4.80a)$$

$$\widehat{\varphi}_2(p) := \widehat{\alpha}_{\min} \left(\left(\frac{p + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1 \right) + \alpha_{\max} \left(\delta \left(\frac{p + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1 \right) + v_R - v_L, \quad (4.80b)$$

where $\delta := \left(\frac{p_{\min} + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_M - \gamma_m}{2\gamma_m \gamma_M}}$. Observe that,

$$\begin{aligned} \left(\frac{p + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_m - 1}{2\gamma_m}} &\geq \left(\frac{p + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_M - 1}{2\gamma_M}} \left(\frac{p_{\max} + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_m - 1}{2\gamma_m}} \\ &= (p + p_\infty)^{\frac{\gamma_M - 1}{2\gamma_M}} (p_{\max} + p_\infty)^{\frac{\gamma_m - 1}{2\gamma_m} - \frac{\gamma_M - 1}{2\gamma_M}} (p_{\min} + p_\infty)^{-\frac{\gamma_m - 1}{2\gamma_m}} \\ &= \left(\frac{p + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_M - 1}{2\gamma_M}} \left(\frac{p_{\min} + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_M - \gamma_m}{2\gamma_m \gamma_M}}. \end{aligned}$$

Therefore, $\widehat{\varphi}_1(p) \leq \widehat{\varphi}(p)$ for $p \in [p_{\min}, p_{\max}]$. We can apply the same reasoning to show $\widehat{\varphi}_2(p) \leq \widehat{\varphi}(p)$. Thus $\max(\widehat{\varphi}_1(p), \widehat{\varphi}_2(p)) \leq \widehat{\varphi}(p)$ for $p \in [p_{\min}, p_{\max}]$. The roots of (4.80a) and (4.80b) are,

$$\widetilde{p}_1^* = \left(\frac{\widehat{\alpha}_{\min} + \alpha_{\max} - (v_R - v_L)}{\delta \widehat{\alpha}_{\min} (p_{\min} + p_\infty)^{-\frac{\gamma_M - 1}{2\gamma_M}} + \alpha_{\max} (p_{\max} + p_\infty)^{-\frac{\gamma_M - 1}{2\gamma_M}}} \right)^{\frac{2\gamma_M}{\gamma_M - 1}} - p_\infty \quad (4.81a)$$

$$\widetilde{p}_2^* = \left(\frac{\widehat{\alpha}_{\min} + \alpha_{\max} - (v_R - v_L)}{\widehat{\alpha}_{\min} (p_{\min} + p_\infty)^{-\frac{\gamma_m - 1}{2\gamma_m}} + \delta \alpha_{\max} (p_{\max} + p_\infty)^{-\frac{\gamma_m - 1}{2\gamma_m}}} \right)^{\frac{2\gamma_m}{\gamma_m - 1}} - p_\infty \quad (4.81b)$$

Thus, our approximation for p^* is $\widehat{p}^* := \min(\widetilde{p}_1^*, \widetilde{p}_2^*, p_{\max})$.

4.6.4 Case 2b: $\gamma_{\min} = \gamma_M$

This case is very similar to Case 2a. Since $\gamma_{\min} = \gamma_M$ and $\gamma_{\max} = \gamma_M$ we have that the following lower bounds on $\widehat{\varphi}$,

$$\widehat{\varphi}_1(p) := \widehat{\alpha}_{\min} \left(\left(\frac{p + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1 \right) + \alpha_{\max} \left(\left(\frac{p + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1 \right) + v_R - v_L, \quad (4.82a)$$

$$\widehat{\varphi}_2(p) := \widehat{\alpha}_{\min} \left(\left(\frac{p + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1 \right) + \alpha_{\max} \left(\left(\frac{p + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1 \right) + v_R - v_L, \quad (4.82b)$$

and they have the corresponding roots,

$$\widetilde{p}_1^* = \left(\frac{\widehat{\alpha}_{\min} + \alpha_{\max} - (v_R - v_L)}{\widehat{\alpha}_{\min}(p_{\min} + p_\infty)^{-\frac{\gamma_m - 1}{2\gamma_m}} + \alpha_{\max}(p_{\max} + p_\infty)^{-\frac{\gamma_m - 1}{2\gamma_m}}} \right)^{\frac{2\gamma_m}{\gamma_m - 1}} - p_\infty, \quad (4.83a)$$

$$\widetilde{p}_2^* = \left(\frac{\widehat{\alpha}_{\min} + \alpha_{\max} - (v_R - v_L)}{\widehat{\alpha}_{\min}(p_{\min} + p_\infty)^{-\frac{\gamma_M - 1}{2\gamma_M}} + \alpha_{\max}(p_{\max} + p_\infty)^{-\frac{\gamma_M - 1}{2\gamma_M}}} \right)^{\frac{2\gamma_M}{\gamma_M - 1}} - p_\infty. \quad (4.83b)$$

Therefore our upper bound on p^* is $\widehat{p}^* := \min(\widetilde{p}_1^*, \widetilde{p}_2^*, p_{\max})$.

4.6.4.1 Case 3: $\varphi(p_{\max}) < 0$

In this case, we have $f_{\min}^{\text{shock}}(p) \geq c(\gamma_{\min})f_{\min}^{\text{exp}}$ and $f_{\max}^{\text{shock}}(p) \geq c(\gamma_{\max})f_{\max}^{\text{exp}}$ for all $p \geq p_{\max}$.

Therefore, a lower approximation of φ is,

$$\widehat{\varphi}_{RR}(p) := \widehat{\alpha}_{\min} \left(\left(\frac{p + p_\infty}{p_{\min} + p_\infty} \right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1 \right) + \widehat{\alpha}_{\max} \left(\left(\frac{p + p_\infty}{p_{\max} + p_\infty} \right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1 \right) + v_R - v_L, \quad (4.84)$$

with the corresponding root,

$$\widetilde{p}_1^* = \left(\frac{\widehat{\alpha}_{\min} + \widehat{\alpha}_{\max} - (v_R - v_L)}{\widehat{\alpha}_{\min}(p_{\min} + p_\infty)^{-\frac{\gamma_m - 1}{2\gamma_m}} + \widehat{\alpha}_{\max}(p_{\max} + p_\infty)^{-\frac{\gamma_m - 1}{2\gamma_m}}} \right)^{\frac{2\gamma_m}{\gamma_m - 1}} - p_\infty. \quad (4.85)$$

Alternatively, we can make use of the fact that φ is composed of the two shock curves, f_L^{shock} and f_R^{shock} and derive an alternative lower bound. For $p \in (p_{\max}, \infty)$ we have that $B_Z \leq B_Z p p_{\max}^{-1}$. Therefore,

$$f_Z^{\text{shock}}(p) = (p - p_Z) \sqrt{\frac{A_Z}{p + B_Z}} \geq \frac{p - p_Z}{\sqrt{p}} \sqrt{\frac{A_Z}{1 + B_Z p_{\max}^{-1}}} \quad (4.86)$$

and hence we can define the following lower bound on φ ,

$$\widehat{\varphi}_{SS}(p) := \frac{p - p_L}{\sqrt{p}} \sqrt{\frac{A_L}{1 + B_L p_{\max}^{-1}}} + \frac{p - p_R}{\sqrt{p}} \sqrt{\frac{A_R}{1 + B_R p_{\max}^{-1}}} + v_R - v_L. \quad (4.87)$$

Solving $\widehat{\varphi}_{SS}(p) = 0$ is equivalent to solving the quadratic equation,

$$(x_L + x_R)p + (v_R - v_L)\sqrt{p} - (p_L x_L + p_R x_R) = 0, \quad (4.88)$$

where $x_Z = \sqrt{\frac{A_Z}{1 + B_Z p_{\max}^{-1}}}$. Define $a := x_L + x_R$, $b := v_R - v_L$, $c := -p_L x_L - p_R x_R$. Then the root of this equation is,

$$\widetilde{p}_2^* = \left(\frac{-b + \sqrt{b^2 - 4ac}}{2a} \right)^2 \quad (4.89)$$

Therefore, the upper bound on p^* is $\widehat{p}^* := \min(\widetilde{p}_1^*, \widetilde{p}_2^*)$.

This completes all possible cases for determining \widehat{p}^* such that $\widehat{p}^* \geq p^*$. The psuedocode for computing the maximum wave speed with \widehat{p}^* is given in Algorithm 1.

5. HIGH ORDER APPROXIMATION OF THE EULER EQUATIONS WITH TABULATED EOS*

In this chapter, we develop a second order method for the Euler equations with a tabulated EOS. This builds on the first order method described in Chapter 4. The foundations of this high order method are based on Guermond et. al. [60] and we will extend the ideas presented there to hold for a tabulated or arbitrary equation of state. This extension is taken from Clayton et. al. [1].

5.1 The Use of the Consistent Mass Matrix

In the low order method we used the *lumped* mass matrix, $(m_i)_{i \in \mathcal{V}}$ as opposed to the consistent mass matrix, $(m_{ij})_{i,j \in \mathcal{V}}$, in order to prove the method was invariant-domain preserving. Furthermore, it was shown in [61], that the use of the consistent mass matrix always violates the maximum principle in a scalar conservation law. However, it is necessary to use the consistent mass matrix for the high order method as using the lumped mass matrix increases dispersion errors. This can be seen in [62]. Additionally, the use of the consistent mass matrix can provide superconvergence effects; see [63]. See also [64]. The high order update is,

$$\frac{1}{\Delta t_n} \sum_{j \in \mathcal{G}(i)} m_{ij} (\mathbf{u}_j^{n+1} - \mathbf{u}_j^n) + \sum_{j \in \mathcal{G}(i)} \mathbb{f}(\mathbf{u}_j^n) \mathbf{c}_{ij} - d_{ij}^{\text{H},n} (\mathbf{u}_j^n - \mathbf{u}_i^n) = \mathbf{0}, \quad (5.1)$$

where we use a higher order graph viscosity, $d_{ij}^{\text{H},n}$. The choice of this graph viscosity is discussed in Section 5.3. A problem with the use of the consistent mass matrix is that it requires numerical matrix solvers to determine the solution. This can increase the computational time even with the use of a preconditioner. We avoid the use of matrix solvers, by using a Neumann series expansion.

* A majority of this chapter is a modification of the work done in [1] and is reprinted with permission from [1].

Proposition 5.1.1 (Neumann Series Expansion). *Let $(X, \|\cdot\|)$ be a Banach space and $T : X \rightarrow X$ a bounded linear operator. If $\|T\| < 1$, then*

$$(I - T)^{-1} = \sum_{n=0}^{\infty} T^n \quad (5.2)$$

where I is the identity map and $T^0 = I$.

The application of the Neumann series to (5.1) can be seen in [65, Sec. 3.4]. We review this application of the Neumann series expansion. Rewrite (5.1) as,

$$\sum_{j \in \mathcal{G}(i)} \frac{m_{ij}}{m_j} \frac{m_j}{\Delta t_n} (\mathbf{U}_j^{n+1} - \mathbf{U}_j^n) + \sum_{j \in \mathcal{G}(i)} \mathbb{f}(\mathbf{U}_j^n) \mathbf{c}_{ij} - d_{ij}^{\text{H},n} (\mathbf{U}_j^n - \mathbf{U}_i^n) = \mathbf{0}, \quad (5.3)$$

Let \mathbb{M} be the matrix with entries $(m_{ij}/m_j)_{i,j \in \mathcal{V}}$. Then the solution, \mathbf{U}_i^{n+1} , can be found inverting \mathbb{M} . By Proposition 5.1.1, we have $\mathbb{M}^{-1} = (\mathbb{I} - (\mathbb{I} - \mathbb{M}))^{-1} = \sum_{n=0}^{\infty} (\mathbb{I} - \mathbb{M})^n$, since. This series converges if $\|\mathbb{I} - \mathbb{M}\| < 1$. Note for small enough mesh size, h , we have that $\sup_{\|\mathbf{x}\|_{\ell^\infty}=1} \|(\mathbb{I} - \mathbb{M})\mathbf{x}\|_{\ell^\infty} < 1$. We approximate this Neumann series expansion by taking only the first two terms of the summation; that is, $\mathbb{M}^{-1} \approx \mathbb{I} + \mathbb{B}$ where $\mathbb{B} = (b_{ij})_{i,j \in \mathcal{V}}$ with $b_{ij} := \delta_{ij} - \frac{m_{ij}}{m_j}$. Next, let $I := \text{card}(\mathcal{V})$ and $\mathbf{R} \in \mathbb{R}^I$. Then we have the following,

$$(\mathbb{M}^{-1}\mathbf{R})_i \approx (\mathbf{R} + \mathbb{B}\mathbf{R})_i = \mathbf{R}_i + \sum_{j \in \mathcal{G}(i)} b_{ij} \mathbf{R}_j = \mathbf{R}_i + \sum_{j \in \mathcal{G}(i)} (b_{ij} \mathbf{R}_j - b_{ji} \mathbf{R}_i), \quad (5.4)$$

Notice that $\sum_{j \in \mathcal{G}(i)} b_{ji} \mathbf{R}_i = 0$ since $\sum_{j \in \mathcal{G}(i)} m_{ji} = m_i$ implies $\sum_{j \in \mathcal{G}(i)} b_{ji} = 0$. Applying this to (5.3) the second order update is given by,

$$\frac{m_i}{\Delta t} (\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^n) = \mathbf{R}_i^n + \sum_{j \in \mathcal{G}(i)} (b_{ij} \mathbf{R}_j^n - b_{ji} \mathbf{R}_i^n), \quad (5.5a)$$

$$\text{where } \mathbf{R}_i^n := \sum_{j \in \mathcal{G}(i)} (-\mathbb{f}(\mathbf{U}_j^n) \mathbf{c}_{ij} + d_{ij}^{\text{H},n} (\mathbf{U}_j^n - \mathbf{U}_i^n)) \quad (5.5b)$$

A more generalized version of the Neumann series approximation is proved by Guermond

& Pasquetti in [62, Sec. 3.1] in the context of transport type problems.

5.2 A Review of the Entropy Solution

In order to define a higher order graph viscosity, we are motivated by the method of the “entropy viscosity commutator” introduced in Guermond et. al. [60, Section 3.4]. The idea is to measure the discrete relative error in the entropy commutator, $\nabla \cdot \mathbf{F}(\mathbf{u}) - (D\eta(\mathbf{u}))^\top \nabla \cdot \mathbf{f}(\mathbf{u})$ where (η, \mathbf{F}) is the entropy pair, see Theorem 5.2.1. We first review some results regarding entropy solutions.

Definition 5.2.1 (Entropy, Entropy-Flux Pair). We say that $(\eta(\mathbf{u}), \mathbf{F}(\mathbf{u}))$ is an *entropy, entropy-flux pair* for the Euler equations, (1.14) if

$$\nabla \cdot \mathbf{F}(\mathbf{u}) = (D\eta(\mathbf{u}))^\top \nabla \cdot \mathbf{f}(\mathbf{u}), \quad \forall \mathbf{u} \in \mathcal{B}(b). \quad (5.6)$$

□

Definition 5.2.2. A weak solution $\mathbf{u} \in [L^\infty(\mathbb{R}^d \times \mathbb{R}_+)]^m$ to the conservation law,

$$\partial_t \mathbf{u} + \nabla \cdot \mathbf{f}(\mathbf{u}) = \mathbf{0}, \quad (5.7)$$

for $\mathbf{g} \in C^1(\mathbb{R}^m; \mathbb{R}^{m \times d})$ is said to be an *entropy solution* if,

$$\partial_t \eta(\mathbf{u}) + \nabla \cdot \mathbf{F}(\mathbf{u}) \geq 0, \quad (5.8)$$

holds in the sense of distributions for every concave entropy $\eta : \mathbb{R}^m \rightarrow \mathbb{R}$ and associated entropy-flux, $\mathbf{F} : \mathbb{R}^m \rightarrow \mathbb{R}^d$. □

Theorem 5.2.1 (Smooth Entropy Solutions). *Every smooth solution to the compressible Euler equations, (1.14), is also an entropy solution.*

Remark 5.2.1 (Mathematical Entropy). Definition 5.2.2 can be equivalently framed by assuming $\partial_t \eta(\mathbf{u}) + \nabla \cdot \mathbf{F}(\mathbf{u}) \leq 0$ holds in the sense of distributions for η convex. See [66].

For numerical purposes, the distinction is not necessary as we will be measuring the discrete error of the quantity, $\nabla \cdot \mathbf{F}(\mathbf{u}) - (D\eta(\mathbf{u}))^\top \nabla \cdot \mathbf{f}(\mathbf{u})$. \square

For the compressible Euler equations with ideal gas law, it is unknown if there is a characterization of all the entropy pairs (η, \mathbf{F}) like the class of Kruzhkov entropies for scalar conservation laws, see Kruzhkov [29]. However, we will exploit a class of known entropy, entropy-flux pairs. They are the so-called *generalized entropy pairs*, defined by $\eta(\mathbf{u}) = -\rho f(\mathbf{s}(\mathbf{u}))$ where $f \in C^2(\mathbb{R})$, $\mathbf{F}(\mathbf{u}) = \frac{\mathbf{m}}{\rho} \eta(\mathbf{u})$, where $f'(s) > 0$, $f'(s)c_p^{-1} - f''(s) > 0$ for all $(\rho, e) \in \mathcal{A}$. In this case η is convex (see [67] and [66]).

The simplest non-trivial entropy pair is when $f(s) = 1$; that is, $\eta(\mathbf{u}) = -\rho$ and $\mathbf{F}(\mathbf{u}) = \mathbf{m}$. However, the entropy, η is not strictly convex. For $f(s) = s$, $\eta(\mathbf{u}) = -\rho \mathbf{s}(\mathbf{u})$ which is the mathematical entropy (the negative of the physical entropy). Another example is the so-called *Harten entropy* given by $\eta(\mathbf{u}) = \rho \exp(\frac{\gamma-1}{\gamma+\alpha} \mathbf{s}(\mathbf{u}))$ for $\alpha > 0$. This was shown for the ideal EOS in [67]; however, it also holds for the covolume and NASG as $c_p = \frac{\gamma}{\gamma-1}$ for each of these EOS. It should be noted that these examples are not entropy pairs for every equation of state. Notably, if the equation of state is non-convex, then the entropies proposed are not convex.

If the interpolatory EOS is covolume, then the specific entropy is given by, $\mathbf{s}(\mathbf{u}) := \log((\mathbf{e}(\mathbf{u})^{\frac{1}{\gamma-1}}(\rho^{-1} - b)))$. Then the Harten entropy can be written in the form,

$$\eta(\mathbf{u}) = \left(\frac{\rho^{\alpha+1} \mathbf{e}(\mathbf{u})}{(1 - b\rho)^{1-\gamma}} \right)^{\frac{1}{\gamma+\alpha}}. \quad (5.9)$$

Note that a nice choice for the Harten entropy is to set $\alpha = 1$, since $\rho^2 \mathbf{e}(\mathbf{u}) = \rho E - \frac{1}{2} \|\mathbf{m}\|_{\ell^2}^2$, which is a quadratic function of the conserved variables.

5.3 The Entropy Viscosity Method

Motivated by notion of the entropy commutator, $\nabla \cdot \mathbf{F}(\mathbf{u}) - (D\eta(\mathbf{u}))^\top \nabla \cdot \mathbf{f}(\mathbf{u})$, we propose to find a local entropy $\eta^{n,i}$ at every node $i \in \mathcal{V}$ and every time t_n . We define the local flux,

$$\mathbf{f}^{n,i}(\mathbf{u}) := \begin{pmatrix} \mathbf{m} \\ \mathbf{v} \otimes \mathbf{m} + \rho_{\text{nasg}}^{n,i}(\mathbf{u}) \mathbb{I}_d \\ \mathbf{v}(E + \rho_{\text{nasg}}^{n,i}(\mathbf{u})) \end{pmatrix} \quad (5.10)$$

where $\rho_{\text{nasg}}^{n,i}(\mathbf{u}) := (\gamma_i^{\min,n} - 1) \left(\frac{\rho(\mathbf{e}(\mathbf{u}) - q)}{1 - b\rho} - p_\infty \right) - p_\infty$, $\gamma_i^{\min,n} := \min_{j \in \mathcal{G}(i)} \gamma_j^n$, and

$$\gamma_j^n := \frac{(p_j^n + p_\infty)(1 - b\rho_j^n)}{\rho_j^n(\mathbf{e}(\mathbf{U}_j^n) - q) - (1 - b\rho_j^n)p_\infty} + 1. \quad (5.11)$$

Note that, $\rho_{\text{nasg}}^{n,i}$ is the NASG EOS with $\gamma = \gamma_i^{\min,n}$. Then two possible entropy pairs for the Euler equations with flux $\mathbf{f}^{n,i}$ are,

$$\eta_1^{i,n}(\mathbf{u}) := -\rho \log \left((e(\mathbf{u}) - q - p_\infty(\rho^{-1} - b)) \frac{1}{\gamma_i^{\min,n-1}} (\rho^{-1} - b) \right) - \frac{\rho}{\rho_i^n} \eta_{\text{ref},1}^{i,n}, \quad (5.12a)$$

$$\mathbf{F}_1^{n,i}(\mathbf{u}) := \mathbf{v} \eta_1^{n,i}(\mathbf{u}), \quad (5.12b)$$

where $\eta_{\text{ref},1}^{i,n} := -\rho_i^n \log \left((e_i^n - q - p_\infty((\rho_i^n)^{-1} - b)) \frac{1}{\gamma_i^{\min,n-1}} ((\rho_i^n)^{-1} - b) \right)$ and

$$\eta_2^{n,i}(\mathbf{u}) := \left(\frac{\rho^{\alpha+1}(\mathbf{e}(\mathbf{u}) - q) - p_\infty \rho^\alpha (1 - b\rho)}{(1 - b\rho)^{1 - \gamma_i^{\min,n}}} \right) \frac{1}{\gamma_i^{\min,n+\alpha}} - \frac{\rho}{\rho_i^n} \eta_{\text{ref},2}^{n,i}, \quad (5.13a)$$

$$\mathbf{F}_2^{n,i}(\mathbf{u}) := \mathbf{v} \eta_2^{n,i}(\mathbf{u}), \quad (5.13b)$$

where $\eta_{\text{ref},2}^{n,i} = \left(\frac{(\rho_i^n)^{\alpha+1}(\mathbf{e}(\mathbf{U}_i^n) - q) - p_\infty (\rho_i^n)^\alpha (1 - b\rho_i^n)}{(1 - b\rho_i^n)^{1 - \gamma_i^{\min,n}}} \right) \frac{1}{\gamma_i^{\min,n+\alpha}}$, for $\alpha > 0$. Note that each entropy has been shifted by $\rho \eta_{\text{ref},k}^{n,i} / \rho_i^n$ so that $\eta_k^{n,i}(\mathbf{U}_i^n) = 0$ for $k = 1, 2$. This shift is valid since $\eta(\mathbf{u}) = \rho(f(\mathbf{s}(\mathbf{u})) - c)$ is also an entropy for any constant c with the corresponding entropy-flux, $\mathbf{F}(\mathbf{u}) := \mathbf{m}(f(\mathbf{s}(\mathbf{u})) - c)$.

Remark 5.3.1 (Remarks on $\eta^{n,i}$ and $\gamma_i^{\min,n}$). Note that $\eta_2^{n,i}$ is the (shifted) Harten entropy. The choice of $\alpha = 1$ is also convenient as $\rho^2(\mathbf{e}(\mathbf{u}) - q) - p_\infty \rho(1 - b\rho) = \rho E - \frac{1}{2} \|\mathbf{m}\|_{\ell^2}^2 - p_\infty \rho(1 - b\rho)$ which is a quadratic function of the conserved variables; this simplifies the computation of $D\eta^{n,i}$.

In the definition of the entropy pairs, (5.12) and (5.13), we could have alternatively used $\gamma_i^{\max,n}$ or simply γ_i^n , as each of these choices will recover the expected Harten entropy if the pressure is given by the ideal, covolume, or NASG equations of state. We choose $\gamma_i^{\min,n}$ as it is already necessary to compute for the limiting of the surrogate physical entropy, see Theorem 6.4.3. \square

We may select either of the two entropy, entropy-flux pairs, so we drop the subscript notation for $\eta_k^{i,n}$. With this collection of entropy pairs, $\{(\eta^{n,i}, \mathbf{F}^{n,i})\}_{i \in \mathcal{V}}$ for $n \in \mathbb{N}$, we can measure the local error in the entropy viscosity commutator by approximating

$$\int_{\Omega} ((\nabla \cdot \mathbf{F}^{n,i}(\mathbf{u}) - (D\eta^{n,i}(\mathbf{u}))^\top \nabla \cdot \mathbf{f}^{n,i}(\mathbf{u})) \varphi_i(\mathbf{x}) \, d\mathbf{x} \quad (5.14)$$

with our numerical solution $\mathbf{u}_h^n = \sum_{i \in \mathcal{V}} \mathbf{U}_i^n \varphi_i$. As mentioned in Section 3.2, the divergence of the flux is approximated by, $\nabla \cdot \mathbf{f}^{n,i}(\mathbf{u}_h^n) \approx \sum_{i \in \mathcal{V}} \mathbf{f}^{n,i}(\mathbf{U}_i^n) \nabla \varphi_i$. Similarly, for the entropy flux, $\nabla \cdot \mathbf{F}^{n,i}(\mathbf{u}_h^n) \approx \sum_{i \in \mathcal{V}} \mathbf{F}^{n,i}(\mathbf{U}_i^n) \cdot \nabla \varphi_i$. The discretization of (5.14) is thus given by,

$$N_i^n := \sum_{j \in \mathcal{G}(i)} (\mathbf{F}^{n,i}(\mathbf{U}_j^n) - (D\eta(\mathbf{U}_i^n))^\top \mathbf{f}^{n,i}(\mathbf{U}_j^n)) \cdot \mathbf{c}_{ij}. \quad (5.15)$$

We now define the *entropy residual* as

$$R_i^n := \frac{|N_i^n|}{D_i^n + \varepsilon \max_{j \in \mathcal{V}} D_j^n + \epsilon}, \quad \varepsilon = 10^{-1}, \quad \epsilon = 10^{-14}, \quad (5.16)$$

$$D_i^n := \sum_{j \in \mathcal{G}(i)} |\mathbf{F}^{n,i}(\mathbf{U}_j^n) \cdot \mathbf{c}_{ij}| + \sum_{j \in \mathcal{G}(i)} |(D\eta^{n,i}(\mathbf{U}_i^n))^\top \mathbf{f}(\mathbf{U}_i^n) \cdot \mathbf{c}_{ij}|. \quad (5.17)$$

The parameter ϵ is a machine precision parameter used to avoid division by zero.

Remark 5.3.2 (Concave/Convex Entropy). Note that there is no assumption made on the concavity or convexity of $\eta^{n,i}(\mathbf{u})$ need not be a concave entropy. We only need that (η, \mathbf{F}) be an entropy, entropy-flux pair so as to measure the smoothness of the solution. Therefore a non-convex entropy is valid in the construction of the entropy residual. \square

The high order viscosity is then defined as,

$$d_{ij}^{\text{H},n} := \max(R_i^n, R_j^n) d_{ij}^{\text{L},n}. \quad (5.18)$$

Another alternative for the definition can be,

$$d_{ij}^{\text{H},n} := \frac{R_i^n + R_j^n}{2} d_{ij}^{\text{L},n}. \quad (5.19)$$

Remark 5.3.3 (Thresholding Function). One can further emphasize the result of the residual. That is, define the function,

$$\psi(x) := \frac{4x_0^3 - (x + x_0)(x - 2x_0)((x - 2x_0) - \text{ReLU}(x - 2x_0))}{4x_0^3} \quad (5.20)$$

where $\text{ReLU}(x) = (x + |x|)/2$ and $x_0 \in (0, 0.5]$ parameter. Note that ψ satisfies, $\psi(0) = 0$, $\psi(x_0) = \frac{1}{2}$, and $\psi(x) = 1$ for all $x \in [2x_0, 1]$. The fixed point, $\psi(x^*) = x^*$ is computed by $x^* = x_0(\frac{3}{2} - \frac{1}{2}\sqrt{9 - 16x_0})$. Thus we can alternatively define the high order viscosity as,

$$d_{ij}^{\text{H},n} = \max(\psi(R_i^n), \psi(R_j^n)) d_{ij}^{\text{L},n}. \quad (5.21)$$

or using an average of the residuals,

$$d_{ij}^{\text{H},n} := \frac{1}{2}(\psi(R_i^n) + \psi(R_j^n)) d_{ij}^{\text{L},n}. \quad (5.22)$$

\square

To conclude, the high order viscosity, $d_{ij}^{H,n}$ detects whether the solution is not smooth and if so, takes on values close to the low order artificial viscosity.

6. QUASICONCAVE LIMITING*

Since the second order method is **not** invariant-domain preserving, we must perform a limiting process on the solution to keep it in the invariant domain. The concept of limiting was originally developed by Boris & Book, named the Flux Corrected Transport (FCT) method. This development took course over the three papers [68, 69, 70] and was also extended to multiple dimensions Zalesak in [71]. Numerous work has been done on this methodology; for more resources, see [72].

For this chapter, the limiting is done with respect to some local bounds (see Section 6.1) that are satisfied by the first order method. This limiting process is performed on quasiconcave functionals which was first developed by Guermond et. al. in [60]. We justify in Section 6.3 that limiting with respect to these local bounds preserves the invariant domain properties.

6.1 Local Bounds

Since the high order update is known to violate the invariant domain properties, we need to limit our solution so that it is no longer non-physical. For example, if $\rho_i^{\text{H},n+1} < 0$, then we would like to somehow, pull this high order solution closer to the low order solution. That is, a mapping, $\mathcal{L} : \rho_i^{\text{H},n+1} \mapsto \rho_i^{n+1}$, such that $\rho_i^{n+1} > 0$. The challenge is to do this globally, efficiently and only when necessary. The technique we employ is based on local bounds that are satisfied by the low order solution, which consequently preserve the invariant-domain properties.

Theorem 6.1.1. *The low order update, $\mathbf{U}_i^{\text{L},n+1}$, defined in (4.2), satisfies the following local*

* A majority of this chapter is a modification of the work done in [1] and is reprinted with permission from [1].

bounds,

$$\rho_i^{\min,n} := \min_{j \in \mathcal{G}(i)} (\bar{\rho}_{ij}^n) \leq \rho_i^{\text{L},n+1} \leq \max_{j \in \mathcal{G}(i)} (\bar{\rho}_{ij}^n) =: \rho_i^{\max,n} \quad (6.1)$$

$$\mathbf{E}_i^{\min,n} := \min_{j \in \mathcal{G}(i)} (\bar{\mathbf{E}}_{ij}^n) \leq \mathbf{E}_i^{\text{L},n+1} \leq \max_{j \in \mathcal{G}(i)} (\bar{\mathbf{E}}_{ij}^n) =: \mathbf{E}_i^{\max,n} \quad (6.2)$$

$$\Upsilon_i^{\min} := \min_{j \in \mathcal{G}(i)} (\rho \mathbf{e})(\bar{\mathbf{U}}_{ij}^n) \leq (\rho \mathbf{e})(\mathbf{U}_i^{\text{L},n+1}) \quad (6.3)$$

for all $i \in \mathcal{V}$.

Proof. The result for (6.1)–(6.2) immediately follows from the definition of the low order states as a convex combination of bar state and that $\bar{\rho}_{ij}^n > 0$ and $\bar{\mathbf{E}}_{ij}^n > 0$ for all $i \in \mathcal{V}$ and $j \in \mathcal{G}(i)$. The internal energy minimum, (6.3), is due to the fact that internal energy is a concave function of the conservative variables, therefore,

$$(\rho \mathbf{e})(\mathbf{U}_i^{\text{L},n+1}) = (\rho \mathbf{e})\left(\sum_{j \in \mathcal{G}(i)} \alpha_j \bar{\mathbf{U}}_{ij}^n\right) \geq \sum_{j \in \mathcal{G}(i)} \alpha_j (\rho \mathbf{e})(\bar{\mathbf{U}}_{ij}^n) \geq \min_{j \in \mathcal{G}(i)} (\rho \mathbf{e})(\bar{\mathbf{U}}_{ij}^n), \quad (6.4)$$

where α_j are the coefficients of the bar states shown in (4.2). \square

Remark 6.1.1 (Local Upper Bound on the Density). Note from Theorem 6.1.1, if the oracle has a maximum compression constant, b^{-1} ; for example, the van der Waals or covolume EOS, then we also have that $0 < 1 - b\rho_i^{\max,n} \leq 1 - b\rho_i^{\text{L},n+1} \leq 1 - b\rho_i^{\min,n}$. This follows since $\mathcal{B}(b)$ is an invariant domain. \square

6.1.1 Relaxation on the Density Bounds

When performing limiting on the density as described in Section 6.2 or Section 6.3.4 the high order solution can actually be reduced to first order. The methodology for correcting this issue to loosen the local bounds (6.1). That is, the solution after the limiting process can violate these bounds. This relaxation of the bounds is done so that the error is second order and still preserves $\rho > 0$ and $1 - b\rho > 0$.

First, we have the following lemma in regards to the maximum value of ρ in the solution

to the Riemann problem, (4.8).

Lemma 6.1.1 (Maximum Density Bound). *The following is true.*

1. *The density in the solution the extended Riemann problem (4.8) satisfies the following*

$$\rho \leq \max_{Z \in \{i,j\}} \frac{1}{\tau_Z^\infty} = \max_{Z \in \{i,j\}} \frac{(\gamma_Z + 1)\rho_Z}{(\gamma_Z - 1) + 2b\rho_Z} \quad (6.5)$$

2. *Under the CFL condition stated in Theorem 4.1.1, the low order update satisfies*

$$\rho_i^{\mathbf{L},n+1} \leq \frac{(\gamma_i^{\min,n} + 1)\rho_i^{\max,n}}{(\gamma_i^{\min,n} - 1) + 2b\rho_i^{\max,n}} \quad (6.6)$$

Proof. 1. Assume the Z wave is an expansion, then the density decreases across the expansion, hence $\rho \leq \rho_Z$. If the Z wave is a shock, then from Lemma 6.4.1, we know that $\tau \in (\tau_Z^\infty, \tau_Z]$. Hence $\rho \in (\frac{1}{\tau_Z}, \frac{1}{\tau_Z^\infty}]$. This completes the proof.

2. From (6.1) we have that

$$\rho_i^{\mathbf{L},n+1} \leq \max_{j \in \mathcal{G}(i)} \bar{\rho}_{ij}^n \leq \max_{j \in \mathcal{G}(i)} \frac{(\gamma_j + 1)\rho_j}{(\gamma_j - 1) + 2b\rho_j} \quad (6.7)$$

where we have also applied (6.5) to $\bar{\rho}_{ij}^n$ as the average of the solution to the Riemann problem must also satisfy that inequality. Since $\frac{(\gamma+1)\rho}{(\gamma-1)+2b\rho}$ is an increasing function of ρ and a decreasing function of γ , we conclude that,

$$\rho_i^{\mathbf{L},n+1} \leq \frac{(\gamma_i^{\min,n} + 1)\rho_i^{\max,n}}{(\gamma_i^{\min,n} - 1) + 2b\rho_i^{\max,n}}. \quad (6.8)$$

□

Next we introduce an approximation of the local curvature of the density which will also

be used in the relaxation. For each $i \in \mathcal{V}$, define the following,

$$\Delta^2 \rho_i^n := \frac{\sum_{j \in \mathcal{G}(i) \setminus \{i\}} \beta_{ij} (\rho_i^n - \rho_j^n)}{\sum_{j \in \mathcal{G}(i) \setminus \{i\}} \beta_{ij}} \quad (6.9a)$$

$$\overline{\Delta^2 \rho_i^n} := \frac{1}{2 \text{card}(\mathcal{G}(i))} \sum_{j \in \mathcal{G}(i) \setminus \{i\}} \left(\frac{1}{2} \Delta^2 \rho_i^n + \frac{1}{2} \Delta^2 \rho_j^n \right), \quad (6.9b)$$

where $\beta_{ij} := \int_D \nabla \varphi_i \cdot \nabla \varphi_j \, d\mathbf{x}$ are the stiffness coefficients of the Laplace operator. Recall that φ_i are the global shape functions; see Chapter 3. These definitions are well defined since from the partition of unity property, we have that $\sum_{j \in \mathcal{G}(i) \setminus \{i\}} \beta_{ij} = -\beta_{ii} = -\int_D \|\nabla \varphi_i\|_{\ell^2}^2 \, d\mathbf{x} \neq 0$. Notice that $\overline{\Delta^2 \rho_i^n}$ is an estimate of the local curvature in a neighborhood of node i .

The relaxation on the density is now defined as,

$$\widetilde{\rho_i^{\min, n}} := \max((1 - r_h) \rho_i^{\min, n}, \rho_i^{\min, n} - \overline{\Delta^2 \rho_i^n}), \quad (6.10a)$$

$$\widetilde{\rho_i^{\max, n}} := \min\left((1 + r_h) \rho_i^{\max, n}, \rho_i^{\max, n} + \overline{\Delta^2 \rho_i^n}, \frac{(\gamma+1) \rho_i^{\max, n}}{\gamma-1+2b\rho_i^{\max, n}}\right), \quad (6.10b)$$

where $r_h := (m_i/|D|)^{1.5/d}$.

6.2 The Flux Corrected Transport Method

To demonstrate the FCT method, we apply it to the update on the density, ρ . Let the low order and high order update be given as $\rho_i^{\text{L}, n+1}$ and $\rho_i^{\text{H}, n+1}$, respectively. These two quantities are computed using first and second order methods described in Chapters 4 and 5, respectively. Note that the high order update is related to the low order update by,

$$\rho_i^{\text{H}, n+1} = \rho_i^{\text{L}, n+1} + \sum_{j \in \mathcal{G}(i)} \mathbf{P}_{ij}^\rho. \quad (6.11)$$

From Section 6.1 we know that the first order update satisfies, $\rho_i^{\min} \leq \rho_i^{\text{L}, n+1} \leq \rho_i^{\max}$ and the second order update does not necessarily satisfy these local bounds. The FCT method begins by splitting $\sum_{j \in \mathcal{G}(i)} \mathbf{P}_{ij}^\rho$ into the negative and positive parts. That is, let $\mathcal{G}(i)^+ := \{j \in \mathcal{G}(i) : \mathbf{P}_{ij}^\rho > 0\}$ and $\mathcal{G}(i)^- := \{j \in \mathcal{G}(i) : \mathbf{P}_{ij}^\rho < 0\}$. Then the goal is to find $\ell_i^+, \ell_i^- \in [0, 1]$

so that new update defined by,

$$\rho_i^{n+1} = \rho_i^{L,n+1} + \ell_i^+ \sum_{j \in \mathcal{G}(i)^+} P_{ij}^\rho + \ell_i^- \sum_{j \in \mathcal{G}(i)^-} P_{ij}^\rho, \quad (6.12)$$

satisfies, $\rho_i^{\min} \leq \rho_i^{n+1} \leq \rho_i^{\max}$.

Note that for any $\ell_i^+, \ell_i^- \in [0, 1]$ the new update satisfies,

$$\rho_i^{L,n+1} + \ell_i^- \sum_{j \in \mathcal{G}(i)^-} P_{ij}^\rho \leq \rho_i^{n+1} \leq \rho_i^{L,n+1} + \ell_i^+ \sum_{j \in \mathcal{G}(i)^+} P_{ij}^\rho. \quad (6.13)$$

Thus under the choice of,

$$\ell_i^- := \min \left(\frac{\rho_i^{\min} - \rho_i^{L,n+1}}{\sum_{j \in \mathcal{G}(i)^-} P_{ij}^\rho}, 1 \right) \quad \text{and} \quad \ell_i^+ := \min \left(\frac{\rho_i^{\max} - \rho_i^{L,n+1}}{\sum_{j \in \mathcal{G}(i)^+} P_{ij}^\rho}, 1 \right), \quad (6.14)$$

we have that $\rho_i^{\min,n} \leq \rho_i^{n+1} \leq \rho_i^{\max,n}$. For more recent work in the FCT literature, in the context of the compressible Euler equations, see [73].

Remark 6.2.1 (Limitations of the FCT Method). The FCT method is inherently linear. This is perfectly fine for preserving positivity of the density; however, if one attempts to apply it to the specific internal energy, $\mathbf{e}(\mathbf{U}_i)$, then difficulties immediately arise as \mathbf{e} is a nonlinear functional. \square

6.3 Quasiconcave Limiting

In this section we discuss a new method introduced in [60] which allows us to perform limiting on quasiconcave functionals. This methodology allows us to address the issue of maintaining positivity of the internal energy.

6.3.1 Quasiconcave Functionals

Definition 6.3.1 (Quasiconcavity). For \mathcal{C} a convex set, $\Psi : \mathcal{C} \rightarrow \mathbb{R}$ is said to be quasiconcave if $L_\lambda(\Psi) := \{\mathbf{U} \in \mathcal{C} : \Psi(\mathbf{U}) \geq \lambda\}$ is a convex set for any $\lambda \in \mathbb{R}$. That is, the upper level sets are convex. \square

As might be inferred from the name “quasiconcave”, a functional which is quasiconcave may not be concave. For example, if $\mathcal{C} = [-1, 1]$ then $f(x) = x^3$ is certainly not concave, but it is quasiconcave since for any $\lambda \in \text{Range}(f)$, we have, $L_\lambda(f) = [\sqrt[3]{\lambda}, 1]$ which is a convex set. However, every concave functional is also quasiconcave.

Lemma 6.3.1 (Concave Implies Quasiconcave). *Let $\mathcal{C} \subset \mathbb{R}^m$ be convex and $\Psi : \mathcal{C} \rightarrow \mathbb{R}$ be a concave functional, then Ψ is quasiconcave.*

Proof. Assume that λ is in the range of Ψ ; that is, $L_\lambda(\Psi) \neq \emptyset$, otherwise the result is vacuous. Let $\mathbf{U}_1, \mathbf{U}_2 \in L_\lambda(\Psi)$ and $0 \leq t \leq 1$. Then since Ψ is concave, we have,

$$\Psi(t\mathbf{U}_1 + (1-t)\mathbf{U}_2) \geq t\Psi(\mathbf{U}_1) + (1-t)\Psi(\mathbf{U}_2) \geq \lambda. \quad (6.15)$$

Therefore $t\mathbf{U}_1 + (1-t)\mathbf{U}_2 \in L_\lambda(\Psi)$ for all $t \in [0, 1]$ hence $L_\lambda(\Psi)$ is convex and so Ψ is quasiconcave. \square

6.3.2 The Abstract Scheme

Similar to the FCT method, we estimate the difference in the low and high order methods, $\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^{\text{L},n+1}$. Specifically, we look at the difference, $\sum_{j \in \mathcal{G}(i)} m_{ij}(\mathbf{U}_j^{\text{H},n+1} - \mathbf{U}_j^n) - (\mathbf{U}_i^{\text{L},n+1} - \mathbf{U}_i^n)$. Substituting in the definition of the first and second order updates, (4.1) and (5.5), respectively, we have

$$\sum_{j \in \mathcal{G}(i)} m_{ij}(\mathbf{U}_j^{\text{H},n+1} - \mathbf{U}_j^n) - m_i(\mathbf{U}_i^{\text{L},n+1} - \mathbf{U}_i^n) = \Delta t_n \sum_{j \in \mathcal{G}(i)} (d_{ij}^{\text{H},n} - d_{ij}^{\text{L},n})(\mathbf{U}_j^n - \mathbf{U}_i^n). \quad (6.16)$$

Now, add and subtract $m_i(\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^n)$ and refactor as,

$$\begin{aligned} & m_i(\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^{\text{L},n+1}) - m_i(\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^n) + \sum_{j \in \mathcal{G}(i)} m_{ij}(\mathbf{U}_j^{\text{H},n+1} - \mathbf{U}_j^n) \\ &= \Delta t_n \sum_{j \in \mathcal{G}(i)} (d_{ij}^{\text{H},n} - d_{ij}^{\text{L},n})(\mathbf{U}_j^n - \mathbf{U}_i^n). \end{aligned} \quad (6.17)$$

This can all be combined into a nice expression by defining $\Delta_{ij} := m_i \delta_{ij} - m_{ij}$ where $\delta_{ij} = 1$ if $i = j$ and 0 if $i \neq j$. Thus we have,

$$m_i(\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^{\text{L},n+1}) = \sum_{j \in \mathcal{G}(i)} \Delta_{ij}(\mathbf{U}_j^{\text{H},n+1} - \mathbf{U}_j^n) + \Delta t_n (d_{ij}^{\text{H},n} - d_{ij}^{\text{L},n})(\mathbf{U}_j^n - \mathbf{U}_i^n). \quad (6.18)$$

It is important to note that the right hand side of this equation can be expressed as a *skew-symmetric* matrix. Using that $\sum_{j \in \mathcal{G}(i)} \Delta_{ij} = 0$, we can write the following,

$$m_i(\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^{\text{L},n+1}) = \sum_{j \in \mathcal{G}(i)} \mathbf{A}_{ij}^n \quad (6.19a)$$

$$\mathbf{A}_{ij}^n := \Delta_{ij}(\mathbf{U}_j^{\text{H},n+1} - \mathbf{U}_j^n - (\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^n)) + \Delta t_n (d_{ij}^{\text{H},n} - d_{ij}^{\text{L},n})(\mathbf{U}_j^n - \mathbf{U}_i^n), \quad (6.19b)$$

from which it is readily seen to be skew-symmetric ($\mathbf{A}_{ij}^n = -\mathbf{A}_{ji}^n$).

Lemma 6.3.2 (Consistent Conservation). *The total mass of the high order scheme is the same as the low order scheme; that is,*

$$\sum_{i \in \mathcal{V}} m_i \mathbf{U}_i^{\text{H},n+1} = \sum_{i \in \mathcal{V}} m_i \mathbf{U}_i^{\text{L},n+1}. \quad (6.20)$$

Proof. The result is proven if we can show that

$$\sum_{i \in \mathcal{V}} m_i(\mathbf{U}_i^{\text{H},n+1} - \mathbf{U}_i^{\text{L},n+1}) = \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{G}(i)} \mathbf{A}_{ij}^n = 0. \quad (6.21)$$

Using that \mathbf{A}_{ij}^n is skew-symmetric, we have, $\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{G}(i)} \mathbf{A}_{ij}^n = \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{G}(i)} -\mathbf{A}_{ji}^n$. Using Fubini's theorem, we find that, $\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{G}(i)} \mathbf{A}_{ij}^n + \sum_{j \in \mathcal{V}} \sum_{i \in \mathcal{G}(j)} \mathbf{A}_{ji}^n = 0$. But the indices are arbitrary hence, $2 \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{G}(i)} \mathbf{A}_{ij}^n = 0$ and the result is proven. \square

6.3.3 The Limiter

The idea now, is to introduce a symmetric matrix (ℓ_{ij}) which provides an “intermediate” solution between $\mathbf{U}_i^{\text{H},n+1}$ and $\mathbf{U}_i^{\text{L},n+1}$. In particular we define,

$$\mathbf{U}_i^{n+1} := \mathbf{U}_i^{\text{L},n+1} + \frac{1}{m_i} \sum_{j \in \mathcal{G}(i)} \ell_{ij} \mathbf{A}_{ij}^n. \quad (6.22)$$

where $0 \leq \ell_{ij} \leq 1$ for all $i, j \in \mathcal{V}$ and (ℓ_{ij}) is symmetric. The reason for enforcing symmetry on (ℓ_{ij}) , is that $\ell_{ij} \mathbf{A}_{ij}^n$ is still skew-symmetric, hence the total mass of \mathbf{U}_i^{n+1} is the same as $\mathbf{U}_i^{\text{L},n+1}$. Note, if $\ell_{ij} = 0$ for all $j \in \mathcal{G}(i)$, then $\mathbf{U}_i^{n+1} = \mathbf{U}_i^{\text{L},n+1}$ and similarly, if $\ell_{ij} = 1$ for all $j \in \mathcal{G}(i)$, then $\mathbf{U}_i^{n+1} = \mathbf{U}_i^{\text{H},n+1}$. All of the discussion so far on the limiter is still a part of the FCT method as described in, [71].

We now move towards the quasiconcave limiting theory as described in [60, Sec. 4.2]. Introduce $\{\lambda_j^i\}_{j \in \mathcal{G}(i) \setminus \{i\}}$ such that $\sum_{j \in \mathcal{G}(i) \setminus \{i\}} \lambda_j^i = 1$ and $\lambda_j^i > 0$ for all $j \in \mathcal{G}(i)$. The superscript i is just used to indicate the dependence on i . Then, (6.22) can be written as,

$$\mathbf{U}_i^{n+1} = \sum_{j \in \mathcal{G}(i) \setminus \{i\}} \lambda_j^i \left(\mathbf{U}_i^n + \ell_{ij} \mathbf{P}_{ij}^n \right), \quad \text{where } \mathbf{P}_{ij}^n := \frac{1}{m_i \lambda_j^i} \mathbf{A}_{ij}^n. \quad (6.23)$$

Note that $j = i$ is omitted since $\mathbf{A}_{ii}^n = \mathbf{0}$ for all $i \in \mathcal{V}$ since \mathbf{A}_{ij}^n is skew-symmetric. For the remainder of this thesis, we set $\lambda_j^i = \frac{1}{\text{card}(\mathcal{G}(i)) - 1}$.

Recall that we would like to preserve some physical quantity represented by a quasiconcave functional, Ψ . Therefore, the goal is to determine ℓ_{ij} so that $\Psi(\mathbf{U}_i^{n+1}) \geq 0$. The following two lemmas set the stage for the upcoming numerical algorithms.

Lemma 6.3.3. *Let $\Psi : \mathcal{C} \rightarrow \mathbb{R}$ be a quasiconcave functional. Assume $\ell_{ij} \in [0, 1]$ are such that $\Psi(\mathbf{U}_i^{\text{L},n} + \ell_{ij} \mathbf{P}_{ij}) \geq 0$ for all $j \in \mathcal{G}(i) \setminus \{i\}$, then*

$$\Psi \left(\sum_{j \in \mathcal{G}(i) \setminus \{i\}} \lambda_j^i (\mathbf{U}_i^{\text{L},n} + \ell_{ij} \mathbf{P}_{ij}) \right) \geq 0. \quad (6.24)$$

Proof. The proof of this follows from the definition of quasiconcavity. Note that $\mathbf{U}_i^{\mathbf{L},n+1} + \ell_{ij}\mathbf{P}_{ij} \in L_0(\Psi)$ where $L_0(\Psi) = \{\mathbf{U} \in \mathcal{C} : \Psi(\mathbf{U}) \geq 0\}$. Since $L_0(\Psi)$ is a convex set, then the convex combination defined in (6.23) must also belong to $L_0(\Psi)$, hence the result. \square

Lemma 6.3.4 (Symmetrization of ℓ_{ij}). *Define,*

$$\ell_j^i := \begin{cases} 1 & \text{if } \Psi(\mathbf{U}_i^{\mathbf{L},n+1} + \mathbf{P}_{ij}) \geq 0 \\ \max\{\ell \in [0, 1] : \Psi(\mathbf{U}_i^{\mathbf{L},n+1} + \ell\mathbf{P}_{ij}) \geq 0\} & \text{otherwise} \end{cases} \quad (6.25)$$

for all $i \in \mathcal{V}$ and $j \in \mathcal{G}(i)$. Then the following holds,

1. $\Psi(\mathbf{U}_i^{\mathbf{L},n+1} + \ell\mathbf{P}_{ij}) \geq 0$ for all $\ell \in [0, \ell_j^i]$.
2. Let $\ell_{ij} := \min\{\ell_j^i, \ell_i^j\}$, then $\Psi(\mathbf{U}_i^{\mathbf{L},n+1} + \ell_{ij}\mathbf{P}_{ij}) \geq 0$ for all $i \in \mathcal{V}$ and all $j \in \mathcal{G}(i)$.

Remark 6.3.1 (Comments on the Limiter ℓ_{ij}). From Lemma 6.3.3 we see that the determination of ℓ_{ij} only needs to be found for each $j \in \mathcal{G}(i)$ independently. From Lemma 6.3.4 (i) we see that smaller values of the limiter do not violate the non-negativity of the quasiconcave functional. This is intuitively expected, since we are “pulling” our solution closer to the invariant-domain preserving update, $\mathbf{U}_i^{\mathbf{L},n+1}$. Lastly, from Lemma 6.3.4 (ii) we see that (ℓ_{ij}) is symmetric. \square

6.3.4 Limiting on the Density

Limiting on the density is performed exactly as done in [60, Sec 4.4]. We define, $\Psi_+^\rho(\mathbf{U}) := \rho - \rho_i^{\min,n}$ and $\Psi_-^\rho(\mathbf{U}) := \rho_i^{\max,n} - \rho$. The limiter is then,

$$\ell_j^{i,\rho} := \begin{cases} \min\left(\frac{|\rho_i^{\min,n} - \rho_i^{\mathbf{L},n+1}|}{|\mathbf{P}_{ij}^\rho| + \epsilon_i}, 1\right) & \text{if } \rho_i^{\mathbf{L},n+1} + \mathbf{P}_{ij}^\rho < \rho_i^{\min,n}, \\ 1 & \rho_i^{\min,n} \leq \rho_i^{\mathbf{L},n+1} + \mathbf{P}_{ij}^\rho \leq \rho_i^{\max,n}, \\ \min\left(\frac{|\rho_i^{\max,n} - \rho_i^{\mathbf{L},n+1}|}{|\mathbf{P}_{ij}^\rho| + \epsilon_i}, 1\right) & \text{if } \rho_i^{\max,n} < \rho_i^{\mathbf{L},n+1} + \mathbf{P}_{ij}^\rho, \end{cases} \quad (6.26)$$

where $\epsilon_i := \epsilon \rho_i^{\max, n}$ with $\epsilon := 10^{-14}$. It bears repeating that, limiting on the density also implies preservation of the maximum compressibility. That is, the high order solution satisfies $1 - b\rho_i^{n+1} > 0$.

6.3.5 Limiting on the Internal Energy

In order to secure that our high order method is invariant domain preserving, we would like to enforce limiting on the internal energy. However, this is not necessary to do as limiting on the surrogate entropy will enforce this condition, see Section 6.4. Still, we briefly recount the method for limiting the internal energy. That is, we require $\mathbf{e}(\mathbf{U}_i^{n+1}) > 0$. The details provided in this section are not new and are outlined in [60, Section 4.5]. Define $\Psi_i^e(\mathbf{U}) := (\rho \mathbf{e})(\mathbf{U}) - \Upsilon_i^{\min} = \mathbf{E} - \frac{\|\mathbf{M}\|_{\ell^2}^2}{2\rho} - \Upsilon_i^{\min}$. The goal is to find the largest $t_0 \in [0, \ell_j^{i, \rho}]$ such that $\Psi^e(\mathbf{U}_i^{\mathbf{L}, n+1} + t_0 \mathbf{P}_{ij}) \geq 0$ for all $j \in \mathcal{G}(i)$. This can equivalently be solved using with the functional, $\psi_i^e(\mathbf{U}) : \{\mathbf{U} \in \mathbb{R}^{d+2} \mid \rho > 0\}$ where,

$$\psi_i^e(\mathbf{U}) := \rho \Psi_i^e(\mathbf{U}) = \rho \mathbf{E} - \frac{1}{2} \|\mathbf{M}\|_{\ell^2}^2 - \rho \Upsilon_i^{\min}. \quad (6.27)$$

The advantage of working with this functional is that it is quadratic with respect to the conserved variables.

The derivatives of this functional are,

$$D\psi_i^e(\mathbf{U}) = \begin{pmatrix} \mathbf{E} - \Upsilon_i^{\min} \\ -\mathbf{M} \\ \rho \end{pmatrix} \quad \mathbb{H}\psi_i^e(\mathbf{U}) = \begin{pmatrix} 0 & \mathbf{0}^\top & 1 \\ 0 & -\mathbb{I}_d & 0 \\ 1 & \mathbf{0}^\top & 0 \end{pmatrix}. \quad (6.28)$$

Note that $\psi_i^e(\mathbf{U}_i^{\mathbf{L}, n+1} + t\mathbf{P}_{ij})$ is a quadratic function of t ; in particular, set $\psi_i^e(\mathbf{U}_i^{\mathbf{L}, n+1} + t\mathbf{P}_{ij}) = at^2 + bt + c$ where $a = \frac{1}{2} \mathbf{P}_{ij}^\top \mathbb{H}\psi_i^e(\mathbf{U}_i^{\mathbf{L}, n+1}) \mathbf{P}_{ij}$, $b = D\psi_i^e(\mathbf{U}_i^{\mathbf{L}, n+1}) \cdot \mathbf{P}_{ij}$, and $c = \psi_i^e(\mathbf{U}_i^{\mathbf{L}, n+1})$. Let t_0 be the smallest positive root of $\psi_i^e(\mathbf{U}_i^{\mathbf{L}, n+1} + t\mathbf{P}_{ij}) = 0$ and set $t_0 = 1$ if no such roots exist.

We define,

$$\ell_j^{i, e} = \min(t_0, \ell_j^{i, \rho}). \quad (6.29)$$

Lemma 6.3.5. *If $\ell \in [0, \ell_j^{i,e}]$, then $\Psi_i^e(\mathbf{U}_i^{\text{L},n+1} + \ell \mathbf{P}_{ij}) \geq 0$.*

Proof. Since $\psi_i^e(\mathbf{U}_i^{\text{L},n+1} + \ell \mathbf{P}_{ij})$ is a quadratic function with $\psi_i^e(\mathbf{U}_i^{\text{L},n+1}) > 0$. By definition of t_0 , $\psi_i^e(\mathbf{U}_i^{\text{L},n+1} + \ell \mathbf{P}_{ij}) \geq 0$ for all $\ell \in [0, \ell_j^{i,e}]$, otherwise, if there existed $0 < \tilde{\ell} < \ell_j^{i,e}$, such that $\psi_i^e(\mathbf{U}_i^{\text{L},n+1} + \tilde{\ell} \mathbf{P}_{ij}) < 0$, then t_0 would not be the smallest root. Lastly, $\rho_i^{\text{L},n+1} + \ell P_{ij}^\rho > 0$ for all $\ell \in [0, \ell_j^{i,e}] \subset [0, \ell_j^{i,\rho}]$. Hence, we must have that $\Psi_i^e(\mathbf{U}_i^{\text{L},n+1} + \ell \mathbf{P}_{ij}) \geq 0$ for all $\ell \in [0, \ell_j^{i,e}]$. \square

6.4 The Entropy Surrogate

It is well known that the Euler equations with the ideal EOS preserve the minimum principle on the physical entropy, $s(\mathbf{x}, t)$, see [74].

Theorem 6.4.1 (Minimum Principle on the Specific Entropy). *Let \mathbf{u} be an entropy solution to the Euler equations supplied with the ideal EOS, then the specific entropy satisfies*

$$\operatorname{ess\,inf}_{\|\mathbf{x}\|_{\ell^2} \leq R} s(\mathbf{x}, t) \geq \operatorname{ess\,inf}_{\|\mathbf{x}\|_{\ell^2} \leq R + t\mathbf{v}_{\max}} s(\mathbf{x}, 0) \quad (6.30)$$

where $R > 0$, and $\mathbf{v}_{\max} := \max_{(\mathbf{x}, t) \in D \times [0, t]} \|\mathbf{v}(\mathbf{x}, t)\|_{\ell^2}$

However, as mentioned in Chapter 2, we may not have access to the physical entropy for an arbitrary or tabulated EOS. The novel idea of this section which was originally proposed in Clayton et. al. [1], is the introduction of a local surrogate entropy functional, which maintains some similar properties to a physical entropy. In particular, the surrogate entropy functional increases across shocks, see Theorem 6.4.1.

In this section we outline some known facts regarding the physical entropy and approximate these principles with a surrogate physical entropy for our tabulated EOS. This is all done in the context of the Riemann problem as that is where our application lies.

6.4.1 The Entropy for the NASG EOS

Recall, the NASG EOS is given by $p(\rho, e) = (\gamma - 1) \frac{\rho(e - q)}{1 - b\rho} - \gamma p_\infty$. The specific entropy is, $s(\tau, e) = \log \left((e - q - p_\infty(\tau - b))^{\frac{1}{\gamma-1}} (\tau - b) \right)$. Sometimes, we may abuse the notation and write $s = s(\mathbf{u})$. We will also need to work with the variables, ρ and e . So we may express

the entropy as $s(\rho^{-1}, e)$. Note that one can express the physical entropy in an equivalent form, that is,

$$\exp((\gamma - 1)s(\tau, e)) = (e - q - p_\infty(\tau - b))(\tau - b)^{\gamma-1}. \quad (6.31)$$

We now define $S(\mathbf{u}; \gamma) = \exp((\gamma - 1)s(\rho^{-1}, e))$; that is,

$$S(\mathbf{u}; \gamma) := (\rho(e - q) - p_\infty(1 - b\rho)) \frac{(1 - b\rho)^{\gamma-1}}{\rho^\gamma}. \quad (6.32)$$

In order to preserve some notion of a minimum principle on the specific entropy, we will introduce a local functional Ψ_i^s which behaves similar to a physical entropy.

Theorem 6.4.2 (Equivalence of the Minimum Principle on the Specific Entropy). *If the oracle for the Riemann problem 1.8a–1.8b is given by the NASG EOS, $p(\mathbf{u}) := (\gamma - 1) \frac{\rho(\mathbf{e}(\mathbf{u}) - q)}{1 - b\rho} - \gamma p_\infty$. Then, the minimum principle on the specific entropy in the Riemann problem is equivalent to*

$$\Psi^s(\mathbf{u}(x, t)) \geq \min(\Psi^s(\mathbf{u}_L), \Psi^s(\mathbf{u}_R)) \quad (6.33)$$

for all $(x, t) \in \mathbb{R} \times [0, \infty)$.

Proof. First assume that the Riemann solution contains no vacuum states. From the definition of the specific entropy, we have the expression,

$$S(\mathbf{u}) = \exp((\gamma - 1)s(\mathbf{u})) = (\rho(\mathbf{e}(\mathbf{u}) - q) - p_\infty(1 - b\rho)) \frac{(1 - b\rho)^{\gamma-1}}{\rho^\gamma} \quad (6.34)$$

Define, $S^{\min} := \min\{\exp((\gamma - 1)s(\mathbf{u}_L)), \exp((\gamma - 1)s(\mathbf{u}_R))\}$. Let \mathbf{u} be on a Z-wave, $Z \in \{L, R\}$.

If the Z-wave is an expansion, then the specific entropy is constant which implies that $S(\mathbf{u}_Z) = \exp((\gamma - 1)s(\mathbf{u}_Z)) = S(\mathbf{u}(x, t))$ for all (x, t) belonging to the expansion wave.

If the Z-wave is a shock, then the entropy increases, hence, $\exp((\gamma - 1)s(\mathbf{u}_Z)) \leq \exp((\gamma - 1)s(\mathbf{u}(x, t))) = S(\mathbf{u}(x, t))$ for all (x, t) on the shock wave. We conclude that $S(\mathbf{u}(x, t)) \geq S^{\min}$

for any (x, t) on a Z-wave. By continuity, this inequality also holds across the contact, hence $\Psi^s(\mathbf{u}(x, t)) \geq 0$ for all $(x, t) \in \mathbb{R} \times [0, \infty)$.

In the case of vacuum ($\rho = 0$), we have that $\Psi^s(\mathbf{u}(x, t)) = 0$ for (x, t) in the vacuum region. Note also that

$$S^{\min} = \min_{Z \in \{L, R\}} (\rho_Z(e_Z - q) - p_\infty(1 - b\rho_Z)) \frac{(1 - b\rho_Z)^{\gamma-1}}{\rho_Z^\gamma}, \quad (6.35)$$

and therefore $\min(\Psi^s(\mathbf{u}_L), \Psi^s(\mathbf{u}_R)) = 0$. Hence the result also holds for vacuum states. \square

We now prove a result regarding the behavior of these kind of functionals across shock waves.

Lemma 6.4.1 (Behavior across Shocks). *For all $i \in \mathcal{V}$, assume $\mathbf{U}_i^n \in \mathcal{B}(b)$. For all $i \in \mathcal{V}$ and $j \in \mathcal{G}(i)$ let $\gamma_{ij} \in [1, \min(\gamma_i^n, \gamma_j^n)]$. Given left and right data $(\rho_i^n, \mathbf{m}_i^n, \mathcal{E}_i^n, \Gamma_i^n)^\top$ and $(\rho_j^n, \mathbf{m}_j^n, \mathcal{E}_j^n, \Gamma_j^n)^\top$, respectively, the functional,*

$$\Psi_{ij}^s(\mathbf{u}) := \rho(e(\mathbf{u}) - q) - p_\infty(1 - b\rho) - S_{ij}^{\min} \rho^{\gamma_{ij}} (1 - b\rho)^{1-\gamma_{ij}}, \quad (6.36)$$

with $S_{ij}^{\min} := \min(S(\mathbf{U}_i^n; \gamma_{ij}), S(\mathbf{U}_j^n; \gamma_{ij}))$, increases across shock in the solution to the extended Riemann problem (if a shock wave exists).

Proof. To simplify the notation, we omit the superscript n . The solution to the extended Riemann problem, $(\rho, m, \mathcal{E}, \Gamma)^\top(x, t)$ is described in Chapter 4. Recall that $\gamma(x, t) = \gamma_i$ for $x/t < v^*$ and $\gamma(x, t) = \gamma_j$ for $x/t > v^*$ where v^* is the speed of the contact. Let $Z \in \{i, j\}$ and assume that the Z -wave is a shock wave. Let $\mathbf{u}_Z \in \mathcal{B}(b)$ be the state before the shock and $\mathbf{u} \in \mathcal{B}(b)$ be an arbitrary state connected to \mathbf{u}_Z through a shock wave. Since we are concerned with the Z -wave, $\gamma = \gamma_Z$ is constant, and therefore the interpolatory pressure p_{nasg} is an EOS. We abuse the notation by writing, $p_{\text{nasg}}(\tau, e) = p_{\text{nasg}}(\tau^{-1}, e, \gamma_Z)$ to simplify the notation. From the Rankine-Hugoniot equations, we have the following relationship involving only the thermodynamic quantities, see [31, Chapter III, Lemma 2.2].

$$\mathbf{e}(\mathbf{u}) - e_Z + \frac{1}{2}(p_{\text{nasg}}(\tau, \mathbf{e}(\mathbf{u})) + p_{\text{nasg}}(\tau_Z, e_Z))(\tau - \tau_Z) = 0 \quad (6.37)$$

Substituting in the definition of p_{nasg} ; that is, $p_{\text{nasg}} = (\gamma_Z - 1) \frac{e(\mathbf{u}) - q}{\tau - b} - \gamma p_\infty$, we find,

$$\mathbf{e}(\mathbf{u}) - q = (e_Z - q) \frac{1 - \frac{(\gamma_Z - 1)(\tau - \tau_Z)}{2(\tau_Z - b)}}{1 + \frac{(\gamma_Z - 1)(\tau - \tau_Z)}{2(\tau - b)}} =: r(\tau). \quad (6.38)$$

As we can see, $\mathbf{e}(\mathbf{u}) - q$ is only a function of τ along the shock curve. For $\mathbf{e}(\mathbf{u})$ to be well defined, we must have that $\tau \in (\tau_Z^\infty, \infty)$ where

$$\tau_Z^\infty := \frac{(\gamma_Z - 1)\tau_Z + 2b}{\gamma_Z + 1}. \quad (6.39)$$

Furthermore, $\mathbf{e}(\mathbf{u}) > 0$, for $\tau \in (\tau_Z^\infty, \tau_Z^0)$ where $\tau_Z^0 := \frac{(\gamma_Z + 1)\tau_Z - 2b}{\gamma_Z - 1}$. Since $\mathbf{U}_Z \in \mathcal{B}(b)$ we have that $b < \tau_Z^\infty < \tau < \tau_Z^0$.

The goal now is to show that

$$\mathcal{B}(b) \ni \mathbf{u} \mapsto \rho(e(\mathbf{u}) - q) - p_\infty(1 - b\rho) - c\rho^\gamma(1 - b\rho)^{1-\gamma} \quad (6.40)$$

is an increasing function across the shock wave for all $\gamma \in (1, \gamma_Z]$ and $c \in (0, ((e_Z - q) - p_\infty(\tau_Z - b))(\tau_Z - b)^{\gamma-1}]$. This will complete our proof since, setting $\gamma = \gamma_{ij} \leq \gamma_Z$ and $c := S_{ij}^{\min}$ we see that

$$\begin{aligned} 0 \leq S_{ij}^{\min} &\leq (\rho_Z(e_Z - q) - p_\infty(1 - b\rho_Z)) \frac{(1 - b\rho_Z)^{\gamma_{ij}-1}}{\rho_Z^{\gamma_{ij}}} \\ &= ((e_Z - q) - p_\infty(\tau_Z - b))(\tau_Z - b)^{\gamma-1}. \end{aligned} \quad (6.41)$$

Now define, $q(\tau) := (\gamma_Z - 1) \frac{r(\tau)}{\tau - b} - \gamma_Z p_\infty$ which is the pressure along the shock curve. Then $q'(\tau) = \frac{-4(e_Z - q)\gamma_Z(\gamma_Z - 1)}{(\tau + \tau_Z - 2b + \gamma_Z(\tau - \tau_Z))^2} < 0$. So q is a decreasing function of τ ; i.e., the pressure is a strictly monotonic function of ρ along the shock curve. That is, for $\rho \in [\rho_Z, b^{-1})$. In particular, the pressure is finite only in the range, $\tau \in (\tau_Z^\infty, \tau_Z]$.

Next note from equations (6.38) and (6.40), it is equivalent to prove that the following

function is a nonnegative decreasing function on shock curves,

$$(\tau_Z^\infty, \tau_Z] \ni \tau \mapsto g(\tau) := \tau^{-1}r(\tau) - p_\infty(1 - b\tau^{-1}) - c\tau^{-1}(\tau - b)^{1-\gamma}. \quad (6.42)$$

We will use the fact that the physical entropy, $s(\tau, e) = \log \left((e - q - p_\infty(\tau - b))^{\frac{1}{\gamma_Z - 1}} (\tau - b) \right)$, increases across the shock curves. That is, $s(\tau, r(\tau))$ is a decreasing function for $\tau \in (\tau_Z^\infty, \tau_Z]$. Furthermore, this also implies that $(\tau_Z^\infty, Z] \ni \tau \mapsto \varsigma(\tau) := \exp((\gamma_Z - 1)s(\tau, r(\tau))) = (r(\tau) - p_\infty(\tau - b))(\tau - b)^{\gamma_Z - 1}$ is a decreasing function. We will use this fact to prove that $g(\tau)$ is a decreasing function. So notice that

$$g(\tau) = \varsigma(\tau)\tau^{-1}(\tau - b)^{1-\gamma_Z} - c\tau^{-1}(\tau - b)^{1-\gamma}. \quad (6.43)$$

We further simplify the calculus by defining $\tilde{\zeta}(\tau) = \varsigma(\tau)(\tau - b)^{\gamma - \gamma_Z}$. Then,

$$g(\tau) = \frac{(\tau - b)^{1-\gamma}}{\tau}(\tilde{\zeta}(\tau) - c). \quad (6.44)$$

Computing the derivative we find,

$$g'(\tau) = \frac{(\tau - b)^{1-\gamma}}{\tau} \tilde{\zeta}'(\tau) - \frac{(\gamma - 1)\tau + (\tau - b)}{\tau^2(\tau - b)^\gamma} (\tilde{\zeta}(\tau) - c), \quad (6.45)$$

where

$$\tilde{\zeta}'(\tau) = (\tau - b)^{\gamma - \gamma_Z} \varsigma'(\tau) + (\gamma - \gamma_Z)(\tau - b)^{\gamma - \gamma_Z - 1} \varsigma(\tau). \quad (6.46)$$

Note that $\tilde{\zeta}'(\tau) \leq 0$, since $\varsigma'(t) \leq 0$ and $\varsigma(\tau) \geq 0$ for all $\tau \in (\tau_Z^\infty, \tau_Z]$ and $\gamma \leq \gamma_Z$. Also $\inf_{\tau \in (b, \tau_Z]} \tilde{\zeta}(\tau) = \tilde{\zeta}(\tau_Z) = e_Z(\tau_Z - b)^{\gamma - 1} \geq c$ since $\tilde{\zeta}$ is a decreasing function. Thus for $\gamma = \gamma_{ij}$ and $c = S_{ij}^{\min}$ we have that $g(\tau) \geq 0$ and $g'(\tau) \leq 0$. Therefore the mapping $[\rho_Z, \frac{1}{\tau_Z^\infty}) \ni \rho \mapsto \rho \mathbf{e}(\mathbf{u}) - p_\infty(1 - b\rho) - c\rho^\gamma(1 - b\rho)^{1-\gamma}$ is a nonnegative increasing function along the shock curves. \square

We now present one of the main theorems of this chapter.

Theorem 6.4.3 (Surrogate Entropy). *Let*

$$\gamma_i^{\min,n} := \min_{j \in \mathcal{G}(i)} \gamma_j^n, \quad S_i^{\min,n} := \min \left(\min_{j \in \mathcal{G}(i)} S(\mathbf{U}_j^n; \gamma_i^{\min,n}), \min_{j \in \mathcal{G}(i)} S(\bar{\mathbf{U}}_{ij}^n; \gamma_i^{\min,n}) \right), \quad (6.47a)$$

$$\Psi_i^s(\mathbf{u}) := \rho(\mathbf{e}(\mathbf{u}) - q) - p_\infty(1 - b\rho) - S_i^{\min,n} \rho^{\gamma_i^{\min,n}} (1 - b\rho)^{1-\gamma_i^{\min,n}}. \quad (6.47b)$$

Then the following hold for all $i \in \mathcal{V}$:

1. $\Psi_i^s : \mathcal{B}(b) \rightarrow \mathbb{R}$ is a concave functional.
2. For $\mathbf{U}_i^{\mathbf{L},n+1}$ defined in (4.1), we have that $\Psi_i^s(\mathbf{U}_i^{\mathbf{L},n+1}) \geq 0$ under the CFL constraint, $\tau \sum_{j \in \mathcal{G}(i) \setminus \{i\}} \frac{2d_{ij}^{\mathbf{L},n}}{m_i} \leq 1$.
3. Consider the extended Riemann problem 4.20 with left data, $(\rho_i^n, \mathbf{m}_i^n \cdot \mathbf{n}_{ij}, \mathcal{E}_i^n, \Gamma_i^n)$ and right data, $(\rho_j^n, \mathbf{m}_j^n \cdot \mathbf{n}_{ij}, \mathcal{E}_j^n, \Gamma_j^n)$. If the solution to the extended Riemann problem has a shock wave, then $\Psi_i^s(\mathbf{u})$ increases across the shocks.

Proof. 1. The function, $\mathbf{u} \mapsto \rho(\mathbf{e}(\mathbf{u}) - q)$ is a concave functional and $p_\infty(1 - b\rho)$ is a linear function of \mathbf{u} . Then notice that $f(x) := x^{\gamma_i^{\min,n}} (1 - bx)^{1-\gamma_i^{\min,n}} = x(\frac{1}{x} - b)^{\gamma_i^{\min,n}}$. A quick exercise in calculus shows that $f(x)$ is a convex for $\gamma_i^{\min,n}$, hence $-S_i^{\min,n} \rho^{\gamma_i^{\min,n}} (1 - b\rho)^{1-\gamma_i^{\min,n}}$ is a concave functional since $S_i^{\min,n} \geq 0$. Since the sum of concave functionals is concave, we have that $\Psi_i^s(\mathbf{u})$ is concave.

2. From Theorem 4.1.1, under the CFL condition, $\Delta t \leq -\frac{m_i}{2d_{ij}^{\mathbf{L},n}}$ we have that $\mathbf{U}_i^{\mathbf{L},n+1}$ is in the convex hull of the bar states $\{\bar{\mathbf{U}}_{ij}^n\}_{j \in \mathcal{G}(i)}$ and since Ψ_i^s is a concave functional we have that $\Psi_i^s(\mathbf{U}_i^{\mathbf{L},n+1}) \geq \min_{j \in \mathcal{G}(i)} \Psi_i^s(\bar{\mathbf{U}}_{ij}^n)$. To see this is nonnegative, consider,

$$\begin{aligned} \Psi_i^s(\bar{\mathbf{U}}_{ij}^n) &= \rho(\mathbf{e}(\bar{\mathbf{U}}_{ij}^n) - q) - p_\infty(1 - b\bar{\rho}_{ij}^n) - S_i^{\min,n} (\bar{\rho}_{ij}^n)^{\gamma_i^{\min,n}} (1 - b\bar{\rho}_{ij}^n)^{1-\gamma_i^{\min,n}} \\ &= (S(\bar{\mathbf{U}}_{ij}^n; \gamma_i^{\min,n}) - S_i^{\min,n}) (\bar{\rho}_{ij}^n)^{\gamma_i^{\min,n}} (1 - b\bar{\rho}_{ij}^n)^{1-\gamma_i^{\min,n}}, \end{aligned}$$

where $S(\mathbf{U}; \gamma)$ is defined in (6.32). By definition of $S_i^{\min,n}$ we have that $\Psi_i^s(\bar{\mathbf{U}}_{ij}^n) \geq 0$, hence $\Psi_i^s(\mathbf{U}_i^{\mathbf{L},n+1}) \geq 0$.

3. This immediately follows from, Lemma 6.4.1. \square

6.4.2 Limiting on the Surrogate Entropy

Let $\mathcal{A}(b) := \{\mathbf{u} \in \mathbb{R}^{d+2} : 0 < 1 - b\rho < 1\}$. From the limiting on the density described in Section 6.3.4, we have $\ell_j^{i,\rho} \in [0, 1]$. This implies that the $\rho_i^{\mathbf{L},n+1} + \ell \mathbf{P}_{ij} \in \mathcal{A}(b)$ for all $j \in \mathcal{I}(i) \setminus \{i\}$ and $\ell \in [0, \ell_j^{i,\rho}]$. Therefore, we can perform the quasiconcave limiting of Ψ_i^s on $\mathcal{A}(b)$ since Ψ_i^s is concave and $\mathcal{A}(b)$ is a convex set. In particular, we seek the largest $\ell_0 \in [0, \ell_{ij}^\rho]$ such that $\Psi_i^s(\mathbf{U}_i^{\mathbf{L},n+1} + \ell_0 \mathbf{P}_{ij}) \geq 0$. If $\Psi_i^s(\mathbf{U}_i^{\mathbf{L},n+1} + \ell_{ij}^\rho \mathbf{P}_{ij}) \geq 0$, then we set $\ell_0 := \ell_{ij}^\rho$. If this is not the case, then we must solve $h(\ell) := \Psi_i^s(\mathbf{U}_i^{\mathbf{L},n+1} + \ell \mathbf{P}_{ij}) = 0$. In particular, since $h(0) = \Psi_i^s(\mathbf{U}_i^{\mathbf{L},n+1}) \geq 0$, and Ψ_i^s is continuous and concave, there exists at least one solution. If Ψ_i^s is strictly concave, then the solution set is a singleton; otherwise, the solution set is connected.

Remark 6.4.1. Note that $\Psi_i^n(\mathbf{U}_i^{\mathbf{L},n+1}) = 0$ occurs when $\gamma_i^n = \gamma_i^{\min,n}$. \square

Remark 6.4.2 (Specific Internal Energy). It is not necessary to perform any limiting on the internal energy. To see this, let \mathbf{U}_i^{n+1} to be the final update after limiting on the surrogate entropy. Then $\Psi_i^s(\mathbf{U}_i^{n+1}) \geq 0$ implies the following,

$$\begin{aligned} e(\mathbf{U}_i^{n+1}) - q &\geq p_\infty(\tau_i^{n+1} - b) + S_i^{\min,n} \rho_i^{\min,n} (1 - b\rho)^{1-\gamma_i^{\min,n}} \\ &\geq p_\infty(\tau_i^{n+1} - b). \end{aligned} \tag{6.48}$$

This is the invariant domain property that we expect when interpolating with the NASG EOS. As a reminder, if b , q , and p_∞ are unknown then we take $b = q = p_\infty = 0$ which results in positivity of the specific internal energy. \square

6.5 The Quadratic Newton Method

For computational efficiency, the limiting is not actually performed on Ψ_i^s but rather, $\Phi_i^s(\mathbf{u}) := \rho \Psi_i^s(\mathbf{u})$ as this functional has some nicer properties. This is valid, as the solution sets for $\Psi_i^s(\mathbf{U}_i^{\mathbf{L},n+1} + \ell \mathbf{P}_{ij}) = 0$ and $\Phi_i^s(\mathbf{U}_i^{\mathbf{L},n+1} + \ell \mathbf{P}_{ij}) = 0$ are identical since $\rho_i^{\mathbf{L},n+1} + \ell \mathbf{P}_{ij}^\rho > 0$ for all $\ell \in [0, \ell_{ij}^\rho]$.

Lemma 6.5.1 (Sign of $f'''(\ell)$). Let $\mathbf{u}_0 \in \mathcal{A}(b) := \{\mathbf{u} \in \mathbb{R}^{d+2} : 0 < 1 - b\rho < 1\}$ and $\mathbf{p} = (p_1, \dots, p_{d+2})^\top \in \mathbb{R}^{d+2}$. Let $\ell_0 \in [0, 1]$ such that $\mathbf{u}_0 + \ell_0 \mathbf{p} \in \mathcal{A}(b)$. Define $f(\ell) : [0, \ell_0] \ni \ell \mapsto f(\ell) := \Phi_i^s(\mathbf{u}_0 + \ell \mathbf{p})$. Then the sign of $f'''(\ell)$ is constant over $[0, \ell_0]$.

Proof. Using that $E = \rho e + \frac{\|\mathbf{m}\|_{\ell^2}^2}{2\rho}$, we write,

$$\Phi_i^s(\mathbf{u}) = \rho \Psi_i^s(\mathbf{u}) = \rho E - \frac{1}{2} \|\mathbf{m}\|_{\ell^2}^2 - p_\infty \rho (1 - b\rho) - S_{\min} \rho^{\gamma+1} (1 - b\rho)^{1-\gamma}. \quad (6.49)$$

Note that $f(\ell)$ is well defined for all $\ell \in [0, \ell_0]$ since $\mathbf{u}_0 + \ell_0 \mathbf{p} \in \mathcal{A}(b)$ and $\mathcal{A}(b)$ is a convex set. Hence $\mathbf{u}_0 + \ell \mathbf{p} \in \mathcal{A}(b)$ for all $\ell \in [0, \ell_0]$.

Next notice that $\rho E - \frac{1}{2} \|\mathbf{m}\|_{\ell^2}^2 - p_\infty \rho (1 - b\rho)$ is only quadratic in ρ , hence only the last term $S_i^{\min, n} \rho^{\gamma+1} (1 - b\rho)^{1-\gamma}$ remains when computing the third derivative. In particular, we have that,

$$f'(\ell) = D\Phi_i^s(\mathbf{u} + \ell \mathbf{p}) \mathbf{p} \quad \text{and} \quad f''(\ell) = \mathbf{p}^\top \mathbb{H}\Phi_i^s(\mathbf{u} + \ell \mathbf{p}) \mathbf{p}. \quad (6.50)$$

In particular, note that the Hessian of Φ_i^s is,

$$\mathbb{H}\Phi_i^s = \begin{bmatrix} 2p_\infty b - S_i^{\min, n} \frac{\rho^{\gamma-1}}{(1-b\rho)^{\gamma+1}} (\gamma(\gamma+1) - 2b\rho(\gamma+1-b\rho)) & \mathbf{0}^\top & 1 \\ \mathbf{0} & -\mathbb{I}_d & \mathbf{0} \\ 1 & \mathbf{0}^\top & 0 \end{bmatrix}. \quad (6.51)$$

Therefore, when computing the $f'''(\ell)$ the only term that remains when applying the chain rule is $\partial_\rho^3 \Phi_i^s$. Thus,

$$f'''(\ell) = (p_1)^3 \frac{\partial \Phi_i^s}{\partial \rho}(\mathbf{u} + \ell \mathbf{p}) = -(p_1)^3 S_i^{\min, n} \frac{\gamma(\gamma^2 - 1)(\rho + \ell p_1)^{\gamma-2}}{(1 - b(\rho + \ell p_1))^{\gamma+2}}. \quad (6.52)$$

Since $\mathbf{u} + \ell \mathbf{p} \in \mathcal{A}(b)$, for all $\ell \in [0, \ell_0]$ we see that $\text{sgn}(f'''(\ell))$ has the same sign as $-(p_1)^3$; i.e., $\text{sgn}(f'''(\ell))$ is constant. \square

6.5.1 Review of Divided Differences

For the sake of clarity, we briefly review the notation and some results in regards to the Newton divided differences. Let $f \in C^n([a, b])$, then for $x_0, x_1 \in [a, b]$ with $x_0 \neq x_1$, we define the following,

$$f[x_0] := f(x_0), \quad f[x_0, x_1] := \frac{f[x_1] - f[x_0]}{x_1 - x_0}, \quad f[x_0, x_0] := f'(x_0). \quad (6.53)$$

So in general for $\{x_i\}_{i=1}^n \subset [a, b]$ we have,

$$f[x_0, \dots, x_n] = \begin{cases} \frac{1}{n!} f^{(n)}(x_0), & \text{if } x_0 = x_1 = \dots = x_n, \\ \frac{f[x_1, \dots, x_n] - f[x_0, \dots, x_{n-1}]}{x_n - x_0}, & \text{otherwise.} \end{cases} \quad (6.54)$$

In addition, notice that $f[x_0, \dots, x_n]$ is unchanged by permutation. That is, given a permutation, $\sigma \in S(n)$, we have that $f[x_{\sigma(0)}, \dots, x_{\sigma(n)}] = f[x_0, \dots, x_n]$.

6.5.2 The Quadratic Newton Method

The quadratic Newton method was first introduced by Guermond & Popov in [59] and has been adapted for finding the root of $\Phi_i^s(\mathbf{U}_i^{\mathbf{L}, n+1} + \ell \mathbf{P}_{ij}) = 0$. We Using the divided difference notation, we define the “left” and “right” interpolating polynomials,

$$P_L(\ell) := f(\ell_L) + f'(\ell_L)(\ell - \ell_L) + f[\ell_L, \ell_L, \ell_R](\ell - \ell_L)^2 \quad (6.55a)$$

$$P_R(\ell) := f(\ell_R) + f'(\ell_R)(\ell - \ell_R) + f[\ell_L, \ell_R, \ell_R](\ell - \ell_R)^2 \quad (6.55b)$$

For P_L , we see that $P_L(\ell_L) = f(\ell_L)$, $P_L(\ell_R) = f(\ell_R)$, and $P'_L(\ell_L) = f'(\ell_L)$. Similarly for P_R we have $P_R(\ell_L) = f(\ell_L)$, $P_R(\ell_R) = f(\ell_R)$, and $P'_R(\ell_R) = f'(\ell_R)$.

Let ℓ^* denote the root $\Phi_i^s(\mathbf{U}_i^{\mathbf{L}, n+1} + \ell \mathbf{P}_{ij}) = 0$. For the quadratic Newton method, we will show that we can always find $\tilde{\ell} \leq \ell^*$ for which $\tilde{\ell} \uparrow \ell^*$. Hence, we guarantee that $\Phi_i^s(\mathbf{U}_i^{\mathbf{L}, n+1} + \tilde{\ell} \mathbf{P}_{ij}) > 0$ at each step in the quadratic Newton method. Then we have the

following, Lemma 6.5.2, in regards to the approximation of $f(\ell)$ with $P_L(\ell)$ and $P_R(\ell)$. The pseudo-code for the quadratic Newton method is given in Algorithm 2 in Appendix B.

Lemma 6.5.2 (Interpolation Properties of P_L and P_R). *The following holds true:*

1. *The polynomials $P_L(\ell)$ and $P_R(\ell)$ bound the function $f(\ell)$ in the following sense:*

$$\min(P_L(\ell), P_R(\ell)) < f(\ell) < \max(P_L(\ell), P_R(\ell)), \quad \forall \ell \in (\ell_L, \ell_R). \quad (6.56)$$

2. *$P_L(\ell)$ and $P_R(\ell)$ each have a unique zero over the interval (ℓ_L, ℓ_R) respectively given by*

$$\ell^L(\ell_L, \ell_R) := \begin{cases} \ell_L - \frac{2f(\ell_L)}{f'(\ell_L) + \sqrt{f'(\ell_L)^2 - 4f(\ell_L)f[\ell_L, \ell_L, \ell_R]}}, & \text{if } f[\ell_L, \ell_L, \ell_R] < 0, \\ \ell_L - \frac{2f(\ell_L)}{f'(\ell_L) - \sqrt{f'(\ell_L)^2 - 4f(\ell_L)f[\ell_L, \ell_L, \ell_R]}}, & \text{if } f[\ell_L, \ell_L, \ell_R] \geq 0, \end{cases} \quad (6.57a)$$

$$\ell^R(\ell_L, \ell_R) := \begin{cases} \ell_R - \frac{2f(\ell_R)}{f'(\ell_R) + \sqrt{f'(\ell_R)^2 - 4f(\ell_R)f[\ell_L, \ell_R, \ell_R]}}, & \text{if } f[\ell_L, \ell_R, \ell_R] < 0, \\ \ell_R - \frac{2f(\ell_R)}{f'(\ell_R) - \sqrt{f'(\ell_R)^2 - 4f(\ell_R)f[\ell_L, \ell_R, \ell_R]}}, & \text{if } f[\ell_L, \ell_R, \ell_R] \geq 0. \end{cases} \quad (6.57b)$$

3. *Properties 1. and 2. imply that*

$$\min(\ell^L(\ell_L, \ell_R), \ell^R(\ell_L, \ell_R)) < \ell^* < \max(\ell^L(\ell_L, \ell_R), \ell^R(\ell_L, \ell_R)). \quad (6.58)$$

Proof. Without loss of generality, let $P_L \leq P_R$ and $f'''(\ell) > 0$. Then, with a hefty amount of rewriting, one can show that,

$$f(\ell) - P_L(\ell) = f[\ell_R, \ell_L, \ell_L, \ell](\ell_R - \ell)(\ell - \ell_L)^2 \quad (6.59)$$

1. Notice that, $f[\ell_R, \ell_L, \ell_L, \ell_R] = f'''(\xi) > 0$ for some $\xi \in (\ell_L, \ell_R)$. Then we want to show

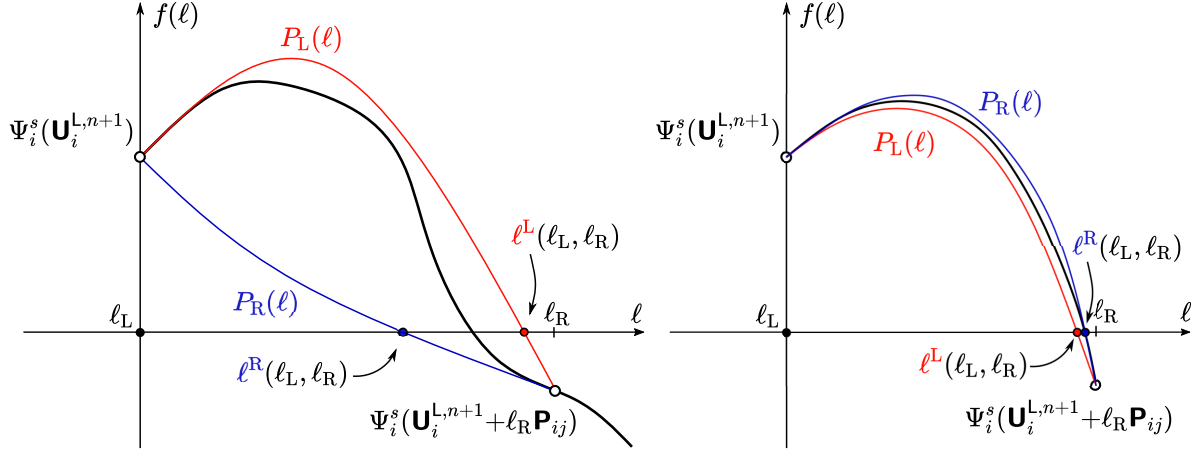


Figure 6.1: Two example visual descriptions of limiting using the quadratic newton method. Left: quasiconcave function that is not concave. Right: concave function.

that $f[\ell_R, \ell_L, \ell_L, \ell] > 0$ for $\ell \in [\ell_L, \ell_R)$. Also note that,

$$f[\ell_R, \ell_L, \ell_L, \ell] = \frac{f[\ell_L, \ell_L, \ell] - f[\ell_R, \ell_L, \ell_L]}{\ell - \ell_R} = \frac{f[\ell_L, \ell_L, \ell_R] - f[\ell, \ell_L, \ell_L]}{\ell_R - \ell}. \quad (6.60)$$

Since $f'''(\ell) > 0$, then $f''(\ell)$ is an increasing function, and therefore, $f[\ell_L, \ell_L, \ell_R] \geq f[\ell, \ell_L, \ell_L]$ for all $\ell \in [\ell_L, \ell_R)$. Therefore, $f[\ell_R, \ell_L, \ell_L, \ell] \geq 0$ for all $\ell \in [\ell_L, \ell_R)$. Dropping this right hand side this implies that, $f(\ell) > P_L(\ell)$. Showing that $P_R(\ell) > f(\ell)$ is similar, and if $f'''(\ell) < 0$ then the P_R is the lower bound and P_L is the upper bound.

2. Note the formula is simply the quadratic formula but rationalizing the numerator to avoid division by zero in the degenerate case. Since both quadratic functions interpolate a positive and a negative values, the function must have a single root. The root which lies in the interval (ℓ_L, ℓ_R) depends on the concavity or convexity of the quadratic function. Hence the result.

3. This follows immediately from 1. and 2. □

6.6 Relaxation on the Surrogate Entropy

We also perform relaxation on the surrogate entropy exactly as described in [60, Sec. 4.7.2]. Specifically, we relax the bound, $S_i^{\min,n}$. Specifically, we define the new bound,

$$\widetilde{S}_i^{\min,n} := \max((1 - r_h)S_i^{\min,n}, S_i^{\min,n} - \Delta S_i^n) \quad (6.61)$$

where $\Delta S_i^n := \max_{j \in \mathcal{G}(i) \setminus \{i\}} (\frac{1}{2}(S_i^n + S_j^n) - S_i^{\min,n})$ and $r_h := (m_i/|D|)^{1.5/d}$. This type of relaxation was originally done by Khobalatte & Perthame in [75, Sec. 3.3] in order to guarantee second order convergence of their numerical method. Since the entropy is constant in regions of smooth flow, we loosen the bound on the order of $\mathcal{O}(h^2)$. However, for the surrogate entropy, we cannot say that $S(\mathbf{u}; \gamma_i^{\min,n})$ will be constant in the region of smooth flow and in fact it may even decrease. If there is a shock wave then the relaxation is irrelevant as the first order scheme is used based on the entropy viscosity method outlined in Chapter 5.

7. NUMERICAL RESULTS*

In this chapter we present a variety of numerical illustrations to demonstrate the robust effectiveness of the proposed methods. This chapter is broken up into two parts one-dimensional and two-dimensional problems. When relevant, first order and second order solutions will be compared.

The numerical results for the one-dimensional problems have been performed on an in-house code written in Fortran originally developed by Dr. Jean-Luc Guermond and uses \mathbb{P}_1 finite elements. This code base will be referred to as the TAMU code. The two-dimensional problems have been performed on a software named `Ryujin` which can be found at <https://github.com/conservation-laws/ryujin>. `Ryujin` is an efficient multithreaded massively parallel C++ code which uses the `deal.II` finite element library, [76]. The design of the efficient and parallel algorithms for `Ryujin` is done by Maier & Kronbichler in [65]. The two-dimensional simulations were performed on the Texas A&M Mathematics department's cluster, nicknamed Whistler, which used anywhere between 10 and 32 compute nodes.

The time step for the TAMU code was computed by the following,

$$\Delta t := \text{CFL} \min_{i \in \mathcal{V}} \left(-\frac{m_i}{2d_{ii}^{L,n}} \right). \quad (7.1)$$

In the 2D simulations, the boundary conditions we enforce are Dirichlet, slip, and outflow. For the slip condition, we want to enforce $\mathbf{v} \cdot \mathbf{n} = 0$. This is achieved by redefining the momentum on the slip boundary with,

$$\mathbf{M}_i^{n+1} := \mathbf{M}_i^{n+1} - (\mathbf{M}_i^{n+1} \cdot \mathbf{n}_i) \mathbf{n}_i. \quad (7.2)$$

where $i \in \mathcal{V}$ is the index of a node on the boundary, ∂D , and \mathbf{n}_i is the outward normal

* A majority of the simulations presented in this section are also reported in [2] and [1] and is reprinted with permission from [2] and [1].

at node \mathbf{x}_i . Since the domain is polygonal, we define $\mathbf{n}_i = \frac{1}{2}(\mathbf{n}_i^- + \mathbf{n}_i^+)$ where \mathbf{n}_i^- and \mathbf{n}_i^+ are the outward normals on the boundary edges which connect to \mathbf{x}_i . For the outflow condition, we use an approximate solution to the Riemann problem, the details of which can be found in Appendix B.3. This is by no means a perfect outflow implementation, but in certain circumstances, it can work properly. It should be noted that the solution can blow up on the boundary with this implementation. However, effective implementations of outflow boundary conditions are not the focus of this thesis nor on the specific numerical demonstrations.

Lastly, the Schlieren plots mentioned in the 2D section are computed with the following formula,

$$\exp\left(-\beta \frac{r_i^n - \min_{j \in \mathcal{V}} r_j^n}{\max_{j \in \mathcal{V}} r_j^n - \min_{j \in \mathcal{V}} r_j^n}\right), \quad \text{where} \quad r_i^n := \frac{1}{m_i} \left\| \sum_{j \in \mathcal{G}(i)} \mathbf{c}_{ij} \rho_j^n \right\|_{\ell^2}. \quad (7.3)$$

We use $\beta = 15$ for all 2D simulations. Generally speaking, this formula plots the gradient of the pressure and is constructed so as to mimic the photography style known as ‘‘Schlieren photography.’’

7.1 Convergence Tests

Let $(\rho_h(t), \mathbf{m}_h(t), E_h(t))$ denote the approximate solution at time t . We define the consolidated error indicator as the sum of the relative errors for the density, momentum, and total energy; that is,

$$\delta_q(t) := \frac{\|\rho_h(t) - \rho(t)\|_{L^q(D)}}{\|\rho(t)\|_{L^q(D)}} + \frac{\|\mathbf{m}_h(t) - \mathbf{m}(t)\|_{L^q(D)}}{\|\mathbf{m}(t)\|_{L^q(D)}} + \frac{\|E_h(t) - E(t)\|_{L^q(D)}}{\|E(t)\|_{L^q(D)}}, \quad (7.4)$$

for $q \in [1, \infty]$.

7.1.1 The Fan-Jump-Fan Composite Wave

We begin by verifying that our numerical method converges. In particular, we show convergence for a Riemann problem with the van der Waals equation of state. Recall the

van der Waals equation of state is, $p(\rho, e) = (\gamma - 1) \frac{\rho e + a\rho^2}{1 - b\rho} - a\rho^2$. We use $\gamma = 1.02$, $a = 1$, and $b = 1$. The left and right states are chosen to be,

$$(\rho_L, v_L, p_L) := (0.10, -0.475504638574729, 0.022084258693080), \quad (7.5)$$

$$(\rho_R, v_R, p_R) := (0.39, -0.121375781741349, 0.039073167077590), \quad (7.6)$$

and the computational domain is $D = (-1, 1)$ and we use CFL = 0.5. The solution to this Riemann problem is a fan-jump-fan composite wave. That is, the solution is a single wave (3-wave) and is composed of an expansion, immediately followed by a shock and into another expansion. The construction of the exact solution to this problem can be found in the supplementary material of Clayton et. al. [2] at [local/web 319KB]. The exact solution to this problem has also been constructed in [77] and [78]; similar constructions can also be seen in [79].

Convergence rates are reported in Table 7.1 and plots of the approximate and exact solutions are shown in Figure 7.1 for the density and pressure.

#dof	$\delta_1(t)$	rate	$\delta_2(t)$	rate
101	2.14E-01	–	2.67E-01	–
201	1.44E-01	0.58	2.07E-01	0.37
401	9.40E-02	0.62	1.58E-01	0.39
801	5.96E-02	0.66	1.20E-01	0.40
1601	3.66E-02	0.70	8.96E-02	0.42
3201	2.18E-02	0.75	6.66E-02	0.43
6401	1.27E-02	0.78	4.93E-02	0.43
12801	7.26E-03	0.81	3.66E-02	0.43
25601	4.09E-03	0.83	2.72E-02	0.43

Table 7.1: Consolidated errors and convergence rates for the fan-jump-fan composite wave. Solution computed at $t = 5.0$. Reprinted with permission from [2].

Remark 7.1.1 (Existence of Composite Waves). For the van der Waals EOS, the composite

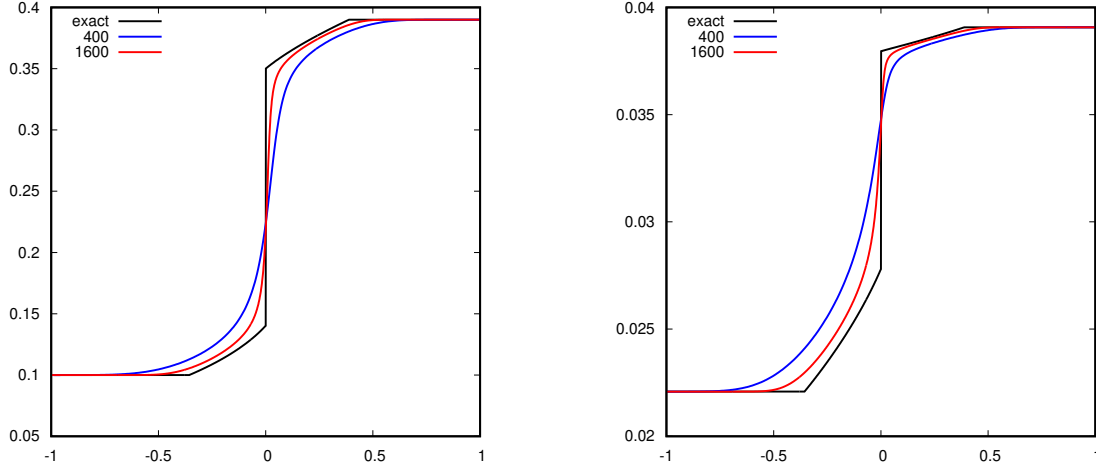


Figure 7.1: Comparison of the exact solution for density (left) and pressure (right) for the fan-jump-fan composite wave. Approximate solutions are computed using 400 and 1600 DoFs and the corresponding mesh sizes are $h = 0.005$ and $h = 0.00125$, respectively.

waves can only exist if $\gamma < 1.06$. This condition implies there is a region in the thermodynamic phase space where the isentropes are non-convex in the (p, τ) diagram. Hence composite waves can emerge. This has been shown in [80, Section 3.1], [78], and [79, Section 6.3]; more information on composite waves can also be found in [37]. \square

7.1.2 Smooth Wave with Various EOS

To further demonstrate the robustness, we show convergence rates for the first and second order methods of a smooth travelling wave. The exact solution for the smooth wave is,

$$\rho(x, t) = \begin{cases} \rho_0 + 2^6(x_1 - x_0)^{-6}(x - v_0t - x_0)^3(x_1 - x + v_0t)^3 & \text{if } x_0 \leq x - v_0t \leq x_1, \\ \rho_0 & \text{otherwise,} \end{cases} \quad (7.7a)$$

$$v(x, t) = v_0, \quad p(x, t) = p_0, \quad (7.7b)$$

where x_0 and x_1 are arbitrary constants with $x_0 < x_1$. The constants, ρ_0 , v_0 , and p_0 are chosen depending on the EOS being used. For our simulations we choose the computational domain to be $D = (0, 1)$ with $x_0 = 0.1$ and $x_1 = 0.3$. All tests were performed using

CFL = 0.1. Note that the solution to this problem is independent of the EOS since the exact solution has constant pressure.

7.1.2.1 *Ideal EOS*

This problem has been performed before in [60] but we repeat it here for comparison. We the parameters, $\gamma = 1.4$, $\rho_0 = v_0 = p_0 = 1$, and for the interpolation parameters, we set $b = q = p_\infty = 0$. The final time is run to $t = 0.6$.

7.1.2.2 *Van der Waals EOS*

We set the parameters for the van der Waals EOS to be, $a = 1$, $b = 0.075$, $\gamma = 1.4$ and $\rho_0 = v_0 = p_0 = 1$. The interpolation parameters are set to be $b = 0.075$ and $q = p_\infty = 0$. The final time is run to $T = 0.6$.

7.1.2.3 *Jones-Wilkins-Lee EOS*

We set the parameters for the Jones-Wilkins-Lee EOS to be, $A = 1$, $B = -1$, $R_1 = 2$, $R_2 = \omega = \rho_0 = v_0 = p_0 = 1$. The interpolation parameters are set to be $b = q = p_\infty = 0$ and the final time is run to $t = 0.6$.

7.1.2.4 *Mie-Gruneisen EOS*

We set the parameters for the Mie-Gruneisen EOS to be, $\tilde{\rho}_0 = 2790$, $c_0 = 5330$, $s = 1.34$, $P_0 = 0$, $e_0 = 0$, $\Gamma_0 = 2.00$, $\rho_0 = 3500$, $v_0 = 1 \times 10^4$ and $p_0 = 1 \times 10^{11}$. The interpolation parameters are set to $b = q = p_\infty = 0$ and the final run time is $t = 6 \times 10^{-5}$.

The convergence results for the second order method are reported in Table 7.2.

7.1.3 **The Isentropic Vortex with van der Waals EOS**

To demonstrate the convergence of the high and low order methods in 2D, we run a simulation for the isentropic vortex using the van der Waals equation of state. The exact

$ \mathcal{V} $	Ideal		VdW		JWL		MG	
	$\delta_\infty(t_{\text{final}})$		$\delta_\infty(t_{\text{final}})$		$\delta_\infty(t_{\text{final}})$		$\delta_\infty(t_{\text{final}})$	
101	1.94e-02	–	1.24e-01	–	7.93e-02	–	1.24e-05	–
201	4.03e-03	2.27	6.24e-03	4.30	2.53e-02	1.65	2.56e-06	2.28
401	7.91e-04	2.35	9.92e-04	2.65	3.61e-03	2.81	5.03e-07	2.35
801	1.44e-04	2.46	1.75e-04	2.51	1.31e-04	4.78	9.17e-08	2.46
1601	2.75e-05	2.39	3.29e-05	2.41	2.51e-05	2.38	1.75e-08	2.39
3201	5.18e-06	2.41	6.17e-06	2.41	4.73e-06	2.41	3.29e-09	2.41
6401	9.69e-07	2.42	1.16e-06	2.42	8.87e-07	2.42	6.22e-10	2.41

Table 7.2: $\delta_\infty(t_{\text{final}})$ error defined in equation (7.4) and corresponding convergence rates with various EOS for the one-dimensional smooth traveling wave problem with exact solution (7.7) under uniform refinement of the interval $D = (0, 1)$. Reprinted with permission from [1].

solution is given by,

$$\rho(\mathbf{x}, t) = \left[\frac{3C}{8a} - \frac{1}{2} \sqrt{\frac{9C^2}{16a^2} + \frac{2}{a} \left(F + \frac{1}{2r_0^2} \psi(\bar{\mathbf{x}})^2 \right)} \right]^2, \quad \mathbf{x} \in \mathbb{R}^2, \quad t > 0, \quad (7.8a)$$

$$\mathbf{v}(\mathbf{x}, t) = \mathbf{v}_\infty + \psi(\bar{\mathbf{x}}) (-\bar{x}_2, \bar{x}_1), \quad \mathbf{x} \in \mathbb{R}^2, \quad t > 0, \quad (7.8b)$$

$$p(\mathbf{x}, t) = C(\gamma - 1)\rho(\mathbf{x}, t)^\gamma - a\rho(\mathbf{x}, t)^2, \quad \mathbf{x} \in \mathbb{R}^2, \quad t > 0. \quad (7.8c)$$

where $\bar{\mathbf{x}} := \mathbf{x} - \mathbf{x}^0 - \mathbf{v}_\infty t = (\bar{x}_1, \bar{x}_2)$, $C := (p_\infty + a\rho_\infty^2)/\rho_\infty^{3/2}$, $F := -a\rho_\infty - 3p_\infty/\rho_\infty$. Here ρ_∞ and p_∞ are the density and pressure in the far field, and $\psi(\mathbf{x}) := \frac{\beta}{2\pi} \exp\left(\frac{1}{2}\left(1 - \frac{\|\mathbf{x}\|_{\ell_2}^2}{r_0^2}\right)\right)$. We set the far field conditions to $\rho_\infty = 0.1$, $p_\infty = 1$ and $\mathbf{v}_\infty = (1, 1)$. We also set $\gamma = \frac{3}{2}$ and $a = 1$. This gives $C = \frac{101}{\sqrt{10}}$ and $F = -\frac{301}{10}$. The rest of the constants are set as follows: $\mathbf{x}^0 = (-1, -1)$, $r_0 = 1$, $\beta = 20$. The derivation of the exact solution is shown in Appendix A.1 for the ideal and van der Waals EOS.

The numerical simulation is performed on the square domain, $D = (-5, 5) \times (-5, 5)$. We impose Dirichlet boundary conditions on all sides of D , since the solution decays rapidly to the far field state, $(\rho_\infty, \mathbf{v}_\infty, p_\infty)$, away from \mathbf{x}^0 . Convergence results are reported in Table 7.3.

$ \mathcal{V} $	$\delta_1(t_{\text{final}})$		$\delta_2(t_{\text{final}})$		$\delta_\infty(t_{\text{final}})$	
289	1.17e-01	–	2.01e-01	–	6.82e-01	–
1089	1.18e-02	3.46	2.65e-02	3.06	1.05e-01	2.82
4225	7.92e-04	3.98	1.96e-03	3.84	7.87e-03	3.82
16641	5.57e-05	3.87	1.32e-04	3.93	5.50e-04	3.88
66049	5.07e-06	3.48	1.20e-05	3.48	7.79e-05	2.83
263169	7.55e-07	2.76	2.25e-06	2.42	2.03e-05	1.95
1050625	1.64e-07	2.20	5.51e-07	2.04	5.52e-06	1.88
4198401	4.08e-08	2.01	1.38e-07	2.00	1.51e-06	1.87

Table 7.3: The consolidated error defined in equation (7.4) and convergence rates for the isentropic vortex problem with the Van der Waals EOS. The exact solution is given in (7.8). Reprinted with permission from [1].

7.2 The Two-Expansion Wave Speed Estimate

As mentioned in Chapter 4 the computation of the max wave speed for a local Riemann problem is extremely difficult if the EOS is complicated. A common heuristic approach is to use the so-called *two-expansion approximation* (or *two-rarefaction approximation*) for the wave speed. That is, an approximate wave speed is defined by,

$$\lambda^{\text{exp}} := \max\{|v_L - a_L|, |v_R + a_R|\}. \quad (7.9)$$

Note that a_L and a_R are the material sound speed for the oracle and should not be confused with the sound speed used in the NASG EOS used in the interpolation in Chapter 4. We display several test problems for which this heuristic estimation is not robust. It can lead to an overestimation or underestimation of the artificial viscosity which leads immediate problems.

The test problems in this section all use the van der Waals EOS with $\gamma = 1.02$, $a = 1$, and $b = 1$ as in Section 7.1.1.

7.2.1 Underestimation of Max Wave Speed: Test 1

In this problem we show that the two-expansion approximation leads to an underestimation of the viscosity which leads to non-physical results. The Riemann data is,

$$(\rho_L, v_L, p_L) := (0.2450, 0, 2.9123894332846005 \times 10^{-2}), \quad (7.10)$$

$$(\rho_R, v_R, p_R) := (0.1225, 0, 2.0685894810791836 \times 10^{-2}), \quad (7.11)$$

with corresponding sound speeds $(a_L, a_R) \approx (0.00399, 0.306)$. The computational domain is $D = (-0.5, 1)$ and we use $\text{CFL} = 0.5$. The simulation is run up to $t = 1.25$. The results for the low order method using Algorithm 1 for computing an upper bound on the max wave speed are seen in Figure 7.2. When we use λ^{exp} for the max wave speed, the computation generates complex sound speed after a few time steps. A comparison of these methods can be seen in Figure 7.3.

7.2.2 Underestimation of Max Wave Speed: Test 2

For the second test, we use the following Riemann data,

$$(\rho_L, v_L, p_L) := (2.5 \times 10^{-1}, 0, 3 \times 10^{-2}), \quad (7.12)$$

$$(\rho_R, v_R, p_R) := (4.9 \times 10^{-5}, 0, 5 \times 10^{-8}), \quad (7.13)$$

The computational domain is $D = (-0.5, 1)$ and we use $\text{CFL} = 0.5$. The final time is $t = 0.4$. As before, using the estimation of \widehat{p}^* from Section 4.6 we have the physically relevant solution shown in Figure 7.4. When we use λ^{exp} as the estimate, the simulation immediately crashes, generating negative specific internal energy; this crash occurs no matter the CFL number.

7.2.3 Overestimation of Max Wave Speed: Test 3

For the third test, we demonstrate that λ^{exp} can be a gross overestimate of the max wave speed resulting in a dramatic increase in computational time. This is due a decrease in the

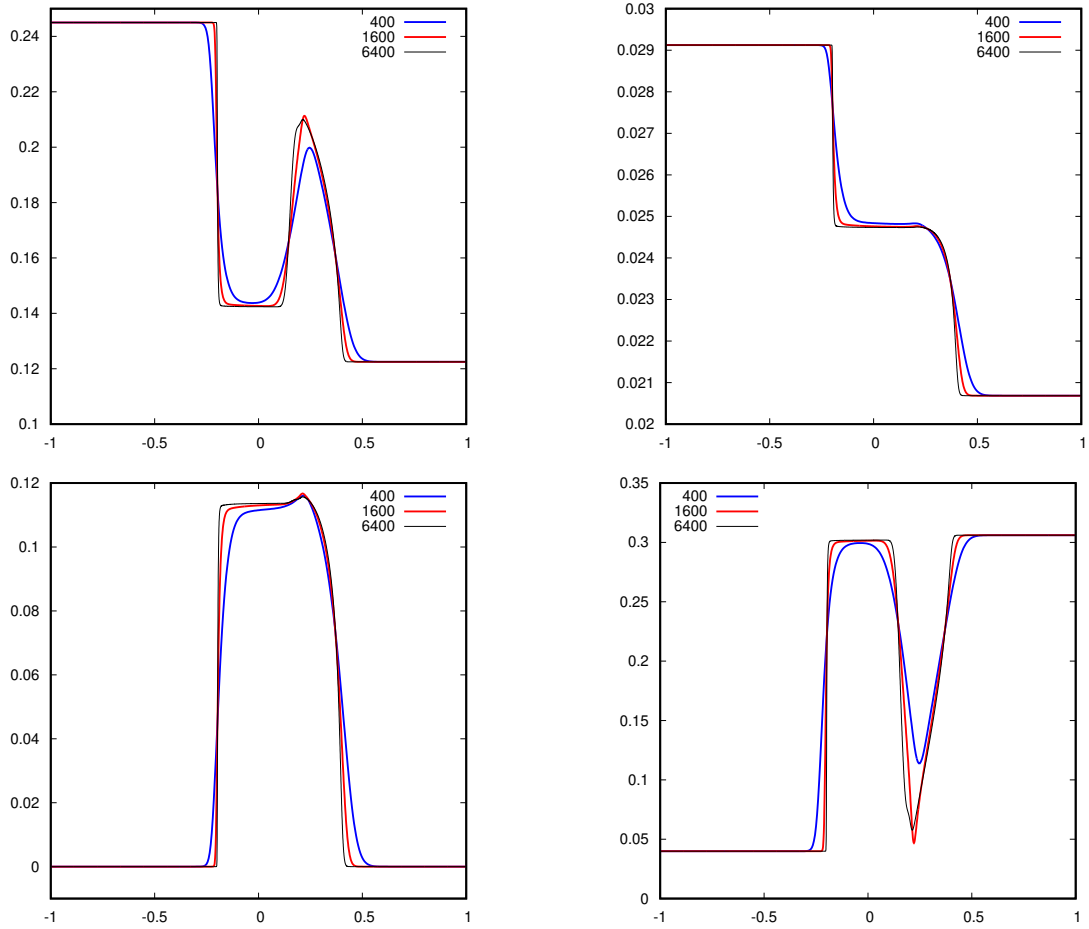


Figure 7.2: Plots for Test 1 of the underestimation problem using \hat{p}^* for computing the maximum wave speed. (Top left): density, (top right): pressure, (bottom left): velocity, (bottom right): sound speed.

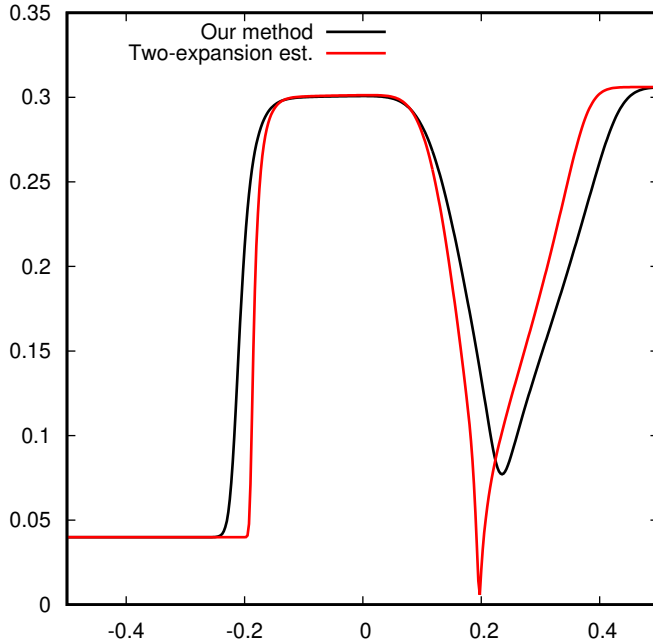


Figure 7.3: Comparison of the sound speed for our method versus the two expansion approximation. This test was run with 800 DoFs up to time $t = 1.2$. The simulation using the two-expansion estimate immediately crashes after $t = 1.2$ generating complex sound speed.

possible range of time step sizes. The Riemann data is,

$$(\rho_L, v_L, p_L) := (0.9932, 3, 2), \quad (7.14)$$

$$(\rho_R, v_R, p_R) := (0.9500, -3, 2). \quad (7.15)$$

The corresponding sound speeds are $(a_L, a_R) \approx (21.2, 7.77)$. The computational domain is $D = (-1.7, 1)$ and the simulation is run to $t_{\text{final}} = 0.005$. When using λ^{exp} we find that the simulation requires a CFL of at most 0.12 in order to maintain positive internal energy. The maximal CFL number we were able to use, using the wave speed, $\hat{\lambda}_{\text{max}} = \lambda_{\text{max}}(\hat{p}^*)$, computed by Algorithm 1, while still preserving positive internal energy and real sound speed was 1.42. R (The sound speed and internal energy are checked at every node and at every time step.) Plots of the solution are shown in Figure 7.5. In this case, using the method with λ^{exp} results in a computational time approximately 12 times greater than the method using $\lambda_{\text{max}}(\hat{p}^*)$.

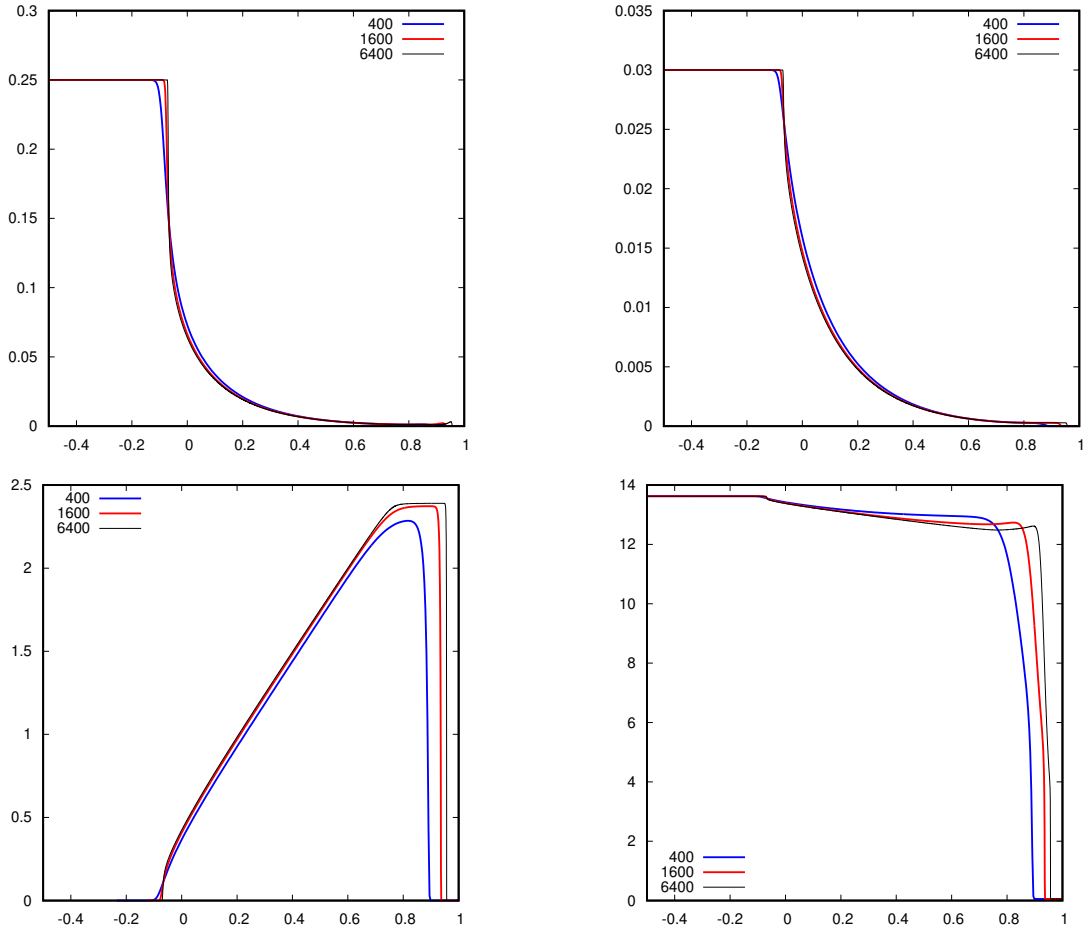


Figure 7.4: Plots for Test 2 of the underestimation problem using \hat{p}^* for computing the maximum wave speed. (Top left): density, (top right): pressure, (bottom left): velocity, (bottom right): specific internal energy. Final time is $t = 0.4$.

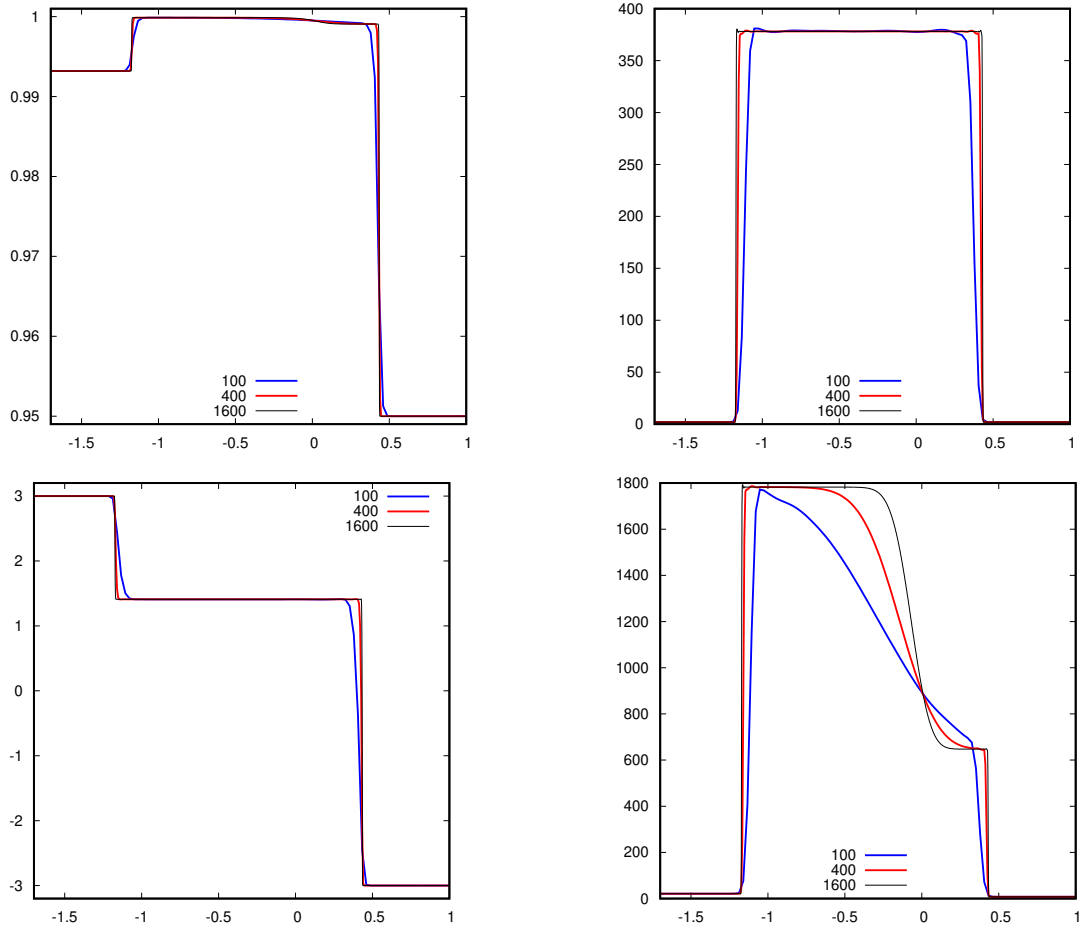


Figure 7.5: Plots for Test 3 of the overestimation problem using \hat{p}^* for computing the maximum wave speed. From left to right: density, pressure, velocity, sound speed. Final time is $t = 0.005$ and using $\text{CFL} = 1.42$. Figures were similar when using λ^{exp} .

7.3 The SESAME Database

In order to demonstrate the robustness of our method, we use the **SESAME** database. The **SESAME** database is a tabulated EOS developed by the Mechanics of Materials Equation of State group (T-1) at Los Alamos National Lab [81]. The **SESAME** database houses data on a large variety of materials from argon to cesium. The construction of such a database relies on experimental data and EOS models fitted to the material. However, such details are beyond the scope of this thesis. Access to this database can be acquired by contacting `sesame@lanl.gov`. The database alone is not usable by itself, we need a way to interface with it. This is done with **EOSPAC6** [82]; all versions of this software can also be found at <https://laws.lanl.gov/projects/data/eos/eospacReleases.php>.

The general process by which **EOSPAC6** functions, is that we provide a material identification number and “table type”. The “table type” indicates the thermodynamic relationship; that is, it specifies the two thermodynamic input quantities and the corresponding output quantity. For our purposes, we provide $\{(\rho_i^n, e_i^n)\}_{i \in \mathcal{V}}$ and **EOSPAC6** returns $\{p_i^n\}_{i \in \mathcal{V}}$. Recall for density, pressure and specific internal energy in the compressible Euler equations are, kg/m^3 , Pa, and J/kg, respectively. However, **EOSPAC6** works with the following respective units, Mg/m^3 , GPa, and MJ/kg. A full description of the functionality of **EOSPAC6** can be found in the user manual in Pimentel [83]. Additionally, the list of the available materials in the **SESAME** database with their corresponding identification numbers are found on the last pages of [81].

We conduct a series of tests utilizing the **SESAME** database. The Riemann problems all use $x = 0$ as the initial discontinuity. Furthermore, we offer comparisons of the different values of $\widehat{\lambda}_{\max}(\mathbf{n}_{ij}, \mathbf{U}_i, \mathbf{U}_j)$. As $\widehat{\lambda}_{\max}(\mathbf{n}_{ij}, \mathbf{U}_i, \mathbf{U}_j)$ depends on the Riemann problem we use the definition,

$$\widetilde{\lambda}_i^{\max} := \max_{j \in \mathcal{G}(i) \setminus \{i\}} \widehat{\lambda}_{\max}(\mathbf{n}_{ij}, \mathbf{U}_i, \mathbf{U}_j), \quad (7.16)$$

in order to plot the maximum wave speed.

7.3.1 Expansion-Contact-Shock Comparison

For the first problem, we run a Riemann problem using dry air (material id: 5030), aluminum (material ID: 3720), vanadium (material ID: 2552), and sulfur hexafluoride (material ID: 7010) with data which generates a wave profile similar to the Sod shock tube. The Riemann data is,

$$(\rho_L, v_L, e_L) := (0.01 \text{ Mg/m}^3, 0 \text{ m/s}, 4000 \text{ MJ/kg}), \quad (7.17)$$

$$(\rho_R, v_R, e_R) := (0.003 \text{ Mg/m}^3, 0 \text{ m/s}, 3400 \text{ MJ/kg}). \quad (7.18)$$

The problem is simulated on $D = (-1 \text{ m}, 1 \text{ m})$ to a final time of $t = 1.2 \times 10^{-5} \text{ s}$ and $\text{CFL} = 0.5$. See Figure 7.6 for the results.

7.4 Benchmark Configurations

7.4.1 EOS Comparison in a Riemann Problem

In this simulation we compare the different effects an EOS has the solution to a Riemann problem. The initial discontinuity is at $x = 0$, the Riemann data we use is,

$$(\rho_L, v_L, p_L) := (1, 1, 2), \quad (7.19)$$

$$(\rho_R, v_R, p_R) := (1, -1, 1), \quad (7.20)$$

on $D = (-0.5, 0.5)$ with $\text{CFL} = 0.5$ to the final time of $t_{\text{final}} = 0.1$. We compare results with the ideal, covolume, van der Waals and EOS. The parameters we use are, $\gamma = 1.4$, $a = 0.5$, $b = 0.5$, $R = 0.4$, and $c_V = 1$ for each of the associated EOS. The results are shown in Figures 7.7.

7.4.2 The Woodward-Colella Blast Wave

A common test problem in the literature is the Woodward-Colella blast wave which was originally proposed in [84]. The initial configuration is composed of three constant states

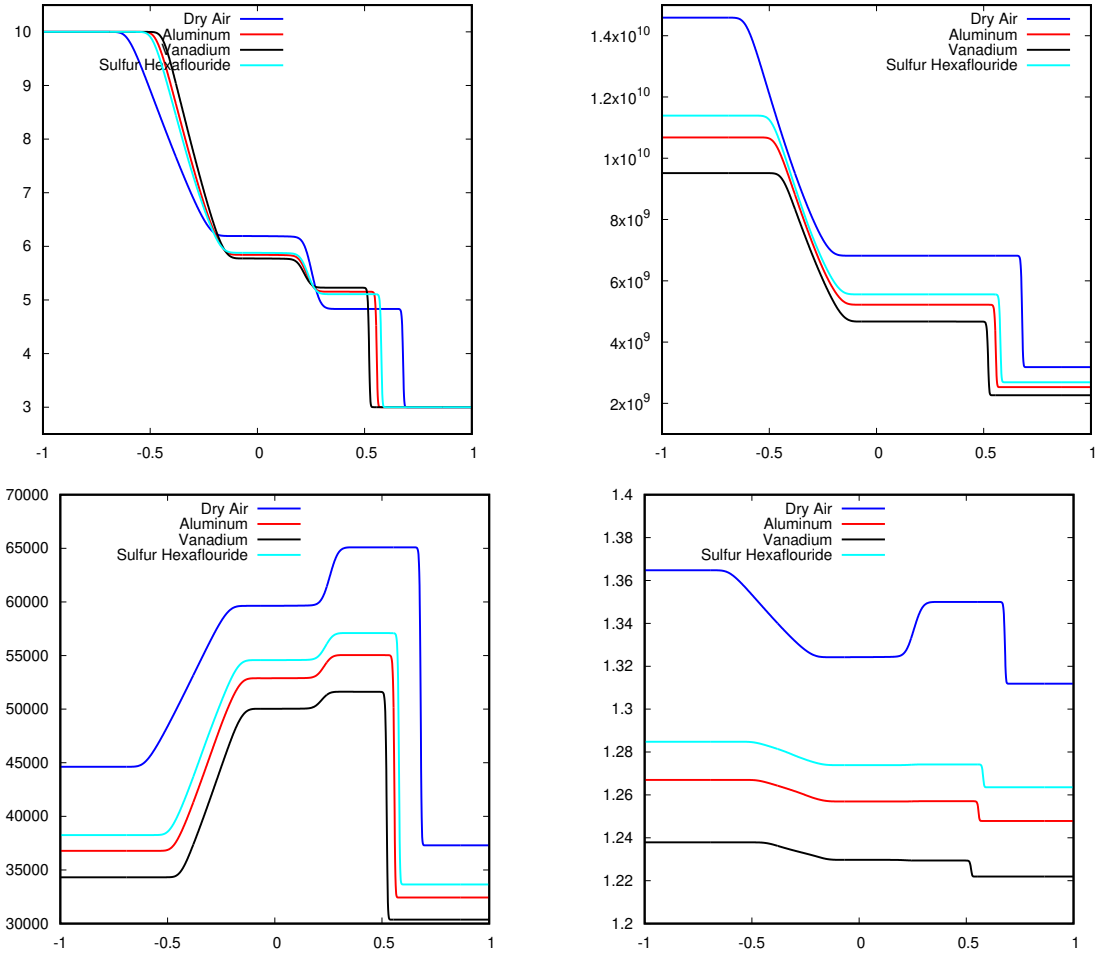


Figure 7.6: Comparison of ρ (top left), p (top right), max wave speed, $\tilde{\lambda}^{\max}$ (bottom left), and γ (bottom right), for the various materials at the final time, $t = 1.2 \times 10^{-5}$ s.

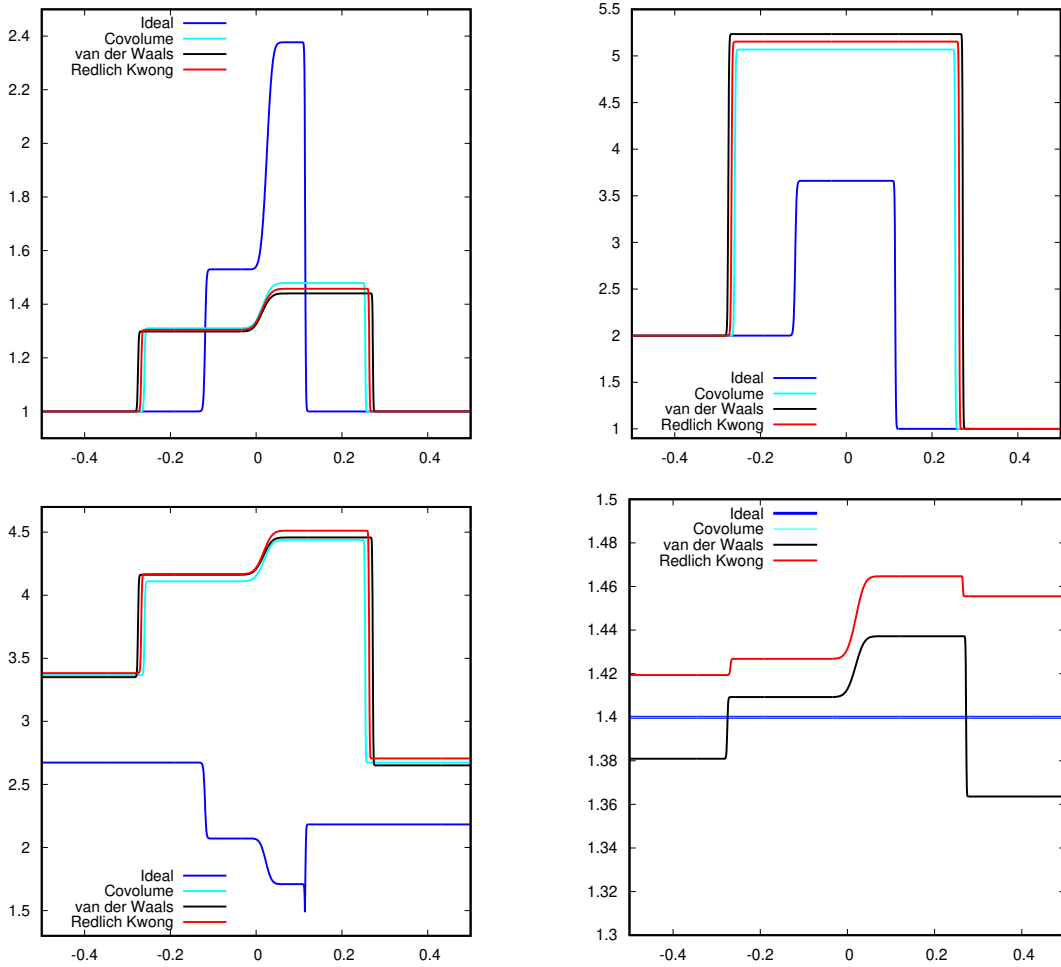


Figure 7.7: Comparison of the density (top left), pressure (top right), max wave speed (bottom left), and γ (bottom right), for the various different EOS at the final time, $t = 0.1$.

which result in the collision of two waves. This collision produces an emergent complex structure. This test problem is normally performed with the ideal EOS; we instead perform two experiments, both using the Jones-Wilkins-Lee EOS. The first experiment uses the parameters proposed in [23, Section 5.1] and the second experiment uses the parameters in [44, Table 2. “HMX”], see Table 7.4 for the values.

The computational domain is $D = (0, 1)$ and the initial state is,

$$(\rho_0(x), v_0(x), p_0(x)) = \begin{cases} (1, 0, 10^3) & \text{if } x \in [0, 0.1], \\ (1, 0, 10^{-2}) & \text{if } x \in (0.1, 0.9), \\ (1, 0, 10^2) & \text{if } x \in [0.9, 1]. \end{cases} \quad (7.21)$$

and the Jones-Wilkins-Lee parameters for the two experiments are given in Table 7.4.

	A	B	R_1	R_2	ω	ρ_0	t_{final}
Case 1	6.321×10^2	-4.472	11.3	1.13	0.8938	1	0.038
Case 2	7.7828×10^{11}	7.071428×10^9	4.2	1.00	0.3000	1891	0.038

Table 7.4: JWL parameters for Woodward-Colella interacting blast wave benchmark. Reprinted with permission from [1].

We use CFL = 0.9. The results are recorded in Figure 7.8 and Figure 7.9.

7.4.3 Shock Collision with Triangular Obstacle

The next benchmark problem is referred to in the literature as the Schardin’s problem, see Schardin [85] for the original physical experiment. This problem is for studying the effects of a shock wave colliding with a triangular obstacle. This physical experiment has been recreated in [86] with a detailed description of the experimental setup.

For our numerical demonstration we use the setup by Toro et. al. in [23, Sec. 5.4]. The computational domain is $D = (-0.65, 0.5) \times (-0.5, 0.5) \setminus K$, where K is the triangle formed

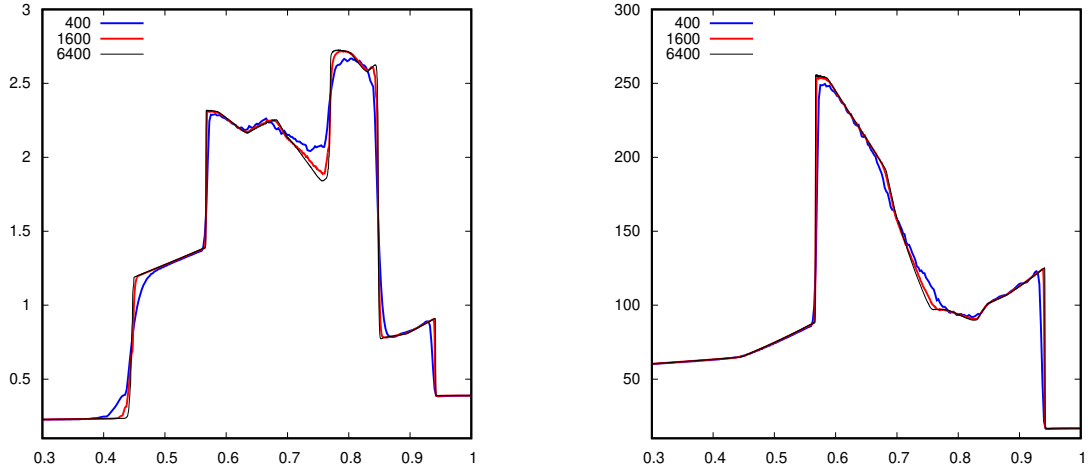


Figure 7.8: Case 1 of the Woodward-Colella blast wave with the JWL EOS. (Left) density, (right) pressure.

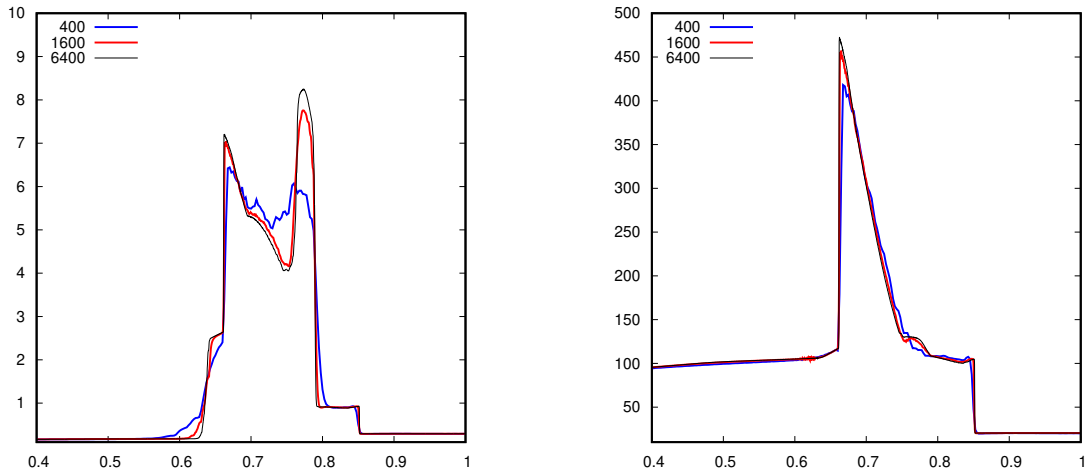


Figure 7.9: Case 2 of the Woodward-Colella blast wave with the JWL EOS. (Left) density, (right) pressure.

by the vertices $(-0.2, 0)$, $(0.1, 1/6)$, and $(0.1, -1/6)$. We enforce Dirichlet conditions on the left boundary, dynamic outflow conditions on the right boundary and slip conditions on the remaining boundaries. Note however that the simulation is stopped before the wave reaches the outflow boundary so one could alternatively enforce Dirichlet on the right boundary. The EOS is van der Waals with $\gamma = 864.7/577.8$, $a = 0.14$ and $b = 3.258 \times 10^{-5}$. For the

interpolation parameters, we use the given b and set $q = p_\infty = 0$.

The ambient state in front of the shock is set to $\mathbf{u}_R = (1.225, 0, 0, 101325)^\top$. We define the Mach speed of the shock to be $M_S := 1.3$ and then compute the state across the shock using the Rankine-Hugoniot conditions, (1.24). The complete initial state is given by,

$$(\rho_0(\mathbf{x}), \mathbf{v}_0(\mathbf{x}), p_0(\mathbf{x})) := \begin{cases} (1.82039, 148.597, 0, 185145), & \text{if } x \leq -0.55, \\ (1.225, 0, 0, 101325), & \text{if } x > -0.55. \end{cases} \quad (7.22)$$

Schlieren plots of the solution can be seen Figure 7.10.

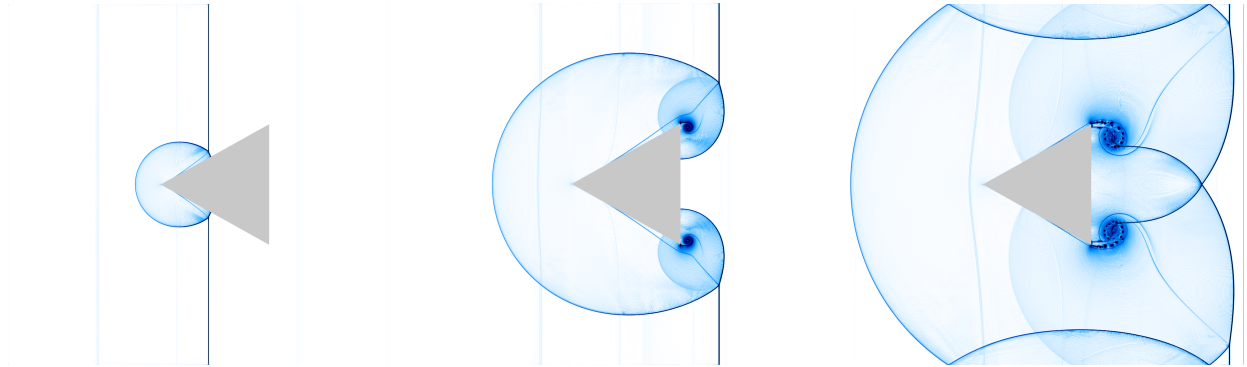


Figure 7.10: Schlieren plot of a shock wave interacting with a triangular obstacle at $t = 1$ ms, 1.6 ms, and 2.2 ms. Reprinted with permission from [1].

7.4.4 Shock Bubble Interaction

The next test problem we simulate is the shock-bubble interaction. Much like the Schardin problem done in Section 7.4.3, we again simulate a shock colliding with an obstacle. This time, the obstacle is not a solid boundary but a bubble. The density of the bubble. This problem is normally performed in the context of multi-fluids, where the bubble is a separate material from the ambient fluid. The ambient fluid is typically chosen to be air and the bubble to be helium. More information on the physical experiments can be found

in [87]. We also point out the paper, by Quirk & Karni [88], where they provide comparison between the numerical and physical experiments.

In our case the bubble only differs from the ambient fluid by a difference in density. We use the setup given Wang & Li in [89, Sec. 5.2.2]. The domain is $D = (0, 3) \times (0, 1)$. The left boundary is Dirichlet, the right boundary is dynamic outflow, and the top and bottom boundaries are slip. Let \mathfrak{B} denote the bubble centered at $(0.5, 0.5)$ with radius 0.25. We use the Jones-Wilkins-Lee EOS with the following parameters,

$$A = 8.545 \times 10^{11}, B = 2.05 \times 10^{10}, R_1 = 4.6, R_2 = 1.35, \omega = 0.25, \rho_0 = 1.84 \times 10^3. \quad (7.23)$$

The primitive states for the ambient fluid and the bubble are,

$$(\rho_{\mathbf{R}}, \mathbf{v}_{\mathbf{R}}, p_{\mathbf{R}}) = (1000, 0, 0, 5 \times 10^{10}) \quad (7.24a)$$

$$(\rho_{\mathfrak{B}}, \mathbf{v}_{\mathfrak{B}}, p_{\mathfrak{B}}) = (2000, 0, 0, 5 \times 10^{10}) \quad (7.24b)$$

We prescribe the pressure across the shock wave to be $p_L = 4.369 \times 10^{11}$ and determine the remaining variables ρ_L and \mathbf{v}_L through the Rankine-Hugoniot conditions, (1.24). The primitive initial state is,

$$(\rho_0(\mathbf{x}), \mathbf{v}_0(\mathbf{x}), p_0(\mathbf{x})) := \begin{cases} (3778.85, 16867.6, 0, 4.369 \times 10^{11}), & \text{if } x < 0.05, \\ (1000, 0, 0, 5 \times 10^{10}), & \text{if } x \geq 0.5 \text{ and } \mathbf{x} \notin \mathfrak{B}, \\ (2000, 0, 0, 5 \times 10^{10}), & \text{if } \mathbf{x} \in \mathfrak{B}. \end{cases} \quad (7.25)$$

A diagram of this initial state is shown in Figure 7.11. The simulation is run to final time $t_{\text{final}} = 100 \mu\text{s}$ with $\text{CFL} = 0.5$. Schlieren plots of the density are shown in Figures 7.12, 7.13, and 7.13 at times $t = 40 \mu\text{s}$, $70 \mu\text{s}$, and $100 \mu\text{s}$, respectively. The simulation was run with 50,348,033 \mathbb{Q}_1 nodes. Our numerical results match well with experiment results shown in [89].

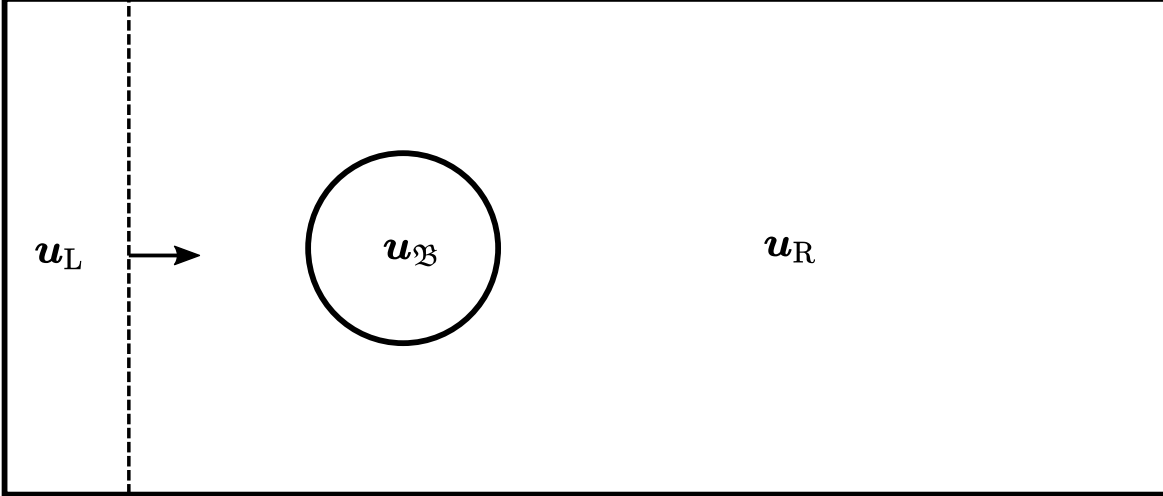


Figure 7.11: Visual description of the initial state for the shock bubble interaction.

7.4.5 Shock Diffraction

This problem involves shock diffracting around a 90° wedge. Physical experiments for a plethora of geometric wedges can be found in [90]. This problem has been studied extensively in the literature for an ideal polytropic gas, see the collection of posters from the 18th International Symposium on Shock Waves (ISSW) [91]. However the results for a non-ideal fluid are much more sparse. Most notably this problem has been simulated for Bethe-Zel'dovich-Thompson (BZT) fluids with the van der Waals equation of state in [92, Problem TD3].

Remark 7.4.1 (BZT Fluids). The Bethe-Zel'dovich-Thompson (BZT) fluids are fluids for which there exists a thermodynamic region where the fundamental derivative,

$$\mathcal{G} := -\frac{\tau}{2} \left[\left(\frac{\partial^2 p}{\partial \tau^2} \right)_s / \left(\frac{\partial p}{\partial \tau} \right)_s \right] = 1 + \frac{\rho}{a} \left(\frac{\partial a}{\partial \rho} \right)_s = \frac{\tau^3}{2a^2} \left(\frac{\partial^2 p}{\partial \tau^2} \right)_s, \quad (7.26)$$

is negative. Fluids that exhibit a negative fundamental region behave in an unusual way. For one, the speed of sound increases as the density decreases. This causes the existence of expansion shocks, also called composite waves, where a single wave consists of an expansion

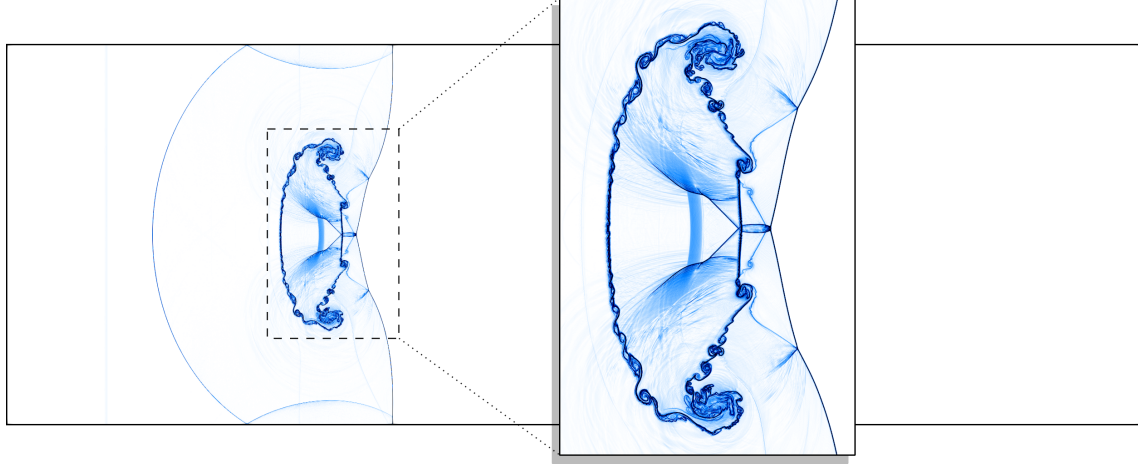


Figure 7.12: Schlieren plots for the shock-bubble interaction benchmark for $t = 40 \mu\text{s}$. Reprinted with permission from [1].

immediately followed by a shock or vice-versa. For more information see [93] and [94]. \square

All problems are simulated on the domain $\Omega = [0, 2]^2 \setminus ([0, 0.5] \times [0, 1])$. The first test is run with the van der Waals equation of state where $a = 2.2 \text{ Pa m}^6 \text{ kg}^{-2}$, $b = 7.25 \times 10^{-4} \text{ m}^3 \text{ kg}^{-1}$, and $\gamma = 1.0125$. We follow the same setup of Brown & Argrow [92, Problem TD3]; however, they use the nondimensionalized forms and we instead work with dimensional variables. It should also be noted that, despite their claim, the left and right states of Problem TD3 do not satisfy the Rankine-Hugoniot conditions. We choose the right state and the left pressure, and derive the remaining state variables through the Rankine-Hugoniot conditions (1.24). The initial condition is,

$$\begin{pmatrix} \rho_0(\mathbf{x}) \\ v_0^1(\mathbf{x}) \\ v_0^2(\mathbf{x}) \\ p_0(\mathbf{x}) \end{pmatrix} := \begin{cases} (373.108 \text{ kg/m}^3, 18.4156 \text{ m/s}, 0 \text{ m/s}, 156\,568 \text{ Pa})^\top, & x < 0.5 \\ (128.736 \text{ kg/m}^3, 0 \text{ m/s}, 0 \text{ m/s}, 89\,910.6 \text{ Pa})^\top, & x \geq 0.5. \end{cases} \quad (7.27)$$

Schlieren plots of the solution are shown in Figure 7.15.

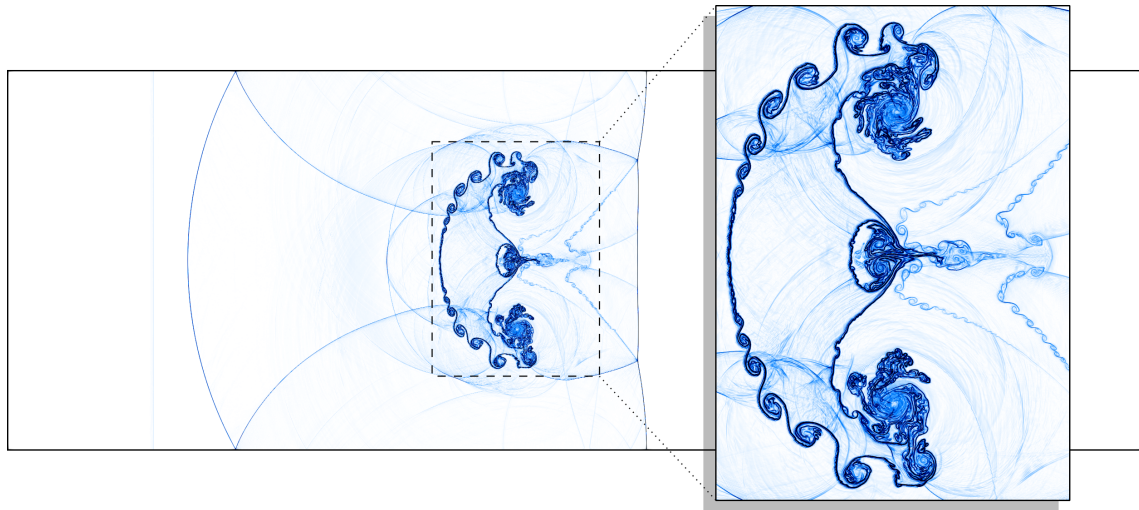


Figure 7.13: Schlieren plots for the shock-bubble interaction benchmark for $t = 70 \mu s$. Reprinted with permission from [1].

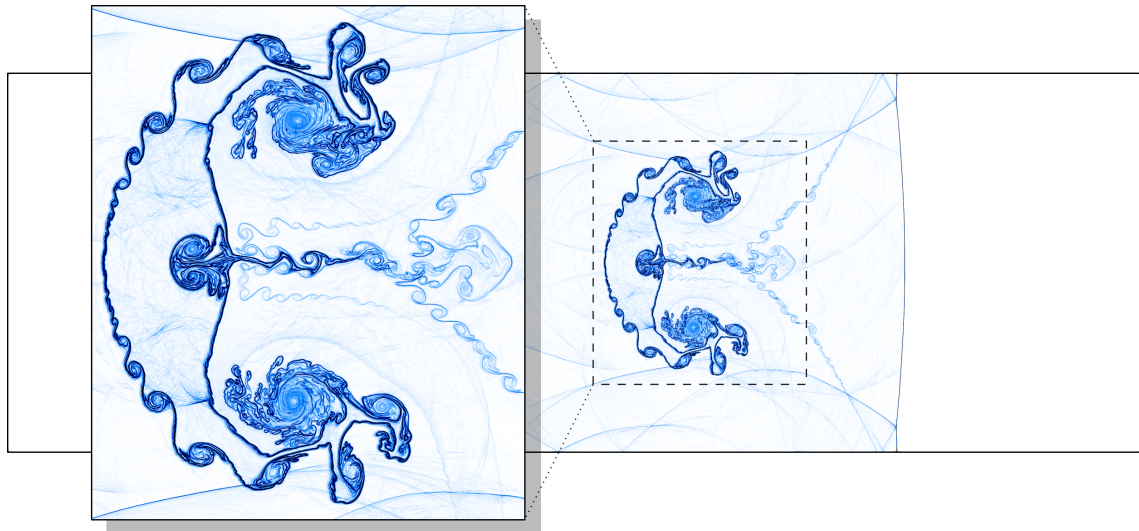


Figure 7.14: Schlieren plots for the shock-bubble interaction benchmark for $t = 100 \mu s$ (bottom). [1].

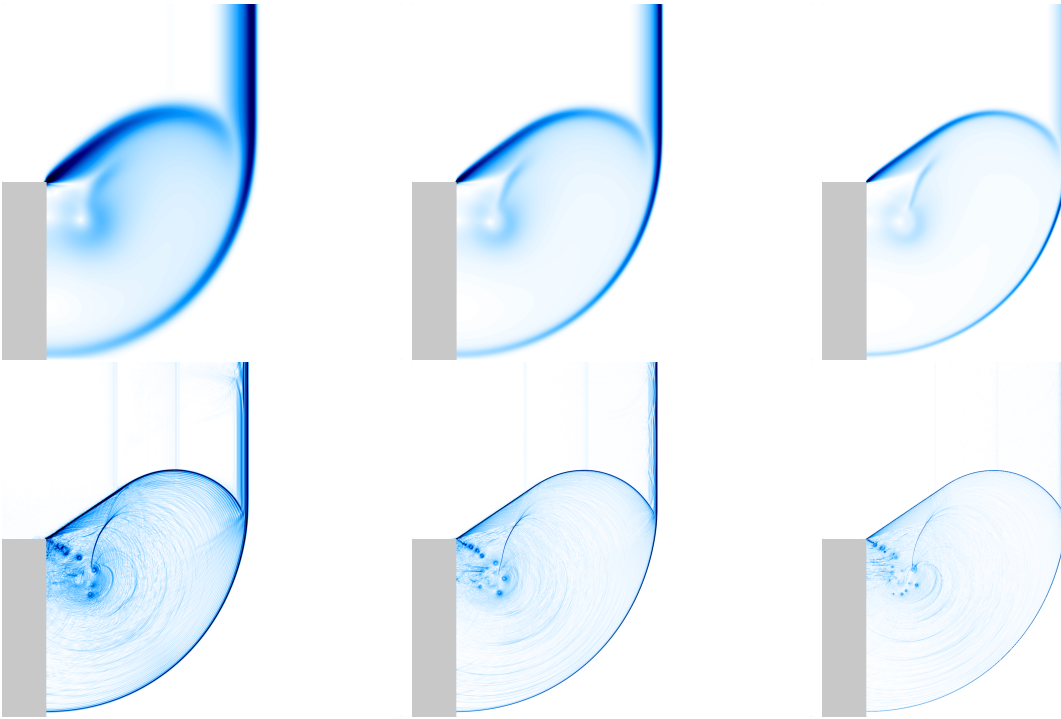


Figure 7.15: Comparison of numerical Schlieren plots for first-order accurate (top row) and second-order accurate (bottom row) solutions at time $t_{\text{final}} = 0.02$ s. The mesh resolution increases from left to right as follows: 1,116,289, 4,460,801, and 17,834,497 Q_1 -nodes. Reprinted with permission from [1].

REFERENCES

- [1] B. Clayton, J.-L. Guermond, M. Maier, B. Popov, and E. J. Tovar, “Robust second-order approximation of the compressible euler equations with an arbitrary equation of state,” Journal of Computational Physics, p. 111926, 2023.
- [2] B. Clayton, J.-L. Guermond, and B. Popov, “Invariant domain-preserving approximations for the Euler equations with tabulated equation of state,” SIAM Journal on Scientific Computing, vol. 44, no. 1, pp. A444–A470, 2022.
- [3] J. von Neumann and R. D. Richtmyer, “A method for the numerical calculation of hydrodynamic shocks,” Journal of applied physics, vol. 21, no. 3, pp. 232–237, 1950.
- [4] P. D. Lax, “Weak solutions of nonlinear hyperbolic equations and their numerical computation,” Communications on pure and applied mathematics, vol. 7, no. 1, pp. 159–193, 1954.
- [5] P. Lax and B. Wendroff, “Systems of conservation laws,” Communications on Pure and Applied Mathematics, vol. 13, no. 2, pp. 217–237, 1960.
- [6] S. K. Godunov and I. Bohachevsky, “Finite difference method for numerical computation of discontinuous solutions of the equations of fluid dynamics,” Matematičeskij sbornik, vol. 47, no. 3, pp. 271–306, 1959.
- [7] R. W. MacCormack, “The effect of viscosity in hypervelocity impact cratering,” Journal of spacecraft and rockets, vol. 40, no. 5, pp. 757–763, 2003.
- [8] J. L. Steger and R. Warming, “Flux vector splitting of the inviscid gasdynamic equations with application to finite-difference methods,” Journal of computational physics, vol. 40, no. 2, pp. 263–293, 1981.
- [9] B. Van Leer, “Flux-vector splitting for the euler equations,” in Eighth International Conference on Numerical Methods in Fluid Dynamics: Proceedings of the Conference, Rheinisch-Westfälische Technische Hochschule Aachen, Germany, June 28–July 2, 1982, pp. 507–512, Springer, 2005.

- [10] P. L. Roe, “Approximate riemann solvers, parameter vectors, and difference schemes,” Journal of computational physics, vol. 43, no. 2, pp. 357–372, 1981.
- [11] A. Harten, P. D. Lax, and B. V. Leer, “On upstream differencing and godunov-type schemes for hyperbolic conservation laws,” SIAM Review, vol. 25, no. 1, pp. 35–61, 1983.
- [12] E. F. Toro, M. Spruce, and W. Speares, “Restoration of the contact surface in the hll-riemann solver,” Shock waves, vol. 4, pp. 25–34, 1994.
- [13] M.-S. Liou, “The evolution of ausm schemes,” Defence Science Journal, vol. 60, no. 6, 2010.
- [14] A. Harten, B. Engquist, S. Osher, and S. R. Chakravarthy, “Uniformly high order accurate essentially non-oscillatory schemes, iii,” Journal of computational physics, vol. 131, no. 1, pp. 3–47, 1997.
- [15] G.-S. Jiang and C.-W. Shu, “Efficient implementation of weighted eno schemes,” Journal of computational physics, vol. 126, no. 1, pp. 202–228, 1996.
- [16] J.-L. Guermond and B. Popov, “Invariant domains and first-order continuous finite element approximation for hyperbolic systems,” SIAM Journal on Numerical Analysis, vol. 54, no. 4, pp. 2466–2489, 2016.
- [17] P. Colella and H. M. Glaz, “Efficient solution algorithms for the Riemann problem for real gases,” J. Comput. Phys., vol. 59, no. 2, pp. 264–289, 1985.
- [18] M. J. Ivings, D. M. Causon, and E. F. Toro, “On Riemann solvers for compressible liquids,” Internat. J. Numer. Methods Fluids, vol. 28, no. 3, pp. 395–418, 1998.
- [19] L. Quartapelle, L. Castelletti, A. Guardone, and G. Quaranta, “Solution of the Riemann problem of classical gasdynamics,” J. Comput. Phys., vol. 190, no. 1, pp. 118–140, 2003.
- [20] J. K. Dukowicz, “A general, noniterative Riemann solver for Godunov’s method,” Journal of Computational Physics, vol. 61, no. 1, pp. 119–137, 1985.
- [21] P. L. Roe and J. Pike, “Efficient construction and utilisation of approximate riemann solutions,” in Proceedings of the Sixth International Symposium on Computing Methods

- in Applied Sciences and Engineering, VI, (Netherlands), pp. 499–518, North-Holland Publishing Co., 1985.
- [22] J. Pike, “Riemann solvers for perfect and near-perfect gases,” AIAA Journal, vol. 31, no. 10, pp. 1801–1808, 1993.
- [23] E. F. Toro, C. E. Castro, and B. J. Lee, “A novel numerical flux for the 3d euler equations with general equation of state,” Journal of Computational Physics, vol. 303, pp. 80–94, 2015.
- [24] K. N. Chueh, C. C. Conley, and J. A. Smoller, “Positively invariant regions for systems of nonlinear diffusion equations,” Indiana University Mathematics Journal, vol. 26, no. 2, pp. 373–392, 1977.
- [25] D. Hoff, “Invariant regions for systems of conservation laws,” Transactions of the American Mathematical Society, vol. 289, no. 2, pp. 591–610, 1985.
- [26] H. Frid, “Maps of convex sets and invariant regions for finite-difference systems of conservation laws,” Archive for rational mechanics and analysis, vol. 160, pp. 245–269, 2001.
- [27] J. Glimm, “Solutions in the large for nonlinear hyperbolic systems of equations,” Communications on pure and applied mathematics, vol. 18, no. 4, pp. 697–715, 1965.
- [28] L. C. Evans, Partial differential equations, vol. 19. American Mathematical Society, 2022.
- [29] S. N. Kružkov, “First order quasilinear equations in several independent variables,” Mathematics of the USSR-Sbornik, vol. 10, pp. 217–243, feb 1970.
- [30] S. Bianchini and A. Bressan, “Vanishing viscosity solutions of nonlinear hyperbolic systems,” Annals of mathematics, pp. 223–342, 2005.
- [31] E. Godlewski and P. Raviart, Numerical Approximation of Hyperbolic Systems of Conservation Laws. Applied Mathematical Sciences, Springer New York, 2021.
- [32] B. Riemann, Über die Fortpflanzung ebener Luftwellen von endlicher Schwingungsweite, vol. 8. Verlag der Dieterichschen Buchhandlung, 1860.

- [33] J. N. Johnson and R. Chéret, Classic papers in shock compression science. Springer Science & Business Media, 2012.
- [34] P. D. Lax, “Hyperbolic systems of conservation laws ii,” Communications on pure and applied mathematics, vol. 10, no. 4, pp. 537–566, 1957.
- [35] C. M. Dafermos, Hyperbolic conservation laws in continuum physics, vol. 3. Springer, 2005.
- [36] H. Callen, Thermodynamics and an introduction to thermostatistics, 2nd ed. John Wiley & Sons, New York, 1985. Second edition.
- [37] R. Menikoff and B. J. Plohr, “The riemann problem for fluid flow of real materials,” Reviews of modern physics, vol. 61, no. 1, p. 75, 1989.
- [38] O. Le Métayer and R. Saurel, “The noble-abel stiffened-gas equation of state,” Physics of Fluids, vol. 28, no. 4, p. 046102, 2016.
- [39] A. Chiapolino and R. Saurel, “Extended noble–abel stiffened-gas equation of state for sub-and-supercritical liquid-gas systems far from the critical point,” Fluids, vol. 3, no. 3, p. 48, 2018.
- [40] I. A. Johnston, “The noble-abel equation of state: Thermodynamic derivations for ballistics modelling,” tech. rep., Defence Science and Technology Organisation Edinburgh (Australia) Weapons . . . , 2005.
- [41] J. D. Van der Waals, Over de Continuïteit van den Gas-en Vloeistofoestand, vol. 1. Sijthoff, 1873.
- [42] J. O. Valderrama, “The state of the cubic equations of state,” Industrial & engineering chemistry research, vol. 42, no. 8, pp. 1603–1618, 2003.
- [43] O. Redlich and J. N. Kwong, “On the thermodynamics of solutions. v. an equation of state. fugacities of gaseous solutions.,” Chemical reviews, vol. 44, no. 1, pp. 233–244, 1949.
- [44] E. L. Lee, H. C. Hornig, and J. W. Kury, “Adiabatic expansion of high explosive detonation products,” tech. rep., Univ. of California Radiation Lab. at Livermore, 1968.

- [45] S. B. Segletes, “An examination of the jwl equation of state,” tech. rep., Army Research Lab Aberdeen Proving Ground, 2018.
- [46] R. Menikoff, “Complete mie-gruneisen equation of state (update),” tech. rep., Los Alamos National Lab, 3 2016.
- [47] R. Menikoff, “Empirical equations of state for solids,” ShockWave Science and Technology Reference Library, pp. 143–188, 2007.
- [48] C. Grossmann, H.-G. Roos, and M. Stynes, Numerical treatment of partial differential equations, vol. 154. Springer, 2007.
- [49] S. Larsson and V. Thomée, Partial differential equations with numerical methods, vol. 45. Springer, 2003.
- [50] P. G. Ciarlet, The Finite Element Method for Elliptic Problems. Society for Industrial and Applied Mathematics, 2002.
- [51] A. Ern and J.-L. Guermond, Finite Elements I: Approximation and Interpolation, vol. 72. Springer, 2021.
- [52] A. Ern and J.-L. Guermond, Finite Elements II: Galerkin Approximation, Elliptic and Mixed PDEs, vol. 73. Springer, 2021.
- [53] A. Ern and J.-L. Guermond, Finite Elements III: First-Order and Time-Dependent PDEs, vol. 74. Springer, 2021.
- [54] A. Ern and J.-L. Guermond, “Invariant-domain-preserving high-order time stepping: I. explicit runge–kutta schemes,” SIAM Journal on Scientific Computing, vol. 44, no. 5, pp. A3366–A3392, 2022.
- [55] J.-L. Guermond, B. Popov, and I. Tomas, “Invariant domain preserving discretization-independent schemes and convex limiting for hyperbolic systems,” Computer Methods in Applied Mechanics and Engineering, vol. 347, pp. 143–175, 2019.
- [56] E. F. Toro, Riemann Solvers and Numerical Methods for Fluid Dynamics. Springer-Verlag Berlin Heidelberg, 2009.
- [57] H. Holden and N. H. Risebro, Front tracking for hyperbolic conservation laws, vol. 152.

Springer, 2015.

- [58] R. Abgrall and S. Karni, “Computations of compressible multifluids,” J. Comput. Phys., vol. 169, no. 2, pp. 594–623, 2001.
- [59] J.-L. Guermond and B. Popov, “Fast estimation of the maximum wave speed in the riemann problem for the euler equations,” arXiv preprint arXiv:1511.02756, 2015.
- [60] J. L. Guermond, M. Nazarov, B. Popov, and I. Tomas, “Second-order invariant domain preserving approximation of the euler equations using convex limiting,” SIAM Journal on Scientific Computing, vol. 40, pp. A3211–A3239, 2018.
- [61] J.-L. Guermond, B. Popov, and Y. Yang, “The effect of the consistent mass matrix on the maximum-principle for scalar conservation equations,” Journal of Scientific Computing, vol. 70, pp. 1358–1366, 2017.
- [62] J.-L. Guermond and R. Pasquetti, “A correction technique for the dispersive effects of mass lumping for transport problems,” Computer Methods in Applied Mechanics and Engineering, vol. 253, pp. 186–198, 2013.
- [63] M. A. Christon, M. J. Martinez, and T. E. Voth, “Generalized Fourier analyses of the advection-diffusion equation-part I: one-dimensional domains,” International Journal for Numerical Methods in Fluids, vol. 45, no. 8, pp. 839–887, 2004.
- [64] T. Thompson, “A discrete commutator theory for the consistency and phase error analysis of semi-discrete c0 finite element approximations to the linear transport equation,” Journal of Computational and Applied Mathematics, vol. 303, pp. 229–248, 2016.
- [65] M. Maier and M. 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, vol. 8, no. 3, pp. 16:1–30, 2021.
- [66] A. Harten, P. D. Lax, C. D. Levermore, and W. J. Morokoff, “Convex entropies and hyperbolicity for general euler equations,” SIAM Journal on Numerical Analysis, vol. 35, no. 6, pp. 2117–2127, 1998.
- [67] A. Harten, “On the symmetric form of systems of conservation laws with entropy,”

- Journal of Computational Physics, vol. 49, no. 1, pp. 151–164, 1983.
- [68] J. P. Boris and D. L. Book, “Flux-corrected transport. i. shasta, a fluid transport algorithm that works,” Journal of Computational Physics, vol. 11, no. 1, pp. 38–69, 1973.
- [69] D. Book, J. Boris, and K. Hain, “Flux-corrected transport ii: Generalizations of the method,” Journal of Computational Physics, vol. 18, no. 3, pp. 248–283, 1975.
- [70] J. Boris and D. Book, “Flux-corrected transport. iii. minimal-error fct algorithms,” Journal of Computational Physics, vol. 20, no. 4, pp. 397–431, 1976.
- [71] S. T. Zalesak, “Fully multidimensional flux-corrected transport algorithms for fluids,” Journal of computational physics, vol. 31, no. 3, pp. 335–362, 1979.
- [72] S. T. Zalesak, The design of Flux-Corrected Transport (FCT) algorithms for structured grids. Springer, 2012.
- [73] C. Lohmann and D. Kuzmin, “Synchronized flux limiting for gas dynamics variables,” Journal of Computational Physics, vol. 326, pp. 973–990, 2016.
- [74] E. Tadmor, “A minimum entropy principle in the gas dynamics equations,” Applied Numerical Mathematics, vol. 2, no. 3, pp. 211–219, 1986. Special Issue in Honor of Milt Rose’s Sixtieth Birthday.
- [75] B. Khobalatte and B. Perthame, “Maximum principle on the entropy and second-order kinetic schemes,” Mathematics of Computation, vol. 62, no. 205, pp. 119–131, 1994.
- [76] D. Arndt, W. Bangerth, M. Feder, M. Fehling, R. Gassmüller, T. Heister, L. Heltai, M. Kronbichler, M. Maier, P. Munch, J.-P. Pelteret, S. Sticko, B. Turcksin, and D. Wells, “The deal.II library, version 9.4,” Journal of Numerical Mathematics, vol. 30, no. 3, pp. 231–246, 2022.
- [77] M. S. Cramer and R. Sen, “Exact solutions for sonic shocks in van der waals gases,” The Physics of fluids, vol. 30, no. 2, pp. 377–385, 1987.
- [78] G. Lai, “Interactions of composite waves of the two-dimensional full euler equations for van der waals gases,” SIAM Journal on Mathematical Analysis, vol. 50, no. 4, pp. 3535–3597, 2018.

- [79] M. Fossati and L. Quartapelle, “The riemann problem for hyperbolic equations under a nonconvex flux with two inflection points,” 2014.
- [80] P. A. Thompson and K. C. Lambrakis, “Negative shock waves,” Journal of Fluid Mechanics, vol. 60, no. 1, p. 187–208, 1973.
- [81] S. P. Lyon, “SESAME: The Los Alamos National Laboratory equation of state database,” Los Alamos National Laboratory report LA-UR-92-3407, 1992.
- [82] K. Thompson, “EOSPAC6.” <https://github.com/KineticTheory/eospac6>. Accessed: April 25, 2023.
- [83] D. A. Pimentel, “EOSPAC user’s manual: Version 6.4.” <https://www.osti.gov/biblio/1489917>, 2019.
- [84] P. Woodward and P. Colella, “The numerical simulation of two-dimensional fluid flow with strong shocks,” Journal of computational physics, vol. 54, no. 1, pp. 115–173, 1984.
- [85] P. H. Schardin, “High frequency cinematography in the shock tube,” The Journal of Photographic Science, vol. 5, no. 2, pp. 17–19, 1957.
- [86] S.-M. Chang and K.-S. Chang, “On the shock–vortex interaction in schardin’s problem,” Shock Waves, vol. 10, no. 5, pp. 333–343, 2000.
- [87] J.-F. Haas and B. Sturtevant, “Interaction of weak shock waves with cylindrical and spherical gas inhomogeneities,” Journal of Fluid Mechanics, vol. 181, pp. 41–76, 1987.
- [88] J. J. Quirk and S. Karni, “On the dynamics of a shock-bubble interaction,” Journal of Fluid Mechanics, vol. 318, pp. 129–163, 1996.
- [89] Y. Wang and J. Li, “Stiffened gas approximation and grp resolution for fluid flows of real materials,” arXiv preprint, arXiv:2108.13780, 2021.
- [90] T. V. Bazhenova, L. G. Gvozdeva, and M. A. Nettleton, “Unsteady interactions of shock waves,” Aerospace ScL, vol. 21, pp. 249–331, 1984.
- [91] K. Takayama and O. Inoue, “Shock wave diffraction over a 90 degree sharp corner posters presented at 18th issw,” Shock Waves, vol. 1, pp. 301–312, 1991.
- [92] B. P. Brown and B. M. Argrow, “Nonclassical dense gas flows for simple geometries,”

- AIAA Journal, vol. 36, pp. 1842–1847, 1998.
- [93] P. A. Thompson, “A fundamental derivative in gasdynamics,” Physics of Fluids, vol. 14, pp. 1843–1849, 1971.
- [94] P. Colonna, N. Nannan, A. Guardone, and T. Van der Stelt, “On the computation of the fundamental derivative of gas dynamics using equations of state,” Fluid phase equilibria, vol. 286, no. 1, pp. 43–54, 2009.

APPENDIX A

Derivation of Exact Solutions for the Euler Equations

We detail the derivation of several exact solutions for the Euler equations.

A.1 Isentropic Vortex*

Theorem A.1.1. *Consider the van der Waals equation of state (2.15) with $a > 0$, $b := 0$, and $\gamma := \frac{3}{2}$ or $\gamma := 2$. Let $\mathbf{x}^0 \in \mathbb{R}^2$, $\beta > 0$, $r_0 > 0$. Let $\rho_\infty > 0$, $\mathbf{v}_\infty \in \mathbb{R}^2$, $p_\infty > 0$ and assume that*

$$p_\infty > \frac{1}{3}a\rho_\infty^2, \quad a\rho_\infty + \frac{3p_\infty}{\rho_\infty} > \frac{\beta^2 e^1}{8r_0^2\pi^2} \quad (\text{A.1})$$

if $\gamma = \frac{3}{2}$. The following density, velocity, and pressure fields solve the compressible Euler equations (1.14a)–(1.14c) with the van der Waals equation of state:

$$\rho(\mathbf{x}, t) := \begin{cases} \left(\frac{3C}{4a} - \frac{1}{2}\sqrt{\frac{9C^2}{4a^2} + \frac{2}{a}\left(F + \frac{1}{2r_0^2}\psi(\bar{\mathbf{x}})^2\right)} \right)^2, & \text{if } \gamma = \frac{3}{2}, \\ \rho_\infty - \frac{\rho_\infty^2}{4p_\infty r_0^2}\psi(\bar{\mathbf{x}}), & \text{if } \gamma = 2, \end{cases} \quad (\text{A.2a})$$

$$\mathbf{v}(\mathbf{x}, t) := \mathbf{v}_\infty + \psi(\bar{\mathbf{x}})(-\bar{x}_2, \bar{x}_1)^\top, \quad (\text{A.2b})$$

$$p(\mathbf{x}, t) := C\rho(\mathbf{x}, t)^\gamma - a\rho(\mathbf{x}, t)^2, \quad (\text{A.2c})$$

with $\psi(\bar{\mathbf{x}}) := \frac{\beta}{2\pi} \exp(\frac{1}{2}(1 - \frac{1}{r_0^2}\|\bar{\mathbf{x}}\|_{\ell_2}^2))$, $(\bar{x}_1, \bar{x}_2) = \bar{\mathbf{x}} := \mathbf{x} - \mathbf{x}^0 - \mathbf{v}_\infty t$, $C = (p_\infty + a\rho_\infty^2)/\rho_\infty^{3/2}$, and $F = -a\rho_\infty - 3p_\infty/\rho_\infty$.

Proof. The derivation of the isentropic vortex begins with the additional assumption that the velocity field is divergence free. That is, $\nabla \cdot \mathbf{v} = 0$. Under this assumption the Euler equations take the following simplified form:

This theorem and proof are reprinted with permission from [1].

$$\partial_t \rho(\mathbf{x}, t) + \mathbf{v}(\mathbf{x}, t) \cdot \nabla \rho(\mathbf{x}, t) = 0, \quad \mathbf{x} \in \mathbb{R}^2, t > 0, \quad (\text{A.3})$$

$$\partial_t \mathbf{v}(\mathbf{x}, t) + (\mathbf{v}(\mathbf{x}, t) \cdot \nabla) \mathbf{v}(\mathbf{x}, t) = -\frac{1}{\rho(\mathbf{x}, t)} \nabla p(\mathbf{x}, t), \quad \mathbf{x} \in \mathbb{R}^2, t > 0 \quad (\text{A.4})$$

$$\partial_t e(\mathbf{x}, t) + \mathbf{v}(\mathbf{x}, t) \cdot \nabla e(\mathbf{x}, t) = 0, \quad \mathbf{x} \in \mathbb{R}^2, t > 0, \quad (\text{A.5})$$

with $\mathbf{x} := (x_1, x_2)$, boundary conditions, $(\rho_\infty, \mathbf{v}_\infty := (v_{1,\infty}, v_{2,\infty})^\top, p_\infty)$ and yet to be determined initial conditions $(\rho_0(\mathbf{x}), \mathbf{v}_0(\mathbf{x}), p_0(\mathbf{x}))$. To keep things general, we make no assumption on the equation of state for $p = p(\rho, e)$.

We write the solution as a perturbation of the far-field state; i.e. we define $\mathbf{v} := \mathbf{v}_\infty + \delta \mathbf{v}$ with

$$\delta \mathbf{v}(\mathbf{x}, t) := \begin{pmatrix} \partial_{x_2} \psi(\mathbf{x} - \mathbf{x}^0 - \mathbf{v}_\infty t) \\ -\partial_{x_1} \psi(\mathbf{x} - \mathbf{x}^0 - \mathbf{v}_\infty t) \end{pmatrix}, \quad (\text{A.6})$$

with the stream function $\psi(\mathbf{x}) := \frac{\beta}{2\pi} \exp(\frac{1}{2}(1 - \frac{\|\mathbf{x}\|_{\ell^2}^2}{r_0^2}))$. Here $\mathbf{x}^0 := (x_1^0, x_2^0) \in \mathbb{R}^2$, β , and r_0 are free parameters. To further simplify notation, define $(\bar{x}_1, \bar{x}_2) = \bar{\mathbf{x}} := \mathbf{x} - \mathbf{x}^0 - \mathbf{v}_\infty t$ and $r^2 := \|\bar{\mathbf{x}}\|_{\ell^2}^2$. Note the following identities which will be used later on:

$$\partial_{x_i} \psi(\bar{\mathbf{x}}) = -\frac{\bar{x}_i}{r_0^2} \psi(\bar{\mathbf{x}}), \quad (\text{A.7})$$

$$\partial_{x_i x_j} \psi(\bar{\mathbf{x}}) = \frac{1}{r_0^2} \left(-\delta_{ij} + \frac{\bar{x}_i \bar{x}_j}{r_0^2} \right) \psi(\bar{\mathbf{x}}), \quad (\text{A.8})$$

$$\partial_{tx_i} \psi(\bar{\mathbf{x}}) = \frac{1}{r_0^2} \left(v_{i,\infty} - \frac{\bar{x}_i \mathbf{v}_\infty \cdot \bar{\mathbf{x}}}{r_0^2} \right) \psi(\bar{\mathbf{x}}), \quad (\text{A.9})$$

where δ_{ij} is the Kronecker symbol and $i, j \in \{1, 2\}$.

Using that $\mathbf{v} = \mathbf{v}_\infty + \delta \mathbf{v}$, the left hand side of (A.4) becomes,

$$\partial_t \mathbf{v} + \mathbf{v} \cdot \nabla \mathbf{v} = \partial_t (\delta \mathbf{v}) + (\mathbf{v}_\infty \cdot \nabla) \delta \mathbf{v} + (\delta \mathbf{v} \cdot \nabla) \delta \mathbf{v}. \quad (\text{A.10})$$

From the definition of $\delta \mathbf{v}$ and the identities (A.7), (A.8) and (A.9), we have,

$$(\delta \mathbf{v} \cdot \nabla) \delta \mathbf{v} = \begin{bmatrix} (\partial_{x_2} \psi)(\partial_{x_1 x_2}^2 \psi) - (\partial_{x_1} \psi)(\partial_{x_2}^2 \psi) \\ -(\partial_{x_2} \psi)(\partial_{x_1}^2 \psi) + (\partial_{x_1} \psi)(\partial_{x_1 x_2}^2 \psi) \end{bmatrix} = -\frac{\bar{\mathbf{x}}}{r_0^4} \psi(\bar{\mathbf{x}})^2 \quad (\text{A.11})$$

$$(\mathbf{v}_\infty \cdot \nabla) \delta \mathbf{v} = \frac{1}{r_0^2} \left(- \begin{bmatrix} v_{2,\infty} \\ -v_{1,\infty} \end{bmatrix} + \frac{\mathbf{v}_\infty \cdot \bar{\mathbf{x}}}{r_0^2} \begin{bmatrix} \bar{x}_2 \\ -\bar{x}_1 \end{bmatrix} \right) \psi(\bar{\mathbf{x}}) \quad (\text{A.12})$$

$$\partial_t(\delta \mathbf{v}) = \frac{1}{r_0^2} \left(\begin{bmatrix} v_{2,\infty} \\ -v_{1,\infty} \end{bmatrix} - \frac{\mathbf{v}_\infty \cdot \bar{\mathbf{x}}}{r_0^2} \begin{bmatrix} \bar{x}_2 \\ -\bar{x}_1 \end{bmatrix} \right) \psi(\bar{\mathbf{x}}) \quad (\text{A.13})$$

Thus equation (A.4) becomes $-\frac{\bar{\mathbf{x}}}{r_0^4} \psi(\bar{\mathbf{x}})^2 = -\frac{1}{\rho} \nabla p$. This identity is furthermore written as,

$$-\frac{1}{2r_0^2} \nabla(\psi(\bar{\mathbf{x}})^2) = \frac{1}{\rho(t, \mathbf{x})} \nabla p(\rho(t, \mathbf{x})). \quad (\text{A.14})$$

Up to this point, we have not made any assumption on the equation of state. We recover the well known isentropic vortex solution if we assume the pressure is given by the ideal gas law; i.e. $p(\rho) = C\rho^\gamma$ for the isentropic flow where $C = p_\infty/\rho_\infty^\gamma$. We now proceed with the van der Waals equation of state. For isentropic flows we have

$$p(\rho) = \frac{C\rho^\gamma}{(1-b\rho)^\gamma} - a\rho^2, \quad (\text{A.15})$$

where C is some constant. (Note, we work with an arbitrary b to keep things general in the beginning.) Following the same process as in the ideal gas case, we compute the indefinite integral, $\int \frac{1}{\rho} \partial_{x_i} p(\rho) dx_i$:

$$\begin{aligned}
-\frac{1}{2r_0^2} \int \partial_{x_i} \psi(\bar{\mathbf{x}})^2 dx_i &= \int \frac{1}{\rho} \partial_{x_i} p(\rho) dx_i = \frac{p(\rho)}{\rho} + \int \frac{p(\rho)}{\rho^2} \partial_{x_i} \rho dx_i \\
&= \frac{p(\rho)}{\rho} + \int \left(\frac{C\rho^{\gamma-2}}{(1-b\rho)^\gamma} - a \right) \rho_{x_i} dx_i = \frac{p(\rho)}{\rho} + \int \frac{\partial}{\partial x_i} \left[\frac{C}{\gamma-1} \left(\frac{\rho}{1-b\rho} \right)^{\gamma-1} - a\rho \right] dx_i \\
&= \frac{C\rho^{\gamma-1}(\gamma-b\rho)}{(\gamma-1)(1-b\rho)^\gamma} - 2a\rho + F.
\end{aligned}$$

Hence, $\rho(\bar{\mathbf{x}})$ can be found by solving the equation,

$$-\frac{1}{2r_0^2} \psi(\bar{\mathbf{x}})^2 = \frac{C\rho^{\gamma-1}(\gamma-b\rho)}{(\gamma-1)(1-b\rho)^\gamma} - 2a\rho + F. \quad (\text{A.16})$$

We have two immediate cases for solutions that can be found explicitly.

Case 1: $\gamma = 3/2$ and $b = 0$: In this case, (A.16) becomes a quadratic equation for $\sqrt{\rho}$,

$$\rho - \frac{3C}{2a} \sqrt{\rho} - \frac{1}{2a} \left(F + \frac{1}{2r_0^2} \psi(\bar{\mathbf{x}})^2 \right) = 0. \quad (\text{A.17})$$

The constants C and F are determined by applying the far field condition to (A.15) and (A.17):

$$C = \frac{p_\infty + a\rho_\infty^2}{\rho_\infty^{3/2}} \quad \text{and} \quad F = -a\rho_\infty - \frac{3p_\infty}{\rho_\infty}. \quad (\text{A.18})$$

However, care must be taken in the choice of p_∞ and ρ_∞ so that the sound speed remains real. Recall that the sound speed for the van der Waals EOS is,

$$c(\rho, p) = \sqrt{\gamma \frac{p + a\rho^2}{\rho(1-b\rho)} - 2a\rho}. \quad (\text{A.19})$$

The hypothesis $p_\infty > \frac{1}{3}a\rho_\infty^2$ guarantees that $c(\rho_\infty, p_\infty)^2 > 0$.

The physical root for equation (A.17) is,

$$\sqrt{\rho} = \frac{3C}{4a} - \frac{1}{2} \sqrt{\frac{9C^2}{4a^2} + \frac{2}{a} \left(F + \frac{1}{2r_0^2} \psi(\bar{\mathbf{x}})^2 \right)}. \quad (\text{A.20})$$

Furthermore, for the root to be real we require that $-F > \frac{1}{2r_0^2}\psi(\bar{\mathbf{x}})^2$ for all $\bar{\mathbf{x}} \in \mathbb{R}^2$. In particular,

$$a\rho_\infty + \frac{3p_\infty}{\rho_\infty} > \frac{\beta^2 e^1}{8r_0^2\pi^2}. \quad (\text{A.21})$$

Lastly, we must justify that the system remains hyperbolic; that is, the sound speed is real for all $(\mathbf{x}, t) \in \mathbb{R}^2 \times [0, \infty)$. Since the flow is isentropic, the sound speed for the van der Waals EOS (with $\gamma = 3/2$ and $b = 0$) is, $f(\rho) := c(p(\rho), \rho)^2 = \frac{3}{2}C\sqrt{\rho} - 2a\rho$. Note that $\lim_{\rho \rightarrow 0^+} f(\rho) = 0$, $f'(\rho) = \frac{3C}{4\sqrt{\rho}} - 2a$, and $\lim_{\rho \rightarrow 0^+} f'(\rho) = \infty$. Therefore, $f(\rho)$ has a maximum at $\rho = \left(\frac{3C}{8a}\right)^2$ and hence $f(\rho) > 0$ for $\rho \in (0, \left(\frac{3C}{8a}\right)^2)$. From the definition of ρ , (A.2a), we see that $0 < \rho < \left(\frac{3C}{4a}\right)^2$. Thus the sound speed is always real.

Case 2: $\gamma = 2$ and $b = 0$: For these choices of parameters, (A.16) becomes,

$$2(C - a)\rho + F + \frac{1}{2r_0^2}\psi(\bar{\mathbf{x}}) = 0. \quad (\text{A.22})$$

Using the far field boundary conditions for (A.2c) and (A.22) we find that $C = \frac{p_\infty}{\rho_\infty^2} + a$ and $F = -2p_\infty/\rho_\infty$, respectively. Solving for ρ in (A.22) we have,

$$\rho = \rho_\infty - \frac{\rho_\infty^2}{4p_\infty r_0^2}\psi(\bar{\mathbf{x}}). \quad (\text{A.23})$$

Note the sound speed is $c(p(\rho), \rho)^2 = 2(C - a)\rho = \frac{2p_\infty}{\rho_\infty}\rho > 0$. □

APPENDIX B

Numerical Algorithms

B.1 Upper Bound on Maximum Wave Speed

In this section we provide the algorithm which computes an upper bound to the maximum wave speed. The details are discussed in Section 4.6.

Algorithm 1 Computing $\lambda_{\max}(\hat{p}^*)$

Require: $\mathbf{u}_L, \mathbf{u}_R, \mathbf{n}_{LR}, p_L, p_R$

compute $\varphi(p_{\min}), \varphi(p_{\max})$ from (4.59)

if $0 \leq \varphi(p_{\min})$ **then**

compute \tilde{p}^* from (4.78)

$\hat{p}^* = \max(p_{\min}, \tilde{p}^*)$

▷ One may also set $\hat{p}^* = 0$

else if $0 \leq \varphi(p_{\max})$ **then**

if $\gamma_{\min} = \gamma_m$ **then**

compute \tilde{p}_1^* and \tilde{p}_2^* from (4.81a) and (4.81b) resp.

else

compute \tilde{p}_1^* and \tilde{p}_2^* from (4.83a) and (4.83b) resp.

end if

$\hat{p}^* = \min\{p_{\max}, \tilde{p}_1^*, \tilde{p}_2^*\}$

else

compute \tilde{p}_1^* and \tilde{p}_2^* from (4.85) and (4.89) resp.

$\hat{p}^* = \min\{\tilde{p}_1^*, \tilde{p}_2^*\}$

end if

return $\lambda_{\max}(\hat{p}^*) := \max\{-\lambda_L^-(\hat{p}^*), \lambda_R^+(\hat{p}^*)\}$

B.2 The Quadratic Newton Method for Limiting the Surrogate Entropy

The algorithm for computing the limiting Φ_i^s is given Algorithm 2.

Algorithm 2 Limiting Φ_i^s with the Quadratic Newton Method

Require: $\mathbf{U}_i^{\text{L},n+1}$, \mathbf{P}_{ij} , $\ell_j^{i,\rho}$, ϵ

- 1: **set** $\ell_{\text{L}}^0 = 0$ and $\ell_{\text{R}}^0 = \ell_j^{i,\rho}$
 - 2: **compute** $f(\ell_Z^0) = \Phi_i^s(\mathbf{U}_i^{\text{L},n+1} + \ell_Z^0 \mathbf{P}_{ij})$ for $Z = \text{L}$ and $Z = \text{R}$.
 - 3: **compute** $f'(\ell_Z^0) = D\Phi_i^s(\mathbf{U}_i^{\text{L},n+1} + \ell_Z^0 \mathbf{P}_{ij}) \cdot \mathbf{P}_{ij}$ for $Z = \text{L}$ and $Z = \text{R}$.
 - 4: **compute** $f[\ell_{\text{L}}^0, \ell_{\text{L}}^0, \ell_{\text{R}}^0]$ and $f[\ell_{\text{L}}^0, \ell_{\text{R}}^0, \ell_{\text{R}}^0]$.
 - 5: **compute** $\ell^{\text{L}}(\ell_{\text{L}}^0, \ell_{\text{R}}^0)$ and $\ell^{\text{R}}(\ell_{\text{L}}^0, \ell_{\text{R}}^0)$ from (6.57a) and (6.57b), resp.
 - 6: **set** $\ell_{\text{L}} = \min\{\ell^{\text{L}}(\ell_{\text{L}}^0, \ell_{\text{R}}^0), \ell^{\text{R}}(\ell_{\text{L}}^0, \ell_{\text{R}}^0)\}$ and $\ell_{\text{R}} = \max\{\ell^{\text{L}}(\ell_{\text{L}}^0, \ell_{\text{R}}^0), \ell^{\text{R}}(\ell_{\text{L}}^0, \ell_{\text{R}}^0)\}$
 - 7: **while** $|\ell_{\text{R}} - \ell_{\text{L}}|/|\ell_{\text{L}}| > \epsilon$ **do**
 - 8: **compute** $f(\ell_Z) = \Phi_i^s(\mathbf{U}_i^{\text{L},n+1} + \ell_Z \mathbf{P}_{ij})$ for $Z = \text{L}$ and $Z = \text{R}$.
 - 9: **compute** $f'(\ell_Z) = D\Phi_i^s(\mathbf{U}_i^{\text{L},n+1} + \ell_Z \mathbf{P}_{ij}) \cdot \mathbf{P}_{ij}$ for $Z = \text{L}$ and $Z = \text{R}$.
 - 10: **compute** $f[\ell_{\text{L}}, \ell_{\text{L}}, \ell_{\text{R}}]$ and $f[\ell_{\text{L}}, \ell_{\text{R}}, \ell_{\text{R}}]$.
 - 11: **compute** $\ell^{\text{L}}(\ell_{\text{L}}, \ell_{\text{R}})$ and $\ell^{\text{R}}(\ell_{\text{L}}, \ell_{\text{R}})$ from (6.57a) and (6.57b), resp.
 - 12: **set** $\ell_{\text{L}} = \min\{\ell^{\text{L}}(\ell_{\text{L}}, \ell_{\text{R}}), \ell^{\text{R}}(\ell_{\text{L}}, \ell_{\text{R}})\}$ and $\ell_{\text{R}} = \max\{\ell^{\text{L}}(\ell_{\text{L}}, \ell_{\text{R}}), \ell^{\text{R}}(\ell_{\text{L}}, \ell_{\text{R}})\}$
 - 13: **end while**
 - 14: **set** $\ell_j^{i,s} := \min\{\ell_{\text{L}}, \ell_{\text{R}}\}$
-

B.3 Approximate Godunov-type Solver

There are several methods for handling outflow boundary conditions in the literature. One way to do this is to exactly solve the Riemann problem at $x = 0$, which requires another state on the interface. This method is referred to as the Godunov method, see Godunov [6] which is used in a vast amount of numerical methods for conservation laws. Due to the costly nature of solving the Riemann problem, we instead propose an approximate solution at $x = 0$.

B.3.1 The Method

In order to solve the extended Riemann problem for the Euler equations we must compute ρ^* by solving $\varphi(\rho^*) = 0$, see (4.59). The velocity, v^* is then computed by $v^* = v_L - f_L(\rho^*) = v_R + f_R(\rho^*)$. We instead propose to use $\widehat{\rho}^* \geq \rho^*$ computed using Algorithm 1. Then, we define the velocity in the star domain by $\widehat{v}_L^* = v_L - f_L(\widehat{\rho}^*)$ and $\widehat{v}_R^* = v_R + f_R(\widehat{\rho}^*)$, and assuming $\widehat{\rho}^* \neq \rho^*$ then $v_L^* \neq v_R^*$. For a state across a shock wave, we compute the density, $\widehat{\rho}_{*L}$ or $\widehat{\rho}_{*R}$ from (4.37) or (4.43); respectively, using $\widehat{\rho}^*$ instead of ρ^* . In this case, we have that $\widehat{\rho}_{*Z} \geq \rho_{*Z}$ for $Z \in \{L, R\}$.

Lemma B.3.1. *If $\widehat{\rho}^* \geq \rho^*$, then $\widehat{v}_L^* \leq \widehat{v}_R^*$ where $\widehat{v}_L^* := v_L - f_L(\widehat{\rho}^*)$ and $\widehat{v}_R^* := v_R + f_R(\widehat{\rho}^*)$.*

Proof. The proof follows immediately from the fact that $\varphi(p)$ is monotonically increasing (see Lemma 4.4.1); that is,

$$\varphi(\widehat{\rho}^*) = f_R(\widehat{\rho}^*) + f_L(\widehat{\rho}^*) + v_R - v_L \geq \varphi(\rho^*) = 0. \quad (\text{B.1})$$

Hence $\widehat{v}_R^* \geq \widehat{v}_L^*$. □

B.3.2 The Algorithm

The choice of solution at $x = 0$ is best described through diagrams as the exact solution can be seen in the algorithm. See Figure B.1 for several different wave profiles that deter-

mine different solutions at $x = 0$ and Algorithm 3 for the full context. Since we need the approximate solution at $x = 0$, we must look at the wave structure.

Algorithm 3 Approximate Riemann Solution

Require: $\mathbf{u}_L, \mathbf{u}_R, \mathbf{n}_{LR}, p_L, p_R$

- 1: **compute** \widehat{p}^* from Algorithm 1
 - 2: **compute** $\varphi(p_L)$ and $\varphi(p_R)$
 - 3: **if** $\varphi(p_L) \geq 0$ and $\varphi(p_R) \geq 0$ **then**
 - 4: **compute** $\lambda_L^- := v_L - a_L$ and $\lambda_R^+ := v_R + a_R$
 - 5: **else if** $\varphi(p_L) \geq 0$ and $\varphi(p_R) < 0$ **then**
 - 6: **compute** $\lambda_L^- := v_L - a_L$ and $\lambda_R^+ := \mathcal{S}_R(\widehat{p}^*)$ from (4.41)
 - 7: **else if** $\varphi(p_L) < 0$ and $\varphi(p_R) \geq 0$ **then**
 - 8: **compute** $\lambda_L^- := \mathcal{S}_L(\widehat{p}^*)$ and $\lambda_R^+ := v_R + a_R$ from (4.38)
 - 9: **else**
 - 10: **compute** $\lambda_L^- := \mathcal{S}_L(\widehat{p}^*)$ and $\lambda_R^+ := \mathcal{S}_R(\widehat{p}^*)$ from (4.38) and (4.41)
 - 11: **end if** ▷ Continued on next page
-

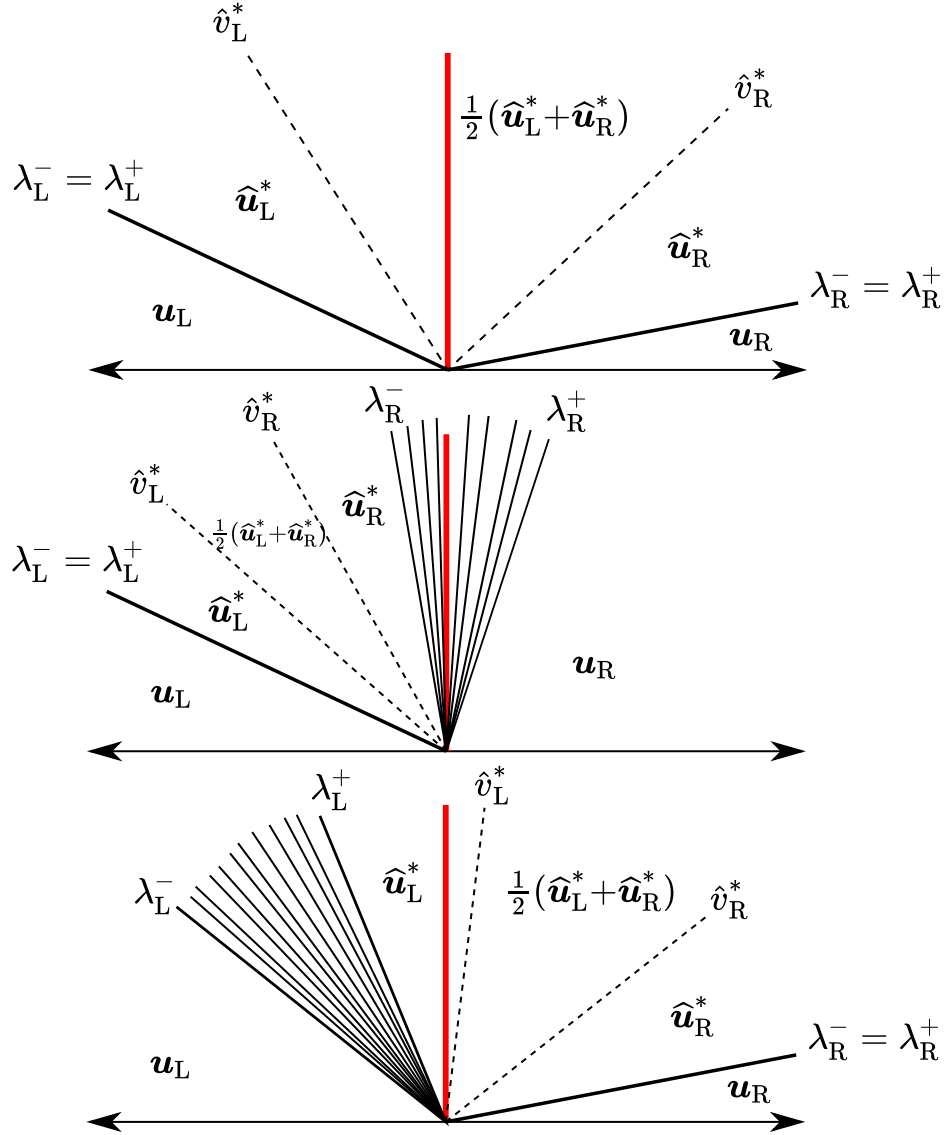


Figure B.1: Figures describing different solutions of $\mathbf{u}(0, t)$. (Top): the solution in this case is $\mathbf{u}(0, t) = \frac{1}{2}(\hat{\mathbf{u}}_L^* + \hat{\mathbf{u}}_R^*)$. (Middle): $\mathbf{u}(0, t)$ is the solution on the expansion wave which connects the two states, \mathbf{u}_R to $\hat{\mathbf{u}}_R^*$ at $x/t = 0$. (Bottom): $\mathbf{u}(0, t) = \hat{\mathbf{u}}_L^*$.

```

12: if  $\lambda_L^- \geq 0$  then
13:   set  $\mathbf{u} := \mathbf{u}_L$ 
14: else if  $\lambda_R^+ \leq 0$  then
15:   set  $\mathbf{u} := \mathbf{u}_R$ 
16: else
17:   compute  $\widehat{v}_L^*$  and  $\widehat{v}_R^*$  from Lemma B.3.1
18:   if  $\widehat{v}_L^* > 0$  then
19:     compute  $\widehat{\rho}_L^*$  from (4.37)
20:     if  $\varphi(p_L) < 0$  then
21:       set  $\mathbf{u} := (\widehat{\rho}_L^*, \widehat{v}_L^*, \mathbf{v}_L^\perp, \widehat{p}^*)$ 
22:     else
23:       compute  $\widehat{a}_L^*$ 
24:       if  $\widehat{v}_L^* - \widehat{a}_L^* > 0$  then
25:         set Riemann fan solution
26:       else
27:         set  $\mathbf{u} := (\widehat{\rho}_L^*, \widehat{v}_L^*, \mathbf{v}_L^\perp, \widehat{p}^*)$ 
28:       end if
29:     end if
30:   else if  $\widehat{v}_R^* < 0$  then
31:     compute  $\widehat{\rho}_R^*$  from (4.43)
32:     if  $\varphi(p_R) < 0$  then
33:       set  $\mathbf{u} := (\widehat{\rho}_R^*, \widehat{v}_R^*, \mathbf{v}_R^\perp, \widehat{p}^*)$ 
34:     else
35:       compute  $\widehat{a}_R^*$ 
36:       if  $\widehat{v}_R^* + \widehat{a}_R^* < 0$  then
37:         set Riemann fan solution
38:       else
39:         set  $\mathbf{u} := (\widehat{\rho}_R^*, \widehat{v}_R^*, \mathbf{v}_R^\perp, \widehat{p}^*)$ 
40:       end if
41:     end if
42:   else
43:     if  $\widehat{v}_L^* + \widehat{v}_R^* < 0$  then
44:       set  $\mathbf{u} := (\frac{1}{2}(\widehat{\rho}_L^* + \widehat{\rho}_R^*), \frac{1}{2}(\widehat{v}_L^* + \widehat{v}_R^*), \mathbf{v}_L^\perp, \widehat{p}^*)$ 
45:     else
46:       set  $\mathbf{u} := (\frac{1}{2}(\widehat{\rho}_L^* + \widehat{\rho}_R^*), \frac{1}{2}(\widehat{v}_L^* + \widehat{v}_R^*), \mathbf{v}_R^\perp, \widehat{p}^*)$ 
47:     end if
48:   end if
49: end if
50: return  $\mathbf{u}$ 

```
