# INVARIANT DOMAIN PRESERVING APPROXIMATIONS FOR THE EULER EQUATIONS WITH TABULATED EQUATION OF STATE

Bennett Clayton, Jean-Luc Guermond, Bojan Popov

# INVARIANT DOMAIN PRESERVING APPROXIMATIONS FOR THE EULER EQUATIONS WITH TABULATED EQUATION OF STATE *

BENNETT CLAYTON†, JEAN-LUC GUERMOND†, AND BOJAN POPOV†

**Abstract.** This paper is concerned with the approximation of the compressible Euler equations supplemented with an equation of state that is either tabulated or is given by an expression that is so involved that solving elementary Riemann problems is hopeless. A robust first-order approximation technique that guarantees that the density and the internal energy are positive is proposed. A key ingredient of the method is a local approximation of the equation of state using a co-volume ansatz from which upper bounds on the maximum wave speed are derived for every elementary Riemann problem.

**Key words.** Compressible Euler equations, tabulated equation of state, maximum wave speed, Riemann problem, Invariant domain preserving approximation, composite waves

**AMS subject classifications.** 65M60, 65M12, 65M22, 35L65

**1. Introduction.** In many important applications, the compressible Euler equations are supplemented with an equation of state that is either tabulated or given by a complicated analytic expression. Throughout the paper, we will refer to this type of equation of state as the 'oracle'. In this case, approximating the Euler equations while guaranteeing positivity of the density and positivity of the internal energy is problematic since no exact solution of elementary Riemann problems can be a priori inferred. Solving a Riemann problem when the equation of state is analytically well defined is feasible, though possibly expensive, (see e.g., Colella and Glaz [5, §1], Ivings et al. [16], Quartapelle et al. [24]). This cannot be efficiently done with an oracle for this requires interpolating/approximating the equation of state, and to the best of our knowledge, there is no clear technique to do so in the literature. Various methods to avoid this problem have been proposed in the literature. For instance, one can use approximate Riemann solvers like in Dukowicz [7], [5, §2], Roe and Pike [26], Pike [23], or simplify the Riemann problem by using flux splitting techniques like in Toro et al. [28]. However, for most of these techniques very little is guaranteed besides positivity of the density, which is not difficult to achieve. The objective of the paper is to address these questions. More precisely, we propose an approximation method to solve the Euler equations equipped with an oracle. This is done by adapting the technique from Guermond and Popov [12] where invariant-domain properties are obtained by ascertaining that they hold true for elementary Riemann problems. The key is to augment each elementary Riemann system with an additional scalar equation and replace the oracle by a covolume equation of state where the coefficient $\gamma$ is variable and obtained as the solution to the additional equation. This idea is adapted from Abgrall and Karni [1]. A variation of this idea is also employed in [5, Eq. (37)] and Pantano et al. [22, Eq. (22)]. The proposed algorithm is explicit in time and preserves the positivity of the density and the internal energy under an appropriate CFL restriction on the time step. Additional properties can be preserved depending on the nature of the oracle. As in Guermond et al. [14], the method is agnostic to the space approximation. An interesting feature of the method is that it automatically recovers the standard co-volume behaviour if the oracle is indeed a covolume equation of state. In

†Department of Mathematics, Texas A&M University 3368 TAMU, College Station, TX 77843, USA.

1

compliance with Godunov's theorem, the method is first-order accurate in space, however, achieving higher-order accuracy in space is out of the scope of this paper. This can be done by implementing the convex limiting technique described in [13, 14]. This work is in progress and will be reported elsewhere.

The paper is organised as follows. The problem and the notation are introduced in §2. The space and time approximation method from [12] is also briefly recalled in this section. The main motivation of the paper is given at the end of §2.2. We introduce an extended Riemann problem in §3. The key point of this section is summarized in Remark 3.1. An exact weak solution to the extended Riemann problem is constructed in §4. It is also shown in this section that this weak solution satisfies the expected invariant domain-properties. The main results of §4 are Lemma 4.4, Lemma 4.5 and Theorem 4.6. An upper bound on the maximum wave speed for the extended Riemann problem is derived in §5. This upper bound is the key piece of information that is needed for practitioners who may have little interest in the Riemann problem theory (see §5.2–§5.5). The fact that this estimate of the maximum wave speed is a guaranteed upper bound implies that the proposed numerical algorithm satisfies the invariant-domain properties stated in Theorem 4.6. The technique introduced in the paper is illustrated in §6 with continuous finite elements and various equations of states. Finally, the paper is supplemented with an appendix collecting technical results. Various pieces of software are made publicly available to guarantee reproducibility (Clayton et al. [3, 4]).

**2. Formulation of the problem.** We formulate the problem and introduce notation in this section. The main motivation for the theory developed in the paper is given at the end of §2.2.

**2.1. The Euler equations.** We consider a compressible inviscid fluid occupying a bounded, polyhedral domain $D$ in $\mathbb{R}^d$. Here $d$ is the space dimension. We assume that the dynamics of the system is modeled by the compressible Euler equations equipped with an equation of state that can be either tabulated or given by a very complicated analytic expression. The dependent variable is $\boldsymbol{u} := (\rho, \boldsymbol{m}, E)^\mathsf{T} \in \mathbb{R}^{d+2}$, where $\rho$ is the density, $\boldsymbol{m}$ the momentum, $E$ the total mechanical energy. In this paper $\boldsymbol{u}$ is considered to be a column vector. The velocity is given by $\boldsymbol{v} := \rho^{-1}\boldsymbol{m}$. The quantity $e(\boldsymbol{u}) := \rho^{-1}E - \frac{1}{2}\|\boldsymbol{v}\|_{\ell^2}^2$ is the specific internal energy. To simplify the notation later on we introduce the flux $\mathbb{f}(\boldsymbol{u}) := (\boldsymbol{m}, \boldsymbol{v} \otimes \boldsymbol{m} + p(\boldsymbol{u})\mathbb{I}_d, \boldsymbol{v}(E+p))^\mathsf{T} \in \mathbb{R}^{(d+2)\times d}$, where $\mathbb{I}_d$ is the $d \times d$ identity matrix. The convention adopted in the paper is that for any vectors $\boldsymbol{a}$, $\boldsymbol{b}$, with entries $\{a_k\}_{k\in\{1:d\}}$, $\{b_k\}_{k\in\{1:d\}}$, the following holds: $(\boldsymbol{a} \otimes \boldsymbol{b})_{kl} = a_k b_l$ and $\nabla\cdot\boldsymbol{a} = \sum_{k=\in\{1:d\}} \partial_{x_k} a_k$. Moreover, for any second-order tensor g with entries $\{g_{kl}\}_{k\in\{1:d+2\}}^{l\in\{1:d\}}$, we define $(\nabla\cdot g)_k = \sum_{l\in\{1:d\}} \partial_{x_l} g_{kl}$.

Given some initial time $t_0$ and initial data $\boldsymbol{u}_0(\boldsymbol{x}) := (\rho_0, \boldsymbol{m}_0, E_0)(\boldsymbol{x})$, we look for $\boldsymbol{u}(\boldsymbol{x},t) := (\rho, \boldsymbol{m}, E)(\boldsymbol{x},t)$ solving the following system in some weak sense:

$$\partial_t \rho + \nabla\cdot(\boldsymbol{v}\rho) = 0 \qquad\qquad \text{a.e. } t > t_0,\ \boldsymbol{x} \in D, \tag{2.1a}$$

$$\partial_t \boldsymbol{m} + \nabla\cdot\big(\boldsymbol{v} \otimes \boldsymbol{m} + p(\boldsymbol{u})\mathbb{I}_d\big) = \boldsymbol{0} \qquad\qquad \text{a.e. } t > t_0,\ \boldsymbol{x} \in D, \tag{2.1b}$$

$$\partial_t E + \nabla\cdot\big(\boldsymbol{v}(E + p(\boldsymbol{u}))\big) = 0 \qquad\qquad \text{a.e. } t > t_0,\ \boldsymbol{x} \in D, \tag{2.1c}$$

where $p : \mathscr{A} \to \mathbb{R}$ is the pressure, and $\mathscr{A}$ is the admissible set:

$$\mathscr{A} := \big\{\boldsymbol{u} = (\rho, \boldsymbol{m}, E) \in \mathbb{R}^{d+2} \mid \rho > 0,\ e(\boldsymbol{u}) > 0\big\}. \tag{2.2}$$

We refer to the mapping $p : \mathscr{A} \to \mathbb{R}$ as the *oracle*. For all $\beta \geq 0$, we introduce the following convex subset of $\mathscr{A}$:

$$\mathscr{B}(\beta) := \big\{\boldsymbol{u} = (\rho, \boldsymbol{m}, E) \in \mathbb{R}^{d+2} \mid \rho > 0,\ 1 - \beta\rho > 0,\ e(\boldsymbol{u}) > 0\big\}. \tag{2.3}$$

82  We further assume in the paper that the oracle is such that there exists a number $b \geq 0$, henceforth
83  called the covolume constant, so that the following holds for all $\boldsymbol{u} \in \mathcal{B}(b)$:

84  (2.4)                                      $$p(\boldsymbol{u}) > 0.$$

85  The inverse of the covolume constant $b$ is the maximal density the fluid can reach. We take $b = 0$
86  if this constant is not a priori known.

87      Our goal in the paper is to approximate (2.1) by adapting the technique described in Guer-
88  mond and Popov [12]. As explained in the next section, this is done by constructing an artificial
89  viscosity that ensures that some relevant invariant-domain properties can be established, thereby
90  guaranteeing that the approximation technique is robust (i.e., satisfies physical bounds under a
91  reasonable CFL condition). The two key difficulties that arise in this endeavor are that it is nearly
92  impossible to construct solutions to elementary Riemann problems (or at least highly nontrivial,
93  see e.g., Quartapelle et al. [24], Fossati and Quartapelle [9]), since the equation of state is either
94  not available or too complicated. We propose a solution to this problem in §3 and §4. Taking
95  inspiration from Colella and Glaz [5], Abgrall and Karni [1], Pantano et al. [22], we introduce a
96  technique consisting of approximating the oracle by a covolume $\gamma$-law, where $\gamma$ solves an additional
97  conservation equation.

98      *Remark* 2.1 (Pressure). In practice there are many equations of state that cannot guarantee
99  (2.4) over the entire set $\mathcal{B}(b)$, but the algorithm proposed in the paper works properly as long as the
100  numerical states stay in a subset of $\mathcal{B}(b)$ where the pressure stays positive. This situation occurs
101  in many realistic applications.                                                                    □

102      **2.2. Space and time approximation.** Let us first recall the space and time approximation
103  technique described in [12]. This method is in some sense a discretization-independent extension
104  of the scheme by Lax [18, p. 163]. Without going into the details, we assume that we have at
105  hand a fully discrete scheme where time is approximated by using the forward Euler time stepping
106  and space is approximated by using some "centered" approximation of (2.1) (i.e., without any
107  artificial viscosity to stabilize the approximation). We denote by $t^n$ the current time, $n \in \mathbb{N}$, and
108  we denote by $\tau$ the current time step size; that is $t^{n+1} := t^n + \tau$. Let us assume that the current
109  approximation is a collection of states $\{\mathbf{U}_i^n\}_{i \in \mathcal{V}}$, where the index set $\mathcal{V}$ is used to enumerate all
110  the degrees of freedom of the approximation. Here $\mathbf{U}_i^n \in \mathbb{R}^{d+2}$ for all $i \in \mathcal{V}$. We assume that the
111  centered update is given by $\mathbf{U}_i^{\mathrm{G},n+1}$ with

112  (2.5)                    $$\frac{m_i}{\tau}(\mathbf{U}_i^{\mathrm{G},n+1} - \mathbf{U}_i^n) + \sum_{j \in \mathcal{I}(i)} \mathbb{f}(\mathbf{U}_j^n)\boldsymbol{c}_{ij} = \mathbf{0}.$$

113  The quantity $m_i$ is called lumped mass and we assume that $m_i > 0$ for all $i \in \mathcal{V}$. The vector
114  $\boldsymbol{c}_{ij} \in \mathbb{R}^d$ encodes the space discretization. The index set $\mathcal{I}(i)$ is called local stencil. This set
115  collects only the degrees of freedom in $\mathcal{V}$ that interact with $i$ (i.e., $j \notin \mathcal{I}(i) \Rightarrow \boldsymbol{c}_{ij} = \mathbf{0}$). We view
116  $\frac{1}{m_i} \sum_{j \in \mathcal{I}(i)} \mathbb{f}(\mathbf{U}_j^n)\boldsymbol{c}_{ij}$ as a Galerkin (or centered or inviscid) approximation of $\nabla \cdot \mathbb{f}(\boldsymbol{u})$ at time $t^n$ at
117  some grid point (or cell) $i \in \mathcal{V}$. The super-index $^{\mathrm{G}}$ is meant to remind us that (2.5) is a Galerkin
118  (or inviscid or centered) approximation of (2.1). That is, we assume that the consistency error
119  in space in (2.5) scales optimally with respect to the meshsize for the considered approximation
120  setting. We do not need to be more specific at this point. The only requirement that we make on

121 the coefficients $\boldsymbol{c}_{ij}$ is that the method is conservative; that is to say, we assume that

122 (2.6)
$$\boldsymbol{c}_{ij} = -\boldsymbol{c}_{ji} \quad \text{and} \quad \sum_{j \in \mathcal{I}(i)} \boldsymbol{c}_{ij} = \mathbf{0}.$$
123

124 An immediate consequence of this assumption is that the total mass is conserved: $\sum_{i \in \mathcal{V}} m_i \mathbf{U}_i^{\mathrm{G},n+1} =$
125 $\sum_{i \in \mathcal{V}} m_i \mathbf{U}_i^n$. Notice that for every $i \in \mathcal{V}$, the update (2.5) invokes the oracle $\mathrm{card}(\mathcal{I}(i))$ times,
126 because computing $\mathbb{f}(\mathbf{U}_j^n)$ requires computing $p(\mathbf{U}_j^n)$ for all $j \in \mathcal{I}(i)$.

127 *Remark* 2.2 (literature). The reader is referred to [12, 13] for realizations of the algorithm (2.5)
128 with continuous finite elements. Realizations of the algorithm with discontinuous elements and
129 with finite volumes are described in [14]. □

130 Of course, the approximation (2.5) is in general not appropriate if the solution to (2.1) is not
131 smooth. To recover some sort of stability (we are going to make a more precise stability statement
132 later in Theorem 4.6), we modify the scheme by adding an artificial graph viscosity based on the
133 stencil $\mathcal{I}(i)$; that is, we compute the stabilized update $\mathbf{U}_i^{n+1}$ by setting:

134 (2.7)
$$\frac{m_i}{\tau}(\mathbf{U}_i^{n+1} - \mathbf{U}_i^n) + \sum_{j \in \mathcal{I}(i)} \mathbb{f}(\mathbf{U}_j^n)\boldsymbol{c}_{ij} - \sum_{j \in \mathcal{I}(i) \setminus \{i\}} d_{ij}^n(\mathbf{U}_j^n - \mathbf{U}_i^n) = \mathbf{0}.$$

135 Here $d_{ij}^n$ is the yet to be defined artificial graph viscosity. We assume that

136
137 (2.8)
$$d_{ij}^n = d_{ji}^n > 0, \quad \text{if} \quad i \neq j.$$

138 The symmetry assumption is essential for the method to be conservative. The question addressed
139 in the paper is the following: how large should $d_{ij}^n$ be for the scheme to preserve invariant sets (and
140 possibly be entropy satisfying for some finite collection of entropies)?
141 One key observation is that one can rewrite (2.7) as follows:

142 (2.9)
$$\mathbf{U}_i^{n+1} = \left(1 - \sum_{j \in \mathcal{I}(i) \setminus \{i\}} \frac{2\tau d_{ij}^n}{m_i}\right)\mathbf{U}_i^n + \sum_{j \in \mathcal{I}(i) \setminus \{i\}} \frac{2\tau d_{ij}^n}{m_i}\overline{\mathbf{U}}_{ij}^n,$$
143

144 with the auxiliary states $\overline{\mathbf{U}}_{ij}^n$ defined as follows:

145 (2.10)
$$\overline{\mathbf{U}}_{ij}^n := \frac{1}{2}(\mathbf{U}_i^n + \mathbf{U}_j^n) - (\mathbb{f}(\mathbf{U}_j^n) - \mathbb{f}(\mathbf{U}_i^n))\boldsymbol{n}_{ij}\frac{\|\boldsymbol{c}_{ij}\|_{\ell^2}}{2d_{ij}^n}.$$
146

147 Hence, if the time step is small enough, (2.9) shows that $\mathbf{U}_i^{n+1}$ is a convex combination of the
148 following states $(\overline{\mathbf{U}}_{ij}^n)_{j \in \mathcal{I}(i)}$ (with the convention $\overline{\mathbf{U}}_{ii}^n := \mathbf{U}_i^n$). Hence if one can prove that the
149 auxiliary states $\overline{\mathbf{U}}_{ij}$ are in the set $\mathcal{B}(b)$ for all $j \in \mathcal{I}(i)$, then the update $\mathbf{U}_i^{n+1}$ is also in $\mathcal{B}(b)$, thereby
150 establishing one important invariant-domain property. (Notice in passing that it is essential here
151 to assume $d_{ij}^n \neq 0$.)
152 The main objective of the paper is to describe a technique to estimate $d_{ij}^n$ that guarantees that
153 $\overline{\mathbf{U}}_{ij}^n \in \mathcal{B}(b)$ provided both states $\mathbf{U}_i^n$ and $\mathbf{U}_j^n$ are in $\mathcal{B}(b)$. This is done by showing that $\overline{\mathbf{U}}_{ij}^n$ is a
154 space average of a solution to a Riemann problem, and by showing that this solution does satisfy
155 the invariant-domain property we are after. Then $d_{ij}^n$ is defined so that $d_{ij}^n \geq \lambda_{ij,\max}\|\boldsymbol{c}_{ij}\|_{\ell^2}$, where
156 $\lambda_{ij,\max}$ is any upper bound on the maximum wave speed in the said Riemann problem.

157    **3. The extended Riemann problem.** An important step in [12] toward proving that the
158    auxiliary state $\overline{\mathbf{U}}_{ij}^n$ defined in (2.10) is a "good" state, if $\lambda_{ij,\max}$ is an upper bound on the maximum
159    wave speed in the Riemann problem, consists of realizing that in this case $\overline{\mathbf{U}}_{ij}^n$ is a space average
160    of the exact solution to the one-dimensional Riemann problem with flux $\mathbb{f}(\boldsymbol{v})\boldsymbol{n}_{ij}$, left data $\mathbf{U}_i$, and
161    right data $\mathbf{U}_j$. The main difficulty we are facing in the present situation is that there is no analytical
162    way to estimate an upper bound $\lambda_{ij,\max}$ since the pressure is given by an oracle. We show in this
163    section how to go around this difficulty.

164    **3.1. Extension of the system and 1D reduction.** To avoid having to refer to particular
165    states $\mathbf{U}_i^n$ and $\mathbf{U}_j^n$, we now assume that we are given a left and a right admissible states, $\boldsymbol{u}_L$ and $\boldsymbol{u}_R$.
166    We also denote $\boldsymbol{n}_{ij}$ by $\boldsymbol{n}$. Instead of considering the Riemann problem where the pressure is given
167    by the oracle, we now consider an extended Riemann problem. First we make a change of basis
168    and introduce $\boldsymbol{t}_1, \ldots, \boldsymbol{t}_{d-1}$ so that $\{\boldsymbol{n}, \boldsymbol{t}_1, \ldots, \boldsymbol{t}_{d-1}\}$ forms an orthonormal basis of $\mathbb{R}^d$. With this
169    new basis we have $\boldsymbol{m} = (m, \boldsymbol{m}^\perp)^\mathsf{T}$, where $m := \rho v$, $v := \boldsymbol{v}\cdot\boldsymbol{n}$, $\boldsymbol{m}^\perp := \rho(\boldsymbol{v}\cdot\boldsymbol{t}_1, \ldots, \boldsymbol{v}\cdot\boldsymbol{t}_{d-1}) := \rho\boldsymbol{v}^\perp$.
170    Second, we augment the system by introducing a new scalar variable $\Gamma$ (and $\gamma := \frac{\Gamma}{\rho}$), the augmented
171    state $\widetilde{\boldsymbol{u}} := (\boldsymbol{u}, \Gamma)^\mathsf{T}$, and the extend the flux as follows:

172    (3.1)            $$\widetilde{\mathbb{f}}(\widetilde{\boldsymbol{u}}) := (\boldsymbol{m}, \boldsymbol{v}\otimes\boldsymbol{m} + \widetilde{p}(\widetilde{\boldsymbol{u}})\mathbb{I}_d, \boldsymbol{v}(E + \widetilde{p}(\widetilde{\boldsymbol{u}})), \boldsymbol{v}\Gamma)^\mathsf{T} = (\mathbb{f}(\widetilde{\boldsymbol{u}}), \boldsymbol{v}\Gamma)^\mathsf{T},$$

173    with the new pressure

174    (3.2)            $$\widetilde{p}(\widetilde{\boldsymbol{u}}) := \frac{(\Gamma - \rho)e(\boldsymbol{u})}{1 - b\rho} = (\gamma - 1)\frac{\rho e(\boldsymbol{u})}{1 - b\rho},$$

175    where $e(\boldsymbol{u}) := \frac{1}{\rho}(E - \frac{\|\boldsymbol{m}\|_{\ell^2}^2}{2\rho})$. Here $b$ is either given to us because this parameter can be measured,
176    or $b$ is set to be zero if one does not have any a priori knowledge on the nature of the fluid. Notice
177    that $\Gamma$ is the last component of the extended variable $\widetilde{\boldsymbol{u}}$; neither $\Gamma$ nor $\gamma = \rho^{-1}\Gamma$ are assumed to be
178    constant. The extended Riemann problem consists of seeking $\widetilde{\boldsymbol{u}} := (\boldsymbol{u}, \Gamma)^\mathsf{T} = (\rho, \boldsymbol{m}, E, \Gamma)^\mathsf{T}$ so that

179    (3.3)            $$\partial_t\widetilde{\boldsymbol{u}} + \partial_x(\widetilde{\mathbb{f}}(\widetilde{\boldsymbol{u}})\boldsymbol{n}) = \boldsymbol{0}, \quad \widetilde{\boldsymbol{u}} = \begin{pmatrix} \rho \\ m \\ \boldsymbol{m}^\perp \\ E \\ \Gamma \end{pmatrix}, \qquad \widetilde{\mathbb{f}}(\widetilde{\boldsymbol{u}})\boldsymbol{n} = \begin{pmatrix} m \\ \frac{1}{\rho}m^2 + \widetilde{p}(\widetilde{\boldsymbol{u}}) \\ v\boldsymbol{m}^\perp \\ v(E + \widetilde{p}(\widetilde{\boldsymbol{u}})) \\ v\Gamma \end{pmatrix},$$

180    with left data and right data $(\rho_Z, \boldsymbol{m}_Z\cdot\boldsymbol{n}, \boldsymbol{m}_Z^\perp, E_Z, \Gamma_Z)^\mathsf{T}$, where $Z \in \{L, R\}$, and $\Gamma_Z$ is defined so
181    that $\widetilde{p}(\widetilde{\boldsymbol{u}}_Z) = p(\boldsymbol{u}_Z) =: p_Z$, i.e., $\Gamma_Z := \rho_Z + \frac{p_Z(1 - b\rho_Z)}{e_Z}$, (notice that this means $\gamma_Z := 1 + \frac{p_Z(1 - b\rho_Z)}{\rho_Z e_Z}$).
182        As usually done in the literature, the above problem can be solved in two steps. First one solves

183    (3.4)            $$\partial_t \begin{pmatrix} \rho \\ m \\ \mathcal{E} \\ \Gamma \end{pmatrix} + \partial_x \begin{pmatrix} m \\ \frac{1}{\rho}m^2 + p \\ \frac{m}{\rho}(\mathcal{E} + p) \\ \frac{m}{\rho}\Gamma \end{pmatrix} = 0, \quad \text{with} \quad p(\rho, m, \mathcal{E}, \Gamma) := \frac{\gamma - 1}{1 - b\rho}\left(\mathcal{E} - \frac{m^2}{2\rho}\right),$$

184    with left data and right data $(\rho_Z, \boldsymbol{m}_Z\cdot\boldsymbol{n}, \mathcal{E}_Z, \Gamma_Z)^\mathsf{T}$, where $\mathcal{E} := E - \frac{\|\boldsymbol{m}^\perp\|_{\ell^2}^2}{2\rho}$. Notice in passing
185    that $E - \frac{\|\boldsymbol{m}\|_{\ell^2}^2}{2\rho} = \mathcal{E} - \frac{m^2}{2\rho}$, i.e., the internal energy does not depend on the change of basis. This,

together with the definition of $\gamma_Z$, implies that $p_Z := \frac{\gamma_Z - 1}{1 - b\rho_Z}(\mathcal{E}_Z - \frac{m_Z^2}{2\rho_Z}) = \frac{\gamma_Z - 1}{1 - b\rho_Z}(E_Z - \frac{\|\boldsymbol{m}_Z\|_{\ell^2}^2}{2\rho_Z}) = p_Z$. Second, one obtains the full solution to the Riemann problem (3.3) by determining $\boldsymbol{m}^\perp$. This field is obtained by solving $\partial_t \boldsymbol{m}^\perp + \partial_x(v\boldsymbol{m}^\perp) = 0$ with the appropriate left and right data. Just like in the case of the Euler equations, one never solves the second step since it does not affect the maximum wave speed and the structure of the Riemann problem. In the rest of this paper we solely focus our attention on the system (3.4).

*Remark* 3.1 (Invariant domain properties).     At this point, it is important to notice that $\widetilde{\mathbb{f}}(\widetilde{\boldsymbol{u}}_Z) = (\mathbb{f}(\boldsymbol{u}_Z), \boldsymbol{v}_Z \Gamma_Z)^\mathsf{T}$ because, as already mentioned above, $\widetilde{p}(\widetilde{\boldsymbol{u}}_Z) = p_Z = p(\boldsymbol{u}_Z)$. Then, recalling (2.10), and setting $\lambda := \frac{d_{ij}^n}{\|\boldsymbol{c}_{ij}\|_{\ell^2}}$ and $\overline{\boldsymbol{u}}_{LR} := \overline{\mathbf{U}}_{ij}$, the extended auxiliary state based on the extended flux $\widetilde{\mathbb{f}}$, say $\overline{\widetilde{\boldsymbol{u}}}_{LR}$, satisfies the following identity:

$$(3.5) \qquad \overline{\widetilde{\boldsymbol{u}}}_{LR} = \begin{pmatrix} \overline{\boldsymbol{u}}_{LR} \\ \frac{1}{2}(\Gamma_L + \Gamma_R) - \frac{1}{2\lambda}(\boldsymbol{v}_R\Gamma_R - \boldsymbol{v}_L\Gamma_L)\cdot\boldsymbol{n}. \end{pmatrix}$$

That is, the density, the momentum, and the total energy of the states $\overline{\widetilde{\boldsymbol{u}}}_{LR}$ and $\overline{\boldsymbol{u}}_{LR}$ are identical. This implies that these two states have the same density and the same internal energy. As a result, if one can prove that the density and the internal energy of the state $\overline{\widetilde{\boldsymbol{u}}}_{LR}$ are both positive, then this conclusion automatically carries over to the state $\overline{\boldsymbol{u}}_{LR}$. This remark is essential, and it is the main motivation for our introducing the extended Riemann problem.     □

**3.2. The invariant domain preserving properties.** We will use the technique of Lax consisting of piecing together elementary waves to construct a weak solution to the extended Riemann problem (3.4). We will show that this weak solution preserves positivity of the density and the internal energy (see Remark 3.1). We will also show that the local gamma constant is uniformly bounded from bellow: $\gamma \geq \min(\gamma_L, \gamma_R)$. The key tool we are going to invoke is the following lemma.

LEMMA 3.2 (Riemann average).     *Let $m$ be a positive integer. Let $\mathcal{A}$ be a subset of $\mathbb{R}^m$. Let $\boldsymbol{g} \in C^1(\mathcal{A}; \mathbb{R}^m)$ be a one-dimensional flux. Let $\boldsymbol{w}_L, \boldsymbol{w}_R \in \mathcal{A}$. Assume that the following Riemann problem*

$$(3.6) \qquad \partial_t \boldsymbol{w} + \partial_x \boldsymbol{g}(\boldsymbol{w}) = \boldsymbol{0}, \quad \boldsymbol{w}(x,0) = \begin{cases} \boldsymbol{w}_L & x < 0, \\ \boldsymbol{w}_R & x > 0, \end{cases}$$

*has a weak solution $\boldsymbol{w}$ in $L^\infty(\mathbb{R} \times (0,\infty); \mathbb{R}^m) \cap C^0([0,\infty); L^1_{\mathrm{loc}}(\mathbb{R}; \mathbb{R}^m))$. Assume that this Riemann solution has a finite maximum wave speed (meaning, there exists $\lambda_{\max} > 0$ s.t. $\boldsymbol{w}(x,t) = \boldsymbol{w}_L$ if $x < -\lambda_{\max}t$ and $\boldsymbol{w}(x,t) = \boldsymbol{w}_R$ if $x > \lambda_{\max}t$.) Let $\mathcal{B}$ be a convex subset of $\mathcal{A}$ and assume that $\boldsymbol{w}(x,t) \in \mathcal{B}$ for a.e. $x \in \mathbb{R}$ and all $t > 0$. Let $\overline{\boldsymbol{w}} := \int_{-\frac{1}{2}}^{\frac{1}{2}} \boldsymbol{w}(x,t)\,\mathrm{d}x$. Then the following holds true for all $t \in (0, \frac{1}{2\lambda_{\max}})$:*
*(i) $\overline{\boldsymbol{w}}(t) = \frac{1}{2}(\boldsymbol{w}_L + \boldsymbol{w}_R) - (\boldsymbol{g}(\boldsymbol{w}_R) - \boldsymbol{g}(\boldsymbol{w}_L))t$;*
*(ii) $\overline{\boldsymbol{w}}(t) \in \mathcal{B}$;*
*(iii) Let $\Psi \in C^1(\mathcal{B}; \mathbb{R})$ be a quasiconcave functional. Assume that $\Psi(\boldsymbol{w}(x,t)) \geq 0$ for a.e. $x \in \mathbb{R}$ and all $t > 0$. Then $\Psi(\overline{\boldsymbol{w}}(t)) \geq 0$.*
*(iv) Let $\Psi \in C^1(\mathcal{B}; \mathbb{R})$ be a concave functional. Assume that $\Psi(\boldsymbol{w}(x,t)) \geq 0$ for a.e. $x \in \mathbb{R}$ and all $t > 0$. Assume that there exists $\lambda_\flat, \lambda_\sharp \in [-\lambda_{\max}, \lambda_{\max}]$, $\lambda_\flat < \lambda_\sharp$, so that $\Psi(\boldsymbol{w}(x,t)) > 0$ for a.e. $\frac{x}{t} \in (\lambda_\flat, \lambda_\sharp)$. Then $\Psi(\overline{\boldsymbol{w}}(t)) > 0$.*

*Proof.* (i) Let $w_1, \ldots, w_m$ be the $m$ components of $\boldsymbol{w}$, and let $g_1, \ldots, g_m$ be the $m$ components of the flux $\boldsymbol{g}$. Let $l \in \{1{:}m\}$. Since $\boldsymbol{w}$ is a weak solution to (3.6), we have

$$0 = \int_{-\infty}^{\infty} \int_0^{\infty} (-w_l \partial_\tau \phi - g_l(\boldsymbol{w}) \partial_x \phi) \, \mathrm{d}\tau \, \mathrm{d}x - w_{l,L} \int_{-\infty}^0 \phi(x, 0) \, \mathrm{d}x - w_{l,R} \int_0^{\infty} \phi(x, 0) \, \mathrm{d}x$$

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

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

Using that $w_l \in C^0([0, \infty); L^1_{\mathrm{loc}}(\mathbb{R}))$, we infer that $\int_{-\infty}^{\infty} \int_0^{\infty} -w_l \partial_\tau \phi_\epsilon \, \mathrm{d}x \, \mathrm{d}\tau \to \int_{-\frac{1}{2}}^{\frac{1}{2}} w_l(x, t) \, \mathrm{d}x$ as $\epsilon \to$ 0. Likewise, we have $\int_{-\infty}^{\infty} \int_0^{\infty} -g_l(\boldsymbol{w}) \partial_x \phi_\epsilon \, \mathrm{d}x \, \mathrm{d}\tau \to \int_0^t (g_l(\boldsymbol{w}_R) - g_l(\boldsymbol{w}_L)) \, \mathrm{d}\tau = (g_l(\boldsymbol{w}_R) - g_l(\boldsymbol{w}_L))t$ as $\epsilon \to 0$. Finally, $-w_{l,L} \int_{-\infty}^0 \phi_\epsilon(x, 0) \, \mathrm{d}x - w_{l,R} \int_0^{\infty} \phi_\epsilon(x, 0) \, \mathrm{d}x \to -\frac{1}{2}(w_{l,L} + w_{l,R})$ as $\epsilon \to 0$. In conclusion, we have established that

$$0 = \overline{\boldsymbol{w}}(t) + (\boldsymbol{g}(\boldsymbol{w}_R) - \boldsymbol{g}(\boldsymbol{w}_L))t - \tfrac{1}{2}(\boldsymbol{w}_L + \boldsymbol{w}_R).$$

(ii) Since $\mathscr{B}$ is convex, $\boldsymbol{w}(x, t) \in \mathscr{B}$ for a.e. $x \in \mathbb{R}$ and all $t > 0$, and the length of the interval $[-\frac{1}{2}, \frac{1}{2}]$ is 1, we infer that $\overline{\boldsymbol{w}}(t) \in \mathscr{B}$.

(iii) Let $\Psi \in C^1(\mathscr{B}; \mathbb{R})$ be a quasiconcave functional. The quasiconcavity implies that $\Psi(\overline{\boldsymbol{w}}(t)) \ge \operatorname{ess\,inf}_{x \in (-\frac{1}{2}, \frac{1}{2})} \Psi(\boldsymbol{w}(x, t)) \ge 0$.

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

$$\Psi(\overline{\boldsymbol{w}}(t)) \ge \int_{-\frac{1}{2}}^{\frac{1}{2}} \Psi(\boldsymbol{w}(x, t)) \, \mathrm{d}x \ge \int_{\lambda_\flat t}^{\lambda_\sharp t} \Psi(\boldsymbol{w}(x, t)) \, \mathrm{d}x > 0,$$

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

*Remark* 3.3 (Weak solution). Notice that Lemma 3.2 only requires us to have access to a weak solution of (3.6) that satisfies an invariant-domain property (i.e., $\boldsymbol{w}(x, t) \in \mathscr{B}$ for a.e. $x \in \mathbb{R}$ and all $t > 0$). No entropy inequality or additional smoothness condition is needed. $\qquad \square$

**4. Solution of the extended Riemann problem.** We now construct a weak solution to the extended Riemann problem (3.4) using the technique described in Lax [19] (we also refer to Holden and Risebro [15, Chap. 5], Godlewski and Raviart [10, Chap. 1], Toro [27, Chap. 4] for further details on the Riemann problem). No originality is claimed on this construction, but we give the details for completeness.

**4.1. Definition of the star states.** We first notice that the Jacobian matrix of (3.4) is diagonalizable and has three distinct eigenvalues. The eigenvalue $\frac{m}{\rho}$ has multiplicity 2. Then, as usual, we postulate that the solution to (3.4) is self-similar and composed of three waves hereafter called L-wave, C-wave, and R-wave. The L-wave and the R-wave are either shocks or expansions. The L-wave will be generated using the covolume equation of state with $\gamma_L$ and the R-wave will be generated by using the covolume equation of state with $\gamma_R$. The C-wave is a contact discontinuity

for the density and $\Gamma$. Compared to the technique described in Toro [27, Chap. 4], the only new feature here is that the dependent variable has a fourth component $\Gamma$. The purpose of this section is to introduce quantities that are useful to define the three waves in question: the intermediate densities $\rho_L^*$, $\rho_R^*$, the intermediate velocities $v_L^*$, $v_R^*$, $v^*$, and the intermediate pressure $p^*$. The actual construction of the solution is done in §4.2 and §4.3.

In the rest of this section we use the primitive variables: density $\rho$, velocity $v$, pressure $p$, and $\gamma := \Gamma/\rho$. We use the symbol $p$ to denote the pressure defined in (3.4). Notice that the oracle is only invoked to compute the two states $p_L$ and $p_R$. We define the primitive state $\boldsymbol{c} := (\rho, v, p, \gamma)^\mathsf{T}$ and set $\boldsymbol{c}_Z := (\rho_Z, v_Z, p_Z, \gamma_L)^\mathsf{T}$. Recalling that we have defined $\gamma_Z := 1 + \frac{p_Z(1-b\rho_Z)}{\rho_Z e_Z}$, the oracle assumption (2.4) implies that $\min(\gamma_L, \gamma_R) > 1$.

We define the covolume sound speed $a_Z := \sqrt{\frac{\gamma_Z p_Z}{\rho_Z(1-b\rho_Z)}}$, the parameters $A_Z := \frac{2(1-b\rho_Z)}{(\gamma_Z+1)\rho_Z}$ and $B_Z := \frac{\gamma_Z-1}{\gamma_Z+1}p_Z$ corresponding to the $Z$ state (see e.g., Toro [27, §4.7], [11]), and introduce the following function

$$(4.1) \qquad f_Z(p) := \begin{cases} f_Z^R(p) := \frac{2a_Z(1-b\rho_Z)}{\gamma_Z-1}\left(\left(\frac{p}{p_Z}\right)^{\frac{\gamma_Z-1}{2\gamma_Z}} - 1\right) & \text{if } 0 \le p < p_Z, \\[3mm] f_Z^S(p) := (p-p_Z)\left(\frac{A_Z}{p+B_Z}\right)^{\frac{1}{2}} & \text{if } p_Z \le p. \end{cases}$$

The definition of $f_Z(p)$ makes sense because $1 < \gamma_Z$ and $0 \le B_Z$. It is shown in Toro [27, §4.3.1] that the function $f_Z(p)$ is in $C^2(\mathbb{R}_+; \mathbb{R})$, monotone increasing, and concave.

We also define the function $\phi \in C^2(\mathbb{R}_+; \mathbb{R})$,

$$(4.2) \qquad \phi(p) := f_L(p) + f_R(p) + v_R - v_L, \qquad p \in [0, \infty).$$

Notice in passing that assuming $\phi(0) < 0$ is equivalent to assuming that the following holds true:

$$(4.3) \qquad v_R - v_L < \frac{2a_L(1-b\rho_L)}{\gamma_L - 1} + \frac{2a_R(1-b\rho_R)}{\gamma_R - 1}.$$

This condition is known in the literature as the non-vacuum condition (see Toro [27, (4.40), p. 127]).

LEMMA 4.1. *If* (4.3) *holds, then $\phi$ has a unique positive root $p^*$.*

*Proof.* Since $\phi(0) = v_R - v_L - \frac{2a_L(1-b\rho_L)}{\gamma_L-1} - \frac{2a_R(1-b\rho_R)}{\gamma_R-1}$, the assumption (4.3) means that $\phi(0) < 0$. We then conclude that $\phi$ has a unique positive root since $\phi(p) \in C^2(\mathbb{R}_+; \mathbb{R})$ is strictly monotone increasing (and concave). $\blacksquare$

DEFINITION 4.2 ($p^*$, $\rho_L^*$, $\rho_R^*$, $v_L^*$, $v_R^*$, $v^*$). (i) *If the non-vacuum condition* (4.3) *holds, we denote by $p^*$ the unique root of $\phi$, and we set $v_L^* := v_L - f_L(p^*)$, $v_R^* := v_R + f_R(p^*)$, $v^* := v_L^* = v_R^*$.* (ii) *If instead there is vacuum, we define $p^* := 0$ and set $v_L^* := v_L - f_L(0)$, $v_R^* := v_R + f_R(0)$.* (iii) *We set $\rho_L^* = \rho_R^* = 0$ if $p^* = 0$; otherwise we set $\rho_Z^* := \left(b + \frac{1-b\rho_Z}{\rho_Z}\left(\frac{p_Z}{p^*}\right)^{\frac{1}{\gamma_Z}}\right)^{-1}$, $Z \in \{L, R\}$.*

Notice that the definition of $v^*$ makes sense if the non-vacuum condition (4.3) holds since in this case $\phi(p^*) = 0 = v_R^* - v_L^*$. The definition of $\rho_Z^*$ is continuous with respect to $p^*$, including at $p^* = 0$. To fully describe our weak solution, we introduce the following wave speeds:

$$\lambda_L^-(p^*) := v_L - a_L\left(1 + \frac{\gamma_L + 1}{2\gamma_L}\left(\frac{p^* - p_L}{p_L}\right)_+\right)^{\frac{1}{2}},$$

290

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

291

$$\lambda_R^+(p^*) := v_R + a_R \left(1 + \frac{\gamma_R+1}{2\gamma_R}\left(\frac{p^*-p_R}{p_R}\right)_+\right)^{\frac{1}{2}},$$

292

293

$$\lambda_R^-(p^*) := \begin{cases} v_R + f_R(p^*) + a_R \frac{1-b\rho_R}{1-b\rho_R^*} \left(\frac{p^*}{p_R}\right)^{\frac{\gamma_R-1}{2\gamma_R}} & \text{if } p^* < p_R, \\ \lambda_R^+(p^*) & \text{if } p_R \le p^*, \end{cases}$$

294

295 LEMMA 4.3 (wave speeds). *Assume* $1 < \min(\gamma_L, \gamma_R)$ *and* $0 < a_L, a_R$. *Then, the following*
296 *holds true:*

297 (4.4) $$\lambda_L^-(p^*) \le \lambda_L^+(p^*) \le v_L^* \le v_R^* \le \lambda_R^-(p^*) \le \lambda_R^+(p^*).$$

298 *Proof.* We will only consider the case $Z = L$; the case $Z = R$ is analogous. There are two
299 possibilities: either $p^* < p_L$ or $p_L \le p^*$. In the first case, $p^* < p_L$, we have

300 $$\lambda_L^-(p^*) = v_L - a_L < \lambda_L^+(p^*) = v_L - f_L(p^*) - a_L \frac{1-b\rho_L}{1-b\rho_L^*}\left(\frac{p^*}{p_L}\right)^{\frac{\gamma_L-1}{2\gamma_L}} \le v_L - f_L(p^*) = v_L^*,$$

301 where we used above that $f_L(p^*) < 0$, $1 < \gamma_L$, $0 < a_L$, $\rho_L^* < \rho_L$, $0 \le p^* < p_L$ and $0 < \frac{1-b\rho_L}{1-b\rho_L^*} \le 1$.
302 In the second case, $p_L \le p^*$, we have

303 $$\lambda_L^-(p^*) = \lambda_L^+(p^*) = v_L - a_L\left(1 + \frac{\gamma_L+1}{2\gamma_L}\left(\frac{p^*-p_L}{p_L}\right)\right)^{\frac{1}{2}}$$

304 and

305 $$v_L^* = v_L - f_L(p^*) = v_L - (p^* - p_L)\left(\frac{A_L}{p^* + B_L}\right)^{\frac{1}{2}}.$$

306 Then proving the inequality $\lambda_L^+(p^*) < v_L^*$ is equivalent to showing that

307 $$\left(\frac{p^*}{p_L} - 1\right)\left(\frac{2(1-b\rho_L)}{\gamma_L(\gamma_L+1)}\frac{\gamma_L p_L}{\rho_L}\frac{1}{\frac{p^*}{p_L}+\frac{\gamma_L-1}{\gamma_L+1}}\right)^{\frac{1}{2}} < a_L\left(1 + \frac{\gamma_L+1}{2\gamma_L}\left(\frac{p^*-p_L}{p_L}\right)\right)^{\frac{1}{2}}.$$

308 Using the substitution $x := \frac{p^*}{p_L} - 1$ and that $a_L := \sqrt{\frac{\gamma_L p_L}{\rho_L(1-b\rho_L)}}$, we derive that the above inequality
309 is equivalent to proving that

310 $$\left(\frac{2}{\gamma_L(\gamma_L+1)}\right)^{\frac{1}{2}} x(1-b\rho_L) < \left(\left(x + \frac{2\gamma_L}{\gamma_L+1}\right)\left(\frac{\gamma_L+1}{2\gamma_L}x+1\right)\right)^{\frac{1}{2}}$$

311 where $x > 0$. Squaring both sides, and recalling that $x > 0$, we observe that the above is equivalent
312 to the inequality

313 $$0 < \left(\frac{\gamma_L+1}{2\gamma_L} - \frac{2(1-b\rho_L)^2}{\gamma_L(\gamma_L+1)}\right)x^2 + 2x + \frac{2\gamma_L}{\gamma_L+1}.$$

314 This inequality holds true for all $x \ge 0$ since we assumed that $1 < \gamma_L$ and $0 \le 1 - b\rho_L \le 1$. □

315    **4.2. Definition of the L-wave and R-wave without vacuum.** We assume in this section
316    that the non-vacuum condition (4.3) holds. The main result of this section is Lemma 4.4. The
317    solution with vacuum is given in §4.3.
318        Recalling the notation from Definition 4.2, the proposed solution to (3.4) is self-similar and has
319    the following form:

320    (4.5)
$$
\boldsymbol{c}(x,t) := \begin{cases} \boldsymbol{c}_L & \text{if } \frac{x}{t} < \lambda_L^-, \\ \boldsymbol{c}_{LL}(\frac{x}{t}) & \text{if } \lambda_L^- \leq \frac{x}{t} < \lambda_L^+, \\ \boldsymbol{c}_L^* & \text{if } \lambda_L^+ \leq \frac{x}{t} < v*, \\ \boldsymbol{c}_R^* & \text{if } v* \leq \frac{x}{t} < \lambda_R^-, \\ \boldsymbol{c}_{RR}(\frac{x}{t}) & \text{if } \lambda_R^- \leq \frac{x}{t} < \lambda_R^+, \\ \boldsymbol{c}_R & \text{if } \lambda_R^+ \leq \frac{x}{t}. \end{cases}
$$

321    with $\boldsymbol{c}_L^* := (\rho_L^*, v^*, p^*, \gamma_L)^\mathsf{T}$ and $\boldsymbol{c}_R^* := (\rho_R^*, v^*, p^*, \gamma_R)^\mathsf{T}$. The parameters $p^*$, $v^*$, $\rho_L^*$, and $\rho_R^*$ are
322    defined in Definition 4.2. The two functions $\boldsymbol{c}_{LL}$, $\boldsymbol{c}_{RR}$ are going to be defined to make sure that
323    (4.5) is indeed a weak solution to (3.4). Notice that $\boldsymbol{c}$ is uniquely defined owing to Lemma 4.3 (i.e.,
324    the waves are well ordered).
325        Let us first construct the L-wave, i.e., we construct the function $\boldsymbol{c}_{LL}(\xi)$ where $\lambda_L^- \leq \xi < \lambda_L^+$.
326    If $p_L \leq p^*$, then $\lambda_L^-(p^*) = \lambda_L^+(p^*)$ and the L-wave is a shock. In this case one does not need to
327    define $\boldsymbol{c}_{LL}$ since the interval $[\lambda_L^-, \lambda_L^+)$ is empty. If $p^* < p_L$, we postulate that the $\gamma$-component of
328    $\boldsymbol{c}_{LL}$ is constant and equal to $\gamma_L$. This means that the L-wave can be computed by assuming that
329    the equation of state is a standard co-volume $\gamma$-law $p(1 - b\rho) = (\gamma_L - 1)\rho e$ (with $e = \frac{1}{\rho}(\mathcal{E} - \frac{m^2}{2\rho})$).
330    In this case the L-wave is an expansion. The construction of this wave is well established, we refer
331    for instance to Toro [27, Chap. 4]. More precisely, the self-similarity parameter $\xi = \frac{x}{t}$ (which is the
332    eigenvalue of the Jacobian of the flux, $v - a$) can be expressed in terms of the parameter $p$:

333    (4.6)        $\xi_L(p) := v_L - f_L(p) - a_L \dfrac{1 - b\rho_L}{1 - b\rho(p)} \left(\dfrac{p}{p_L}\right)^{\frac{\gamma_L - 1}{2\gamma_L}}, \qquad p \in [p^*, p_L],$

334    where $\rho(p)$ is defined as follows:

335        $\dfrac{1}{\rho(p)} - b := \left(\dfrac{1}{\rho_L} - b\right) \left(\dfrac{p_L}{p}\right)^{\frac{1}{\gamma_L}}.$

336    To simplify the notation we use the symbol $\xi(p)$ instead of $\xi_L(p)$ when the context is unambiguous.
337    Notice in passing that $\lambda_L^-(p^*) = \xi(p_L)$ and $\lambda_L^+(p^*) = \xi(p^*)$. Since the function $\xi$ is strictly deceasing
338    in the interval $p \in [p^*, p_L]$, the inverse function theorem implies that $p$ can be uniquely expressed
339    in terms of $\xi$. We abuse the notation and denote by $p(\xi)$ the inverse function. Over the interval
340    $\xi \in [\xi(p_L), \xi(p^*)] = [\lambda_L^-(p^*), \lambda_L^+(p^*)]$, we have (see Toro [27, §4.7.1])

341    (4.7)        $\boldsymbol{c}_{LL}(\xi) := \left(\rho_L\left(b\rho_L + (1 - b\rho_L)\left(\dfrac{p_L}{p(\xi)}\right)^{\frac{1}{\gamma_L}}\right)^{-1}, v_L - f_L(p(\xi)), p(\xi), \gamma_L\right)^\mathsf{T}.$

342        Now we define $\boldsymbol{c}_L^*$. If $p^* < p_L$, the L-wave is an expansion and $\boldsymbol{c}_L^*$ is defined to be the end point
343    of the L-wave: $\boldsymbol{c}_L^* := \boldsymbol{c}_{LL}(\xi(p^*))$. If $p_L \leq p^*$, the L-wave is a shock. We still postulate that the

344  $\gamma$-component of $\boldsymbol{c}$ is equal to $\gamma_L$ for $\frac{x}{t} \leq \lambda_L^+(p^*)$. In this case we define $\boldsymbol{c}_L^*$ so that the Rankine–
345  Hugoniot relation holds between the two state $\boldsymbol{c}_L$ and $\boldsymbol{c}_L^*$ (see Toro [27, §4.7.1]). In conclusion, we
346  have

347  (4.8)
$$\boldsymbol{c}_L^* := \begin{cases} \boldsymbol{c}_{LL}(\xi(p^*)) & \text{if } p^* < p_L, \\ \left( \dfrac{\rho_L\left(\frac{p^*}{p_L} + \frac{\gamma_L-1}{\gamma_L+1}\right)}{\frac{\gamma_L-1+2b\rho_L}{\gamma_L+1}\frac{p^*}{p_L} + \frac{\gamma_L+1-2b\rho_L}{\gamma_L+1}}, v_L - f_L(p^*), p^*, \gamma_L \right)^{\mathsf{T}} & \text{if } p_L \leq p^*. \end{cases}$$

348  We define $\boldsymbol{c}_{RR}(\xi)$ similarly. If $p^* < p_R$, the R-wave is an expansion, otherwise it is a shock.
349  Assuming that $p^* < p_R$, the self-similarity parameter $\xi = \frac{x}{t}$ can be expressed in terms of the
350  parameter $p \in [p^*, p_R]$:

351  (4.9)
$$\xi_R(p) := v_R + f_R(p) + a_R \frac{1 - b\rho_R}{1 - b\rho(p)} \left( \frac{p}{p_R} \right)^{\frac{\gamma_R-1}{2\gamma_R}}$$

352  where we have defined

353
$$\frac{1}{\rho(p)} - b := \left( \frac{1}{\rho_R} - b \right) \left( \frac{p_R}{p} \right)^{\frac{1}{\gamma_R}}.$$

354  To simplify the notation we use the symbol $\xi(p)$ instead of $\xi_R(p)$ when the context is unambiguous.
355  Notice that in this case $\lambda_R^- = \xi(p^*)$, $\lambda_R^+ = \xi(p_R)$, and $\xi$ is a strictly increasing function over the
356  interval $[p^*, p_R]$. Over the interval $\xi \in [\xi(p^*), \xi(p_R)]$, we have

357  (4.10)
$$\boldsymbol{c}_{RR}(\xi) := \left( \rho_R\left( b\rho_R + (1 - b\rho_R)\left( \frac{p_R}{p(\xi)} \right)^{\frac{1}{\gamma_R}} \right)^{-1}, v_R + f_R(p(\xi)), p(\xi), \gamma_R \right)^{\mathsf{T}}.$$

358  Now we define $\boldsymbol{c}_R^*$. If $p^* < p_R$, the R-wave is an expansion and $\boldsymbol{c}_R^*$ is defined to be the end
359  point of the wave: $\boldsymbol{c}_R^* = \boldsymbol{c}_{RR}(\xi(p^*))$. If $p_R \leq p^*$, the R-wave is a shock. We still postulate that
360  the $\gamma$-component of $\boldsymbol{c}$ is equal to $\gamma_R$ for $v^* \leq \frac{x}{t} < \lambda_R^+$. In this case we define $\boldsymbol{c}_R^*$ so that the
361  Rankine–Hugoniot relation holds between the two state $\boldsymbol{c}_R$ and $\boldsymbol{c}_R^*$. In conclusion, we have

362  (4.11)
$$\boldsymbol{c}_R^* = \begin{cases} \boldsymbol{c}_{RR}(\xi(p^*)) & \text{if } p^* < p_R, \\ \left( \dfrac{\rho_R\left(\frac{p^*}{p_R} + \frac{\gamma_R-1}{\gamma_R+1}\right)}{\frac{\gamma_R-1+2b\rho_R}{\gamma_R+1}\frac{p^*}{p_R} + \frac{\gamma_R+1-2b\rho_R}{\gamma_R+1}}, v_R + f_R(p^*), p^*, \gamma_R \right)^{\mathsf{T}} & \text{if } p_R \leq p^*. \end{cases}$$

363  The key result of this section is summarized in the following Lemma.

364  LEMMA 4.4. *Assume that the non-vacuum condition* (4.3) *holds. The field* $(\rho, \boldsymbol{m}, E, \Gamma)^{\mathsf{T}}$ *defined*
365  *by* (4.5) *is a weak solution to* (3.4).

366  *Proof.* In the domain $\{x < v^*t\}$, we have $\gamma = \gamma_L$; hence, $\Gamma = \gamma_L\rho$. This implies that the
367  last equation in (3.4) is equivalent to the first equation (the conservation of mass). Moreover,
368  the first three equations in (3.4) hold true in the weak sense since the field $(\rho, m, \mathcal{E})$ defined in
369  (4.5) is by construction a weak solution to the regular Euler equations with the pressure law
370  $p(1 - b\rho) := (\gamma_L - 1)\left( \mathcal{E} - \frac{m^2}{2\rho} \right)$.
371  Similarly, in the domain $\{x > v^*t\}$, we have $\gamma = \gamma_R$; hence, $\Gamma = \gamma_R\rho$ and the last equation
372  in (3.4) is equivalent to the the conservation of mass equation. The first three equations in (3.4)

hold true in the weak sense because the field $(\rho, m, \mathcal{E})$ defined in (4.5) is by construction a weak solution to the regular Euler equations with a pressure law $p(1 - b\rho) := (\gamma_R - 1)\left(\mathcal{E} - \frac{m^2}{2\rho}\right)$.

To be able to conclude the proof, we now have to make sure that the two states that are separated by the line $\{x = v^*t\}$ satisfy the Rankine–Hugoniot relation. Let $\boldsymbol{c}_L^* = (\rho_L^*, v_L^*, p_L^*, \gamma_L)$ and $\boldsymbol{c}_R^* = (\rho_L^*, v_L^*, p_L^*, \gamma_L)$ be the two constant states defined above. Recall that the construction of $\boldsymbol{c}_L^*$ and $\boldsymbol{c}_R^*$ is such that that $p_L^* = p_R^* = p^*$ (see (4.8) and (4.11)). We have to show that

$$\rho_L^* v_L^* - \rho_R^* v_R^* = v^*(\rho_L^* - \rho_R^*)$$

$$\rho_L^*(v_L^*)^2 + p_L - \rho_R^*(v_R^*)^2 - p_R = v^*(\rho_L^* v_L^* - \rho_R^* v_R^*)$$

$$v_L^*(E_L^* - p_L^*) - v_R^*(E_R^* - p_R^*) = v^*(E_L^* - E_R^*),$$

$$v_L^* \gamma_L - v_R^* \gamma_R = v^*(\gamma_L - \gamma_R).$$

Since the non-vacuum condition (4.3) holds, we have $v^* := v_L^* = v_R^*$ (see Definition 4.2). Then it follows that the above four equations indeed hold true. Therefore the field defined in (4.5) is a weak solution to (3.4). □

**4.3. Definition of the L-wave and R-wave when vacuum is present.** When (4.3) fails, the solution contains a vacuum state. In this case both the L-wave and the R-waves are expansions. Recall that in Definition 4.2 we have set
(4.12)

$$p^* := 0, \qquad v_L^* := v_L - f_L(0) = v_L + \frac{2a_L(1 - b\rho_L)}{\gamma_L - 1}, \qquad v_R^* := v_R + f_R(0) = v_R - \frac{2a_R(1 - b\rho_R)}{\gamma_R - 1}.$$

The solution to the extended Riemann problem (3.4) we propose is as follows:

(4.13)
$$\boldsymbol{c}(x, t) = \begin{cases} \boldsymbol{c}_L & \text{if } \frac{x}{t} < v_L - a_L, \\ \boldsymbol{c}_{LL}(\frac{x}{t}) & \text{if } \lambda_L^- \leq \frac{x}{t} < v_L^*, \\ \frac{v_R^* - \frac{x}{t}}{v_R^* - v_L^*}\boldsymbol{c}_L^* + \frac{\frac{x}{t} - v_L^*}{v_R^* - v_L^*}\boldsymbol{c}_R^* & \text{if } v_L^* \leq \frac{x}{t} < v_R^*, \\ \boldsymbol{c}_{RR}(\frac{x}{t}) & \text{if } v_R^* \leq \frac{x}{t} < v_R + a_R, \\ \boldsymbol{c}_R & \text{if } v_R + a_R \leq \frac{x}{t}. \end{cases}$$

The definitions of the expansion waves $\boldsymbol{c}_{LL}$ and $\boldsymbol{c}_{RR}$ are the same as in the non-vacuum case. We define the states $\boldsymbol{c}_L^*$ and $\boldsymbol{c}_R^*$ as in §4.2 by setting $\boldsymbol{c}_L^* := \boldsymbol{c}_{LL}(v_L^*) = (0, v_L^*, 0, \gamma_L)^\mathsf{T}$ and $\boldsymbol{c}_R^* := \boldsymbol{c}_{RR}(v_R^*) = (0, v_R^*, 0, \gamma_R)^\mathsf{T}$. The key result of this section is the following Lemma.

LEMMA 4.5. *Assume that the vacuum condition holds, i.e., $p^* = 0$. The field $(\rho, \boldsymbol{m}, E, \Gamma)^\mathsf{T}$ defined by (4.13) is a weak solution to (3.4).*

*Proof.* We have already established that, once expressed in conserved variable, (4.13) is a weak solution to (3.4) in the regions $\{x < v_L^* t\} \cup \{v_R^* t < x\}$. In the region $\{v_L^* t < x < v_R^* t\}$, all the conserved variables are zero by construction. Hence, (4.13) rewritten in conserved variables is also weak solution to (3.4) in the region $\{v_L^* t < x < v_R^* t\}$. Let us verify now that the field defined in (4.13) is continuous across the line $\{x = v_L^* t\}$. Denoting $\xi_L(p)$ the function defined in (4.6), we obtain $\xi_L(0) = v_L - f_L(0) =: v_L^*$, i.e., $p(v_L^*) = 0$. Hence $\lim_{\xi \uparrow v_L^*} \boldsymbol{c}_{LL}(\xi) = (0, v_L^*, 0, \gamma_L)$. Moreover, $\lim_{\xi \downarrow v_L^*} \frac{v_R^* - \xi}{v_R^* - v_L^*}\boldsymbol{c}_L^* + \frac{\xi - v_L^*}{v_R^* - v_L^*}\boldsymbol{c}_R^* = (0, v_L^*, 0, \gamma_L)$. This proves the assertion. This in turn establishes that the conserved field is also continuous across $\{x = v_L^* t\}$. The argument to prove continuity across $\{x = v_R^* t\}$ is similar. The conclusion follows readily. □

407     **4.4. Summary.** In Sections §4.2 and §4.3 we have defined a weak solution to the extended
408 Riemann problem (3.4). Notice that this weak solution satisfies the assumption of Lemma 3.2,
409 i.e., it is in $L^\infty(\mathbb{R}\times(0,\infty);\mathbb{R}^m)\cap C^0([0,\infty);L^1_{\mathrm{loc}}(\mathbb{R};\mathbb{R}^m))$ with $m = d + 2$, and the maximum wave
410 speed $\lambda_{\max} = \max(|\lambda_L^-(p^*)|, |\lambda_R^+(p^*)|) = \max(-\lambda_L^-(p^*), \lambda_R^+(p^*))$ is finite. As a result, we can invoke
411 Lemma 3.2 for any quasiconcave functional. The following theorem is the main result of §4.

412     THEOREM 4.6. (i) *Let* $\mathbf{U}_i^n$, $\mathbf{U}_j^n$ *be two states in* $\mathcal{B}(b)$ *(with* $\mathcal{B}(b)$ *defined in* (2.3)*). Let* $p^*$ *be*
413 *defined as in Definition 4.2 with left state* $\mathbf{U}_i^n$ *and right state* $\mathbf{U}_j^n$. *Let* $\widehat{p}^*$ *be any upper bound on* $p^*$
414 *(i.e.,* $\widehat{p}^* \ge p^*$*). Let*

415 (4.14a) $$\widehat{\lambda}(\boldsymbol{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n) := \max(-\lambda_L^-(\widehat{p}^*), \lambda_R^+(\widehat{p}^*)),$$

416 (4.14b) $$d_{ij}^n := \max(\widehat{\lambda}(\boldsymbol{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n)\|\boldsymbol{c}_{ij}\|_{\ell^2}, \widehat{\lambda}(\boldsymbol{n}_{ji}, \mathbf{U}_j^n, \mathbf{U}_i^n)\|\boldsymbol{c}_{ji}\|_{\ell^2}).$$
417

418 *Let* $\overline{\mathbf{U}}_{ij}^n$ *be defined by* (2.10). *Then* $\overline{\mathbf{U}}_{ij}^n \in \mathcal{B}(b)$.
419 (ii) *Let* $i \in \mathcal{V}$. *Assume that* $\mathbf{U}_j^n \in \mathcal{B}(b)$ *for all* $j \in \mathcal{I}(i)$. *Assume that* $d_{ij}^n$ *is defined as above in*
420 (4.14b) *for all* $j \in \mathcal{I}(i)$. *Assume that* $\tau$ *is small enough so that* $\tau \sum_{j\in\mathcal{I}(i)\setminus\{i\}} \frac{2d_{ij}^n}{m_i} \le 1$. *Let* $\mathbf{U}_i^{n+1}$ *be*
421 *the update defined in* (2.7). *Then* $\mathbf{U}_i^{n+1} \in \mathrm{Conv}\{\overline{\mathbf{U}}_{ij}^n \mid j \in \mathcal{I}(i)\} \subset \mathcal{B}(b)$.

422     *Proof.* (i) We first notice that $\widehat{\lambda}(\boldsymbol{n}_{ij}, \mathbf{U}_i^n, \mathbf{U}_j^n) \ge \max(-\lambda_L^-(p^*), \lambda_R^+(p^*)) =: \lambda_{\max}$ since the func-
423 tions $-\lambda_L^-$ and $\lambda_R^+$ are monotone increasing and $\widehat{p}^* \ge p^*$. We now apply Lemma 3.2 with the flux
424 $\boldsymbol{g}(\widetilde{\boldsymbol{w}}) = \widetilde{\mathbb{f}}(\widetilde{\boldsymbol{w}})\boldsymbol{n}$ and the Riemann data $\widetilde{\mathbf{U}}_i^n$, $\widetilde{\mathbf{U}}_j^n$. We observe that the Riemann solution defined in
425 (4.5) and (4.13) has nonnegative density and nonnegative internal energy (recall that the internal
426 energy $\rho e$ is equal to $\frac{1}{(\gamma-1)}(1 - b\rho)p$). Notice also that the only way to have zero density and zero
427 internal energy on a set of nonzero measure is when vacuum is present in the solution and $v_L^* < v_R^*$;
428 in this case, $\lambda_L^- < \lambda_L^+$ and $\lambda_R^- < \lambda_R^+$ and the density and the internal energy are positive in the
429 regions $\frac{x}{t} \in [\lambda_L^-, \lambda_L^+)$, $\frac{x}{t} \in (\lambda_R^-, \lambda_R^+]$. Consider the concave functionals $\widetilde{\Psi}_1 : \widetilde{\boldsymbol{u}} \mapsto \rho$, $\widetilde{\Psi}_2 : \widetilde{\boldsymbol{u}} \mapsto 1 - b\rho$,
430 and $\widetilde{\Psi}_3 : \widetilde{\boldsymbol{u}} \mapsto \rho e$. Notice that $\widetilde{\Psi}_l(\widetilde{\mathbf{U}}_j^n) > 0$ for all $j \in \mathcal{I}(i)$ and all $l \in \{1{:}3\}$ whether vacuum
431 occurs or not. We conclude that $\widetilde{\Psi}_l(\overline{\widetilde{\mathbf{U}}}_{ij}^n) > 0$ for all $l \in \{1{:}3\}$ by invoking Item (iv) in Lemma 3.2.
432 But the identity (3.5) shows that the density and the internal energy of the states $\overline{\widetilde{\mathbf{U}}}_{ij}^n$ and $\overline{\mathbf{U}}_{ij}^n$
433 are identical; as a result, defining $\Psi_1 : \boldsymbol{u} \mapsto \rho$, $\Psi_2 : \boldsymbol{u} \mapsto 1 - b\rho$, and $\Psi_3 : \boldsymbol{u} \mapsto \rho e$, we infer that
434 $\Psi_l(\overline{\mathbf{U}}_{ij}^n) = \widetilde{\Psi}_l(\overline{\widetilde{\mathbf{U}}}_{ij}^n) > 0$ for all $l \in \{1{:}3\}$. This establishes that $\overline{\mathbf{U}}_{ij}^n \in \mathcal{B}(b)$.
435 (ii) The assertion follows from (i), the convexity of $\mathcal{B}(b)$, and the observation that (2.9) implies that
436 $\mathbf{U}_i^{n+1}$ is in the convex hull of $\{\overline{\mathbf{U}}_{ij}^n \mid j \in \mathcal{I}(i)\}$ if $\tau \sum_{j\in\mathcal{I}(i)\setminus\{i\}} \frac{2d_{ij}^n}{m_i} \le 1$. This completes the proof. $\square$

437     Theorem 4.6 says that the algorithm (2.7) is invariant-domain preserving under the appropriate
438 CFL condition. To make this theorem useful, we now need to derive a computable upper bound on
439 the maximum wave speed in the extended Riemann problem (3.4). This task is achieved in §5.

440     **5. Upper bound on the maximum wave speed.** Setting $\lambda_{\max}(p) := \max(-\lambda_L^-(p), \lambda_R^+(p))$,
441 we recall that the maximum wave speed in the Riemann problem (3.4) is given by $\lambda_{\max}(p^*)$. Recall
442 also that $p \mapsto \lambda_{\max}(p)$ is a nondecreasing function. Since we only need an upper bound on $\lambda_{\max}(p^*)$,
443 we derive in this section an explicit upper bound on $p^*$.

444     **5.1. Motivation and notation.** We recall that $p^* = 0$ if vacuum is present, and the max-
445 imum speed of propagation is then $\lambda_{\max}(0) = \max(|v_L - a_L|, |v_R + a_R|)$. (The L-wave and the
446 R-wave are both expansions in this case.) If the non-vacuum condition holds (see (4.3)), $p^*$ solves

the equation

$$\phi(p) = f_L(p) + f_R(p) + v_R - v_L = 0, \qquad p \in (0, \infty).$$

As proved in Guermond and Popov [11, Lem. 4.2], a simple upper bound for $p^*$ can be obtained by using the so called double-rarefaction approximation (see also Pike [23]), which consists of finding the unique root of the modified equation $\phi_{RR}(p) = 0$, where

$$(5.1) \quad \phi_{RR}(p) := \frac{2a_L(1 - b\rho_L)}{\gamma_L - 1}\left(\left(\frac{p}{p_L}\right)^{\frac{\gamma_L - 1}{2\gamma_L}} - 1\right) + \frac{2a_L(1 - b\rho_R)}{\gamma_R - 1}\left(\left(\frac{p}{p_R}\right)^{\frac{\gamma_R - 1}{2\gamma_R}} - 1\right) + v_R - v_L.$$

It can be shown that $\phi_{RR}(p) \le \phi(p)$ for all $p \in [\min(p_L, p_R), \infty)$ if $\max(\gamma_L, \gamma_R) \in (1, \frac{5}{3}]$. Using the notation from (4.1), this result is proved in [11, Lem. 4.2] by showing that $f_Z^S(p) \ge f_Z^R(p)$ for all $p > p_Z$ if $\gamma_Z \in (1, \frac{5}{3}]$. We revisit this idea in the rest of §5 and remove the assumption $\max(\gamma_L, \gamma_R) \in (1, \frac{5}{3}]$. More precisely, we use a result from Theorem A.2 proved in Appendix A: there exists a function $c(\gamma_Z)$ (defined in (A.3)) so that $f_Z^S(p) \ge c(\gamma_Z)f_Z^R(p)$ for all $p > p_Z$. This function is equal to 1 over the range $\gamma_Z \in (1, \frac{5}{3}]$ and decreases monotonically to $\frac{1}{\sqrt{2}}$ as $\gamma_Z$ grows to infinity. To simplify the notation, let us set $\alpha_Z := c(\gamma_Z)\frac{2a_Z(1-b\rho_Z)}{\gamma_Z - 1}$. We then redefine $\phi_{RR}$ for all $\gamma_Z \in (1, \infty)$ by setting

$$(5.2) \qquad \phi_{RR}(p) := \alpha_L\left(\left(\frac{p}{p_L}\right)^{\frac{\gamma_L - 1}{2\gamma_L}} - 1\right) + \alpha_R\left(\left(\frac{p}{p_R}\right)^{\frac{\gamma_R - 1}{2\gamma_R}} - 1\right) + v_R - v_L.$$

We then have $\phi_{RR}(p) \le \phi(p)$ for all $p \in [\min(p_L, p_R), \infty)$ and all $\gamma_Z \in (1, \infty)$.

When $\gamma_L = \gamma_R$ (i.e., the case of the ideal gas law) the equation $\phi_{RR}(p) = 0$ can be easily solved since it is linear up to a trivial change of variable. But solving $\phi_{RR}(p) = 0$ in the general case (i.e., $\gamma_L \ne \gamma_R$) is far more difficult since the equation is nonlinear. In the rest of §5 we extract further lower bounds on $\phi_{RR}$ to derive an explicit upper bound on $p^*$.

To simplify the notation in many of the expressions used below, we introduce two indices in the set $\{L, R\}$ denoted by "min" and "max" and defined as follows:

$$(5.3) \qquad \min := \begin{cases} L & \text{if } p_L \le p_R, \\ R & \text{if } p_L > p_R, \end{cases} \qquad \max := \begin{cases} R & \text{if } p_L \le p_R, \\ L & \text{if } p_L > p_R. \end{cases}$$

Notice that $p_{\min} = \min(p_L, p_R)$, $p_{\max} = \max(p_L, p_R)$. For instance $a_{\min} = a_Z$ and $\gamma_{\min} = \gamma_Z$ if $p_{\min} = p_Z$, and $a_{\max} = a_Z$ and $\gamma_{\max} = \gamma_Z$ if $p_{\max} = p_Z$. We also introduce the two indices $m \in \{L, R\}$ and $M \in \{L, R\}$ defined as follows:

$$(5.4) \qquad m := \begin{cases} L & \text{if } \gamma_L \le \gamma_R, \\ R & \text{if } \gamma_L > \gamma_R, \end{cases} \qquad M := \begin{cases} R & \text{if } \gamma_L \le \gamma_R, \\ L & \text{if } \gamma_L > \gamma_R. \end{cases}$$

Notice that $\gamma_m = \min(\gamma_L, \gamma_R)$ and $\gamma_M := \max(\gamma_L, \gamma_R)$. However, $\gamma_{\min}$ and $\gamma_{\max}$ may not coincide with the values $\gamma_m$ and $\gamma_M$, respectively. We now propose an upper bound on $p^*$ based on the signs of $\phi(p_{\min})$ and $\phi(p_{\max})$.

**5.2. Case 0: vacuum.** If the vacuum condition holds, i.e., $v_R - v_L \ge \frac{2a_L(1-b\rho_L)}{\gamma_L - 1} + \frac{2a_R(1-b\rho_R)}{\gamma_R - 1}$, we have $p^* = 0$ and $\lambda_{\max}(0) = \max(|v_L - a_L|, |v_R + a_R|)$.

**5.3. Case 1: $0 < p^*$ and $0 < \phi(p_{\mathbf{min}})$.** This case corresponds to the L-wave and the R-wave both being expansion waves. In this case $p^* < p_{\min}$, which means that we do not need to compute $p^*$ as we have $\lambda_1^-(p^*) = v_L - a_L$ and $\lambda_3^+(p^*) = v_R + a_R$. But, if for some reason an upper bound for $p^*$ is needed, one can use the root of the function

(5.5) $$\widehat{\phi}_{RR}(p) := \alpha_R\left(\left(\frac{p}{p_R}\right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1\right) + \alpha_L\left(\left(\frac{p}{p_L}\right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1\right) + v_R - v_L.$$

Note that $\widehat{\phi}_{RR}(p) \le \phi_{RR}(p) = \phi(p)$ for all $p \in [0, p_{\min}]$. We give the root for completeness,

(5.6) $$\widehat{p}^* = \left(\frac{\alpha_R + \alpha_L - (v_R - v_L)}{\alpha_R p_R^{-\frac{\gamma_M - 1}{2\gamma_M}} + \alpha_L p_L^{-\frac{\gamma_M - 1}{2\gamma_M}}}\right)^{\frac{2\gamma_M}{\gamma_M - 1}}.$$

We have that $p^* = \widetilde{p}^* \le \widehat{p}^*$. In conclusion, an upper bound on $p^*$ is $\min(p_{\min}, \widehat{p}^*)$. This implies that $0 < p^* \le \min(p_{\min}, \widehat{p}^*)$. Notice in passing that $\lambda_1^-(\min(p_{\min}, \widehat{p}^*)) = v_L - a_L$ and $\lambda_3^+(\min(p_{\min}, \widehat{p}^*)) = v_R + a_R$.

**5.4. Case 2: $\phi(p_{\mathbf{min}}) < 0 < \phi(p_{\mathbf{max}})$.** In this case the min-wave is a shock and the max-wave is an expansion. Here we have $p_{\min} < p^* < p_{\max}$ and so for $p \in (p_{\min}, p_{\max})$ we have that

(5.7) $$\phi_{RR}(p) = \alpha_{\min}\left(\left(\frac{p}{p_{\min}}\right)^{\frac{\gamma_{\min} - 1}{2\gamma_{\min}}} - 1\right) + \alpha_{\max}\left(\left(\frac{p}{p_{\max}}\right)^{\frac{\gamma_{\max} - 1}{2\gamma_{\max}}} - 1\right) + v_R - v_L.$$

We consider two cases to derive a lower bound on $\phi_{RR}(p)$. If $\gamma_{\min} = \gamma_m$, we define

$$\widehat{\phi}_1(p) := \alpha_{\min}\left(\left(\frac{p}{p_{\min}}\right)^{\frac{\gamma_M - 1}{2\gamma_M}} r - 1\right) + \alpha_{\max}\left(\left(\frac{p}{p_{\max}}\right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1\right) + v_R - v_L,$$

$$\widehat{\phi}_2(p) := \alpha_{\min}\left(\left(\frac{p}{p_{\min}}\right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1\right) + \alpha_{\max}\left(\left(\frac{p}{p_{\max}}\right)^{\frac{\gamma_m - 1}{2\gamma_m}} r - 1\right) + v_R - v_L,$$

where $r := \left(\frac{p_{\min}}{p_{\max}}\right)^{\frac{\gamma_M - \gamma_m}{2\gamma_m\gamma_M}}$. We have $\max(\widehat{\phi}_1(p), \widehat{\phi}_2(p)) \le \phi_{RR}(p)$ for all $p \in (p_{\min}, p_{\max})$. Solving $\widehat{\phi}_1(p) = 0$ and $\widehat{\phi}_2(p) = 0$ gives

$$\widehat{p}_1^* = \left(\frac{\alpha_{\min} + \alpha_{\max} - (v_R - v_L)}{r\alpha_{\min}p_{\min}^{-\frac{\gamma_M - 1}{2\gamma_M}} + \alpha_{\max}p_{\max}^{-\frac{\gamma_M - 1}{2\gamma_M}}}\right)^{\frac{2\gamma_M}{\gamma_M - 1}}, \qquad \widehat{p}_2^* = \left(\frac{\alpha_{\min} + \alpha_{\max} - (v_R - v_L)}{\alpha_{\min}p_{\min}^{-\frac{\gamma_m - 1}{2\gamma_m}} + r\alpha_{\max}p_{\max}^{-\frac{\gamma_m - 1}{2\gamma_m}}}\right)^{\frac{2\gamma_m}{\gamma_m - 1}}.$$

Hence, an upper bound on $p^*$ is $\min(p_{\max}, \widehat{p}_1^*, \widehat{p}_2^*)$ if $\gamma_{\min} = \gamma_m$. This implies that $p_{\min} < p^* \le \min(p_{\max}, \widehat{p}_1^*, \widehat{p}_2^*)$. In the other case, $\gamma_{\min} = \gamma_M$, we have $\gamma_{\max} = \gamma_m$ and two lower bounds on $\check{\phi}(p)$ are given by

$$\widehat{\phi}_1(p) := \alpha_{\min}\left(\left(\frac{p}{p_{\min}}\right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1\right) + \alpha_{\max}\left(\left(\frac{p}{p_{\max}}\right)^{\frac{\gamma_m - 1}{2\gamma_m}} - 1\right) + v_R - v_L,$$

$$\widehat{\phi}_2(p) := \alpha_{\min}\left(\left(\frac{p}{p_{\min}}\right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1\right) + \alpha_{\max}\left(\left(\frac{p}{p_{\max}}\right)^{\frac{\gamma_M - 1}{2\gamma_M}} - 1\right) + v_R - v_L.$$

Again, the equations $\widehat{\phi}_1(p) = 0$, $\widehat{\phi}_2(p) = 0$ are linear (up to a change of variable). The roots are

$$\widehat{p}_1^* = \left( \frac{\alpha_{\min} + \alpha_{\max} - (v_R - v_L)}{\alpha_{\min} p_{\min}^{-\frac{\gamma_m-1}{2\gamma_m}} + \alpha_{\max} p_{\max}^{-\frac{\gamma_m-1}{2\gamma_m}}} \right)^{\frac{2\gamma_m}{\gamma_m-1}}, \qquad \widehat{p}_2^* = \left( \frac{\alpha_{\min} + \alpha_{\max} - (v_R - v_L)}{\alpha_{\min} p_{\min}^{-\frac{\gamma_M-1}{2\gamma_M}} + \alpha_{\max} p_{\max}^{-\frac{\gamma_M-1}{2\gamma_M}}} \right)^{\frac{2\gamma_M}{\gamma_M-1}}.$$

An upper bound on $p^*$ is $\min(p_{\max}, \widehat{p}_1^*, \widehat{p}_2^*)$ if $\gamma_{\min} = \gamma_M$. Hence $p_{\min} < p^* \le \min(p_{\max}, \widehat{p}_1^*, \widehat{p}_2^*)$.

**5.5. Case 3: $\phi(p_{\mathbf{max}}) < 0$.** In this case we have $p_{\max} < p^*$ and the L-wave and the R-wave are shocks. We bound $\phi_{RR}(p)$ from below by the function,

$$(5.8) \qquad \widehat{\phi}(p) := \alpha_L \left( \left( \frac{p}{p_L} \right)^{\frac{\gamma_m-1}{2\gamma_m}} - 1 \right) + \alpha_R \left( \left( \frac{p}{p_R} \right)^{\frac{\gamma_m-1}{2\gamma_m}} - 1 \right) + v_R - v_L.$$

The corresponding root for $\widehat{\phi}(p) = 0$ is

$$(5.9) \qquad \widehat{p}_1^* = \left( \frac{\alpha_L + \alpha_R - (v_R - v_L)}{\alpha_L p_L^{-\frac{\gamma_m-1}{2\gamma_m}} + \alpha_R p_R^{-\frac{\gamma_m-1}{2\gamma_m}}} \right)^{\frac{2\gamma_m}{\gamma_m-1}}.$$

Another possibility consists of observing that $\phi$ is the sum of two shock curves plus the constant $v_R - v_L$. Observing that $B_Z \le B_Z p\, p_{\max}^{-1}$ for all $p \in (p_{\max}, \infty)$, we infer that the graph of the following function is also below the graph of $\phi$:

$$(5.10) \qquad \widehat{\phi}(p) := \frac{p - p_L}{\sqrt{p}} \left( \frac{A_L}{1 + \frac{B_L}{p_{\max}}} \right)^{\frac{1}{2}} + \frac{p - p_R}{\sqrt{p}} \left( \frac{A_R}{1 + \frac{B_R}{p_{\max}}} \right)^{\frac{1}{2}} + v_R - v_L.$$

Let $x_Z := \left( \frac{A_Z}{1 + B_Z p_{\max}^{-1}} \right)^{\frac{1}{2}}$, $a := x_L + x_R$, $b := v_R - v_L$, $c := -p_L x_L - p_R x_R$, then the only positive root of $\widehat{\phi}$ is

$$(5.11) \qquad \widehat{p}_2^* = \left( \frac{-b + (b^2 - 4ac)^{\frac{1}{2}}}{2a} \right)^2.$$

An upper bound on $p^*$ is $\min(\widehat{p}_1^*, \widehat{p}_2^*)$. Hence $p_{\max} < p^* \le \min(\widehat{p}_1^*, \widehat{p}_2^*)$.

**5.6. Iterative solution.** Another possibility to estimate $p^*$ from above consists of solving $\phi(p) = 0$ by using the iterative quadratic Newton method described in Guermond and Popov [11, Alg. 1]. The method is guaranteed to be convergent since the function $\phi$ defined in (4.2) is concave. Using the lower and upper bounds provided in §5.3–§5.5, the method is also guaranteed to deliver an upper bound on $p^*$ for every termination threshold since $\phi'''(\xi) > 0$ for all $\xi > 0$ (see the proof of Lemma 4.5 in [11]). A source code for this method is publicly available at [3].

**6. Numerical Results.** We numerically illustrate in this section the algorithm (2.7) with the viscosity defined in Theorem 4.6 using the explicit upper bound $\widehat{p}^*$ defined in §5.2–5.5.

**6.1. Convergence tests.** We use the van der Waals equation of state as the oracle to validate the method. More precisely, we consider the solution to a Riemann problem and compare it to the numerical approximation (2.7) where the viscosity $d_{ij}^n$ is defined in (4.14b) with $\widehat{p}^*$ being the upper

534  bound on $p^*$ derived in §5.2–§5.5. Recall that for the van der Waals equation of state, the pressure
535  is given by $p(\rho, e) := (\gamma - 1)\frac{\rho e + a\rho^2}{1 - b\rho} - a\rho^2$, where $\gamma$, $a$ and $b$ are constants depending on the nature
536  of the fluid (see e.g., Callen [2, §3.5], Fossati and Quartapelle [9, §6.3]). We select the parameters $\gamma$,
537  $a$, $b$ so that the problem is hyperbolic and the solution exhibits a composite wave structure: we use
538  $\gamma = 1.02$, $a = 1$, $b = 1$. With these parameters the isentropes in the $(p, \frac{1}{\rho})$ diagram are nonconvex.
539  The loss of convexity is necessary for the existence of composite waves. The initial left and right
540  states we choose are:

541  (6.1)
$$(\rho_L, v_L, p_L) := (0.10, -0.475504638574729, 0.022084258693080),$$
$$(\rho_R, v_R, p_R) := (0.39, -0.121375781741349, 0.039073167077590).$$

542  The exact solution is a 3-wave composed of an expansion fan, a shock, and another expansion fan.
543  The details of the construction of the solution can be found in Cramer and Sen [6], Lai [17], and
544  Fossati and Quartapelle [9, §6.4]. For completeness and reproducibility, the construction of the
545  exact solution is given in the supplementary material and a code computing the exact solution is
546  available at Clayton et al. [4].

| #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 1: Consolidated errors and convergence rates. Solution computed at $t = 5.0$.

547  We approximate the solution with $\mathbb{P}_1$ continuous finite elements in one dimension. The com-
548  putational domain is $D := (-1, 1)$ with CFL=0.5. The estimation of the maximum wave speed
549  (see (4.14a)) is done by using $\widehat{p}^*$ as explained in §5.2–§5.5. A series of computations is done on nested
550  uniform meshes to estimate the convergence rate of the method. Denoting by $(\rho_h(t), \boldsymbol{m}_h(t), E_h(t))$
551  the approximation at time $t$, we compute a consolidated error indicator by adding the relative error
552  in the $L^q$-norm on the density, the momentum, and the total energy as follows:

553  (6.2)
$$\delta_q(t) := \frac{\|\rho_h(t) - \rho(t)\|_{L^q(D)}}{\|\rho(t)\|_{L^q(D)}} + \frac{\|\boldsymbol{m}_h(t) - \boldsymbol{m}(t)\|_{\boldsymbol{L}^q(D)}}{\|\boldsymbol{m}(t)\|_{\boldsymbol{L}^q(D)}} + \frac{\|E_h(t) - E(t)\|_{L^q(D)}}{\|E(t)\|_{L^q(D)}}.$$
554

555  The results of the convergence tests are reported in Table 1. The number of grid points is reported
556  in the leftmost column. The errors are computed at $t = 0.5$. We observe that the method is
557  convergent, and the convergence rates are consistent with the approximation being formally first-
558  order accurate.

559  **6.2. The two-expansion-wave-speed estimate.** It is often reported in the literature that,
560  for practical purpose, one can use the two expansion wave speeds, $v_L - c_L$, $v_R + c_R$, to estimate

561  the maximum wave speed. Using the covolume equation of state, we have shown in [11, App. B]
562  that $\max(|v_L - c_L|, |v_R + c_R|)$ is not an upper bound on the maximum wave speed in the Riemann
563  problem. But the reader could legitimately be skeptical about this kind of theoretical result and
564  may wonder whether these academic arguments have any impact on practical computations. We
565  now illustrate that the two-expansion-wave-speed estimate is not robust: it can either lead to an
566  underestimation or to an overestimation of the viscosity with severe consequences in both cases.



Fig. 1: Test with the data (6.3), $t = 1.25$. From left to right: density, pressure, sound speed.



Fig. 2: Test with the data (6.4), $t = 0.4$. From left to right: density, pressure, sound speed.

567      We start by showing that $\max(|v_L - c_L|, |v_R + c_R|)$ can lead to an underestimation of the
568  viscosity and therefore lead to violations of important properties. Our oracle is the van der Waals
569  equation of state with $a = 1$, $b = 1$, $\gamma = 1.02$. We solve two Riemann problems. The first one is
570  equipped with the following data set:

571  (6.3)
$$(\rho_L, v_L, p_L) := (0.2450, 0, 2.9123894332846005 \times 10^{-2}),$$
$$(\rho_R, v_R, p_R) := (0.1225, 0, 2.0685894810791836 \times 10^{-2}),$$

572  which gives the sound speeds $(c_L, c_R) \approx (0.00399, 0.306)$. The second one is equipped with the

following data set:

$$
\begin{aligned}
(\rho_L, v_L, p_L) &:= (2.5\times 10^{-1}, 0, 3\times 10^{-2}), \\
(\rho_R, v_R, p_R) &:= (4.9\times 10^{-5}, 0, 5\times 10^{-8}),
\end{aligned}
\tag{6.4}
$$

which gives the sound speeds $(c_L, c_R) \approx (0.057, 0.031)$. For each data set, we perform two series of computations on the domain $D = (-0.5, 1)$. The computations are done up to $t = 1.25$ for the first data set and up to $t = 0.4$ for the second data set. In both cases we use CFL $= 0.5$. One series of computations is done with the estimation of the maximum wave speed (see (4.14a)) using $\widehat{p}^*$ as explained in §5.2–§5.5 (no iteration is done). The other one is done using the two-expansion-wave-speed estimate $\max(|v_L - c_L(p_L, \rho_L)|, |v_R + c_R(p_L, \rho_L)|)$ with $c(p, \rho) = (\gamma \frac{p + a\rho^2}{\rho(1 - b\rho)} - 2a\rho)^{\frac{1}{2}}$. It turns out that the computations done with the two-expansion-wave-speed estimate violates the invariant domain property after a few time steps for both data sets: one obtains a complex sound speed for the first data set and one obtains a negative internal energy for the second data set. These violations occur no matter how small the CFL number is. The computations done with the method proposed in the paper run without any problem. We show in Figure 1 the density, the pressure and the sound speed profiles for various mesh sizes $(\frac{1.5}{100}, \frac{1.5}{400}, \frac{1.5}{1600}, \frac{1.5}{25600})$ for the data set (6.3). The results for the second data set (6.4) are shown in Figure 2 with the mesh sizes $\frac{1.5}{1600}, \frac{1.5}{6400}$. Notice that in both cases the R-wave is a composite wave composed of an expansion followed by a shock.
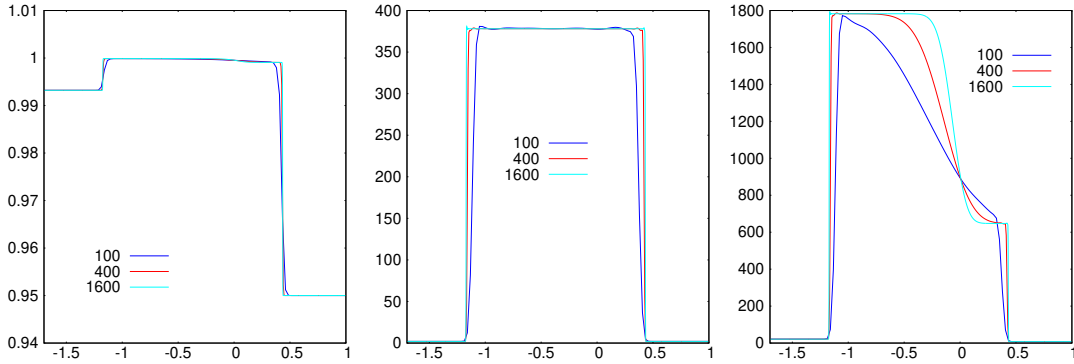


Fig. 3: Test with the data (6.5), $t = 0.005$. From left to right: density, pressure, sound speed.

We now show that the two-expansion-wave-speed estimate can lead to a local overerestimation of the viscosity and thereby to a reduction of the admissible range of time step sizes. We use again the van der Waals equation of state with the same parameters as above for the oracle. We consider the Riemann problem with the following data:

$$
\begin{aligned}
(\rho_L, v_L, p_L) &:= (0.9932, \ 3, 2), \\
(\rho_R, v_R, p_R) &:= (0.9500, -3, 2).
\end{aligned}
\tag{6.5}
$$

The corresponding sound speeds are $(c_L, c_R) \approx (21.2, 7.77)$. The computational domain is $D = (-1.7, 1)$ and the computations are done up to $t = 0.005$. For the computation with the two-expansion-wave-speed, the CFL number needed to avoid producing negative internal energy is about 0.06. The maximal admissible CFL number for the present method is about 0.71 (i.e., below

this CFL number the sound speed is real and the internal energy is positive at every grid point and for every time step). As a result the computational cost of the method using the two-expansion-wave-speed estimate is almost 12 times higher than that of the present method. We show in Figure 3 the density, the pressure and the sound speed for various meshes using the present method. The results obtained with the two-expansion-wave-speed estimate are almost identical (not shown).

**6.3. Further illustrations.** We continue by illustrating the proposed method by using a cubic equation of state as the oracle, see Redlich and Kwong [25], Valderrama [29]. We refer the reader to Dumbser and Casulli [8] where series of tests are done with this type of equation of state. For a general cubic equation of state, the pressure is given by

(6.6) $$p(\rho, e) = \frac{R\rho T(\rho, e)}{1 - b\rho} - \frac{\alpha \rho^2}{\sqrt{T(\rho, e)}(1 - br_1\rho)(1 - br_2\rho)},$$

where $T(\rho, e)$ solves the following cubic equation:

(6.7) $$e = c_v T + \frac{3\alpha}{2b\sqrt{T}} \frac{1}{r_1 - r_2} \log\left(\frac{1 - br_1\rho}{1 - br_2\rho}\right).$$
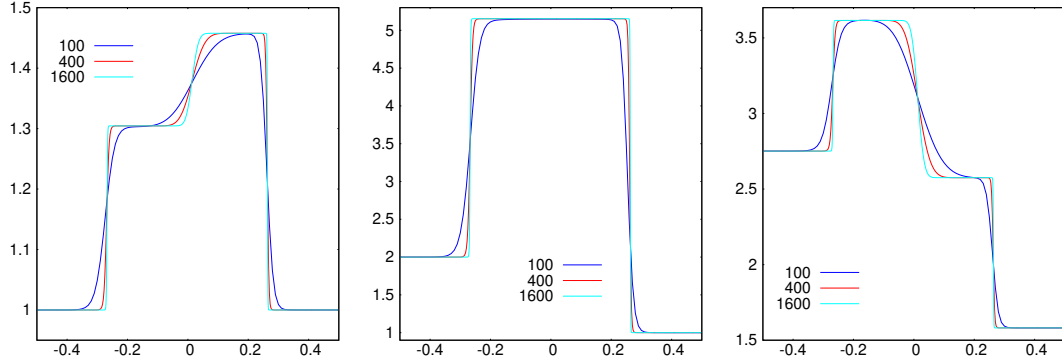


Fig. 4: Test with the data (6.8), $t = 0.1$. From left to right: density, pressure, temperature.

We take $r_1 = 0$ and $r_2 = -1$ (this corresponds to the so-called Redlich-Kwong equation). We solve two of the problems from [8, §3.3] where $R = 0.4$, $\alpha = 0.5$, $b = 0.5$. These are two Riemann problems. For the first problem we take $c_v = 1$ and the initial data are

(6.8) $$\begin{aligned} (\rho_L, v_L, p_L) &:= (1, 1, 2), \\ (\rho_R, v_R, p_R) &:= (1, -1, 1). \end{aligned}$$

The computational domain is $(-0.5, 0.5)$ and the final time is $t = 0.1$. For the second problem we take

(6.9) $$\begin{aligned} (\rho_L, v_L, p_L) &:= (1, 0, 1000), \\ (\rho_R, v_R, p_R) &:= (1, 0, 0.01), \end{aligned}$$

with $c_v = 1.5$ (we suspect there is a typo in [8, §3.3], since the authors say that they use $c_v = 1$ with the above data, but this gives a negative internal energy for the right state.) The computational
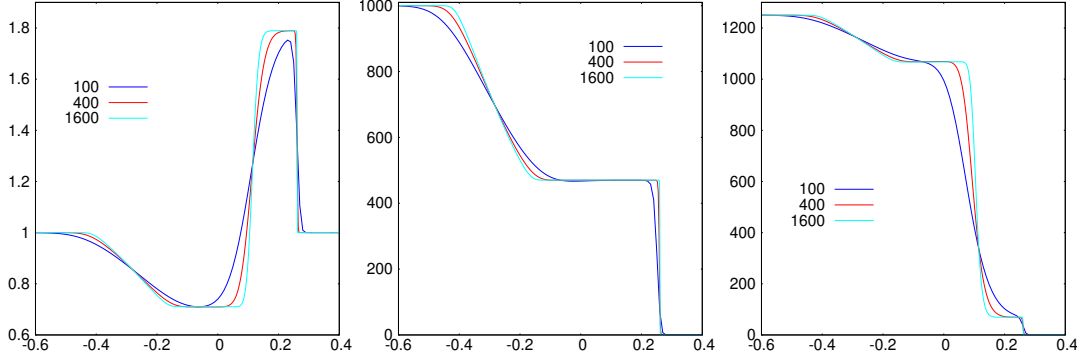
Fig. 5: Test with the data (6.9), $t = 0.008$. From left to right: density, pressure, temperature.

619    domain is $D = (-0.6, 0.4)$ and the final time is $t = 0.008$. In both cases, we take the covolume
620    constant in (3.2) to be $b = 0.5$ (using $b = 0$ in (3.2) gives similar results, not shown). The CFL
621    number is 0.5. The results obtained with various meshes are displayed in Figure 4, for the first
622    case, and in Figure 5, for the second case. In each case, we show the density, the pressure and the
623    temperature. These results are similar to those reported in [8, §3.3].

624    **6.4. Two-dimensional illustration.** To demonstrate that the proposed method is actually
625    independent of the space dimension, we illustrate it by using a finite element code which implements
626    the algorithm (2.7). The documentation of this program is found in Maier and Tomas [21]. We
627    replace the estimation of $\widehat{\lambda}(\boldsymbol{n}_{ij}, \mathbf{U}_i, \mathbf{U}_j)$ used in this code (and described in [11]) by the estimation
628    (4.14a) with $\widehat{p}^*$ computed as explained in §5.2–§5.5. The oracle is the van der Waals equation of
629    state with $\gamma = 1.4$, $a = 0.3215$, and $b = 0.1$. The computation of $\widehat{p}^*$ is done with the assumption
630    that $b = 0$. That is, we assume that the covolume constant $b$ is unknown.
631    We simulate the flow around a cylinder in a two-dimensional channel. The computational
632    domain is $D = (-0.9, 3.1) \times (-1, 1) \backslash C$, with $C$ being the disk of radius 0.15 centered at $(0, 0)$.
633    We enforce the density, the momentum and the total energy at the inflow boundary, $\{x = -0.9\}$:
634    $(\rho, \boldsymbol{m}, E) = (1.4, (4.2, 0)^{\mathsf{T}}, 9.154375)$. The primitive variable corresponding to these data are $\boldsymbol{v} =$
635    $(3, 0)^{\mathsf{T}}$ and $p = 1$. The corresponding Mach number is 3. The slip boundary condition is enforced at
636    the top and at the bottom of the channel. Nothing is done at the outflow boundary condition (this
637    is a supersonic outflow boundary). We use continuous $\mathbb{Q}_1$ finite elements. We refer the reader to
638    Maier and Tomas [21] for the implementation details. We show in Figure 6 the density computed
639    at time $t = 4$ using a Schlieren-like representation. Letting $\sum_{i \in \mathcal{V}} \rho_i^n \varphi_i$ be the approximation of
640    the density, we approximate the Euclidean norm of the gradient of the density as follows $r_i^n :=$
641    $m_i^{-1} \| \sum_{j \in \mathcal{I}(D_i)} \boldsymbol{c}_{ij} \rho_j^n \|_{\ell^2}$, for all $i \in \mathcal{V}$. The values of the Schlieren field are defined at the grid
642    points by $\exp(-\beta(r_i^n - \min_{j \in \mathcal{I}(i)} r_j^n)/(\max_{j \in \mathcal{I}(i)} r_j^n - \min_{j \in \mathcal{I}(i)} r_j^n))$ where $\beta = 10$. For comparison,
643    we also show in the right panel of this figure the density obtained at the same time using the ideal
644    gas equation of state. The inflow boundary data is $(\rho, \boldsymbol{m}, E) := (1.4, (4.2, 0)^{\mathsf{T}}, 8.8)$ and $\gamma = 1.4$.
645    This corresponds to the same primitive state, $\boldsymbol{v} = (3, 0)^{\mathsf{T}}$ and $p = 1$, as the simulation with the van
646    der Waals equation of state. The mesh used for these computations has $1.4 \times 10^6$ grid points.
647    Of course, these simulations are first-order accurate in space. Making the approximation higher-
648    order accurate can be done by implementing the convex limiting technique described in [13, 14].
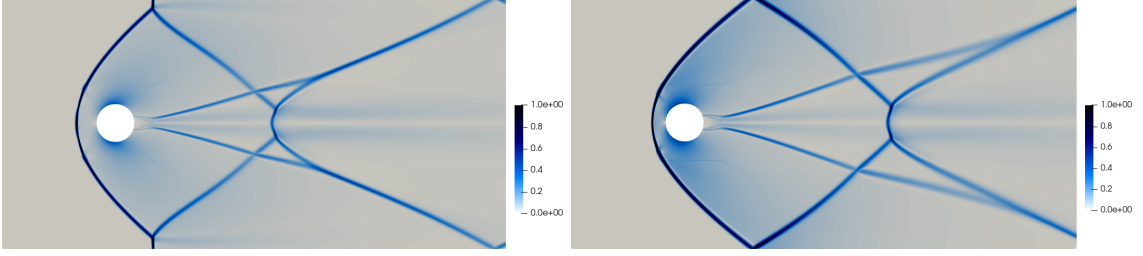
Fig. 6: Cylinder at Mach 3 in a channel. Density, $t = 4$. Left: the oracle is the van der Waals equation of state. Right: the oracle is the ideal gas equation of state with $\gamma = 1.4$.

This however requires developing surrogate entropies functionals for the oracle. This task is under way and the results of this work will be reported elsewhere. We are currently implementing the technique in the massively parallel code documented in Maier and Kronbichler [20].

**7. Conclusions.** We have proposed in the paper an approximation technique for the compressible Euler equations where the equation of state is given by an oracle. The key feature is an artificial graph viscosity using an estimate on the maximum wave speed on each elementary Riemann problem that guarantees the positivity of the density and of the internal energy. This estimate also guarantees an upper bound on the density when a covolume constant in known. The main theoretical result of the paper is Theorem 4.6. The guaranteed bounds developed in §5.2–§5.5 are easy to compute. These upper bounds can be used in any algorithm that is based on approximate Riemann solvers. A computer code implementing all these bounds is freely available at Clayton et al. [3]. All the simulations reported in the paper have been done with this code.

**Acknowledgments.** The authors thank Matthias Maier and Eric Tovar for stimulating discussions and the help they provided at various stages of this project.

**Appendix A. Improvement on the $\gamma > \frac{5}{3}$ estimates.** The objective of this a appendix is to prove that $\phi_{RR}(p) \leq \phi(p)$ for all $p \in [\min(p_L, p_R), \infty)$, where we recall that the function $\phi$ is defined in (4.2), the function $\phi_{RR}$ is defined in (5.1). For future reference we also recall that

(A.1) $$f_Z^S(p) := (p - p_Z)\sqrt{\frac{2}{(\gamma_Z + 1)\rho_Z}}\left(p + \frac{\gamma_Z - 1}{\gamma_Z + 1}p_Z\right)^{-\frac{1}{2}}\sqrt{1 - b\rho_Z},$$

(A.2) $$f_Z^R(p) := \frac{2\sqrt{\frac{\gamma_Z p_Z}{\rho_Z}}}{\gamma_Z - 1}\left(\left(\frac{p}{p_Z}\right)^{\frac{\gamma_Z - 1}{2\gamma_Z}} - 1\right)\sqrt{1 - b\rho_Z}.$$

The functions $f_Z^S(p)$ and $f_Z^R$ are, respectively, the shock and rarefaction curves introduced in (4.1). The following lemma is one of the main result established in Guermond and Popov [11]:

LEMMA A.1 ([11, Lem. 4.2]). *Let $p_Z > 0$, $\rho_Z$ be such that $0 < 1 - b\rho_Z < 1$, and $\gamma_Z \in (1, \infty)$. Assume that $\gamma \in (1, \frac{5}{3}]$. Then $f_R(p) < f_S(p)$ for all $p \in (p_Z, \infty)$ and $f_R(p_Z) = f_S(p_Z)$, i.e., the shock curve is above the rarefaction curve.*

THEOREM A.2. *Assume $\gamma \in (1, \frac{5}{3}]$. Let $p_{\min}$ and $p_{\max}$ be defined as in §5.1. For any $p \geq 0$,*

675  *the graph of $\phi(p)$ is above the graph of $\phi_{RR}(p)$; more precisely, $\phi_{RR}(p) = \phi(p)$ for all $p \in [0, p_{\min}]$*
676  *and $\phi_{RR}(p) < \phi(p)$ for all $p \in (p_{\min}, \infty)$.*

677  *Proof.* Note that the two curves $(p, \phi(p))$ and $(p, \phi_{RR}(p))$ coincide if $p \leq p_{\min}$ because both $\phi$
678  and $\phi_{RR}$ are the sum of the two rarefaction curves plus the constant $v_R - v_L$. If $p_{\min} < p \leq p_{\max}$
679  the function $\phi(p)$ is the sum of one rarefaction curve and one shock curve plus the constant $v_R - v_L$.
680  We then conclude by invoking Lemma A.1 with $(p_Z, \rho_Z) = (p_{\min}, \rho_{\min})$. If $p_{\max} \leq p$ the function
681  $\phi(p)$ is the sum of two shock curves plus the constant $v_R - v_L$. Now we invoke Lemma A.1 twice
682  to complete the proof, once with $(p_Z, \rho_Z) = (p_{\min}, \rho_{\min})$ and once with $(p_Z, \rho_Z) = (p_{\max}, \rho_{\max})$. $\square$

683  The assertion in Lemma A.1 is false when $\frac{5}{3} < \gamma_Z$. To remedy this deficiency, we now define a
684  new function that is guaranteed to be always under $\phi(p)$ for all $\gamma_Z \in (1, \infty)$ and all $p \in (p_{\min}, \infty)$.
685  Consider

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

687  Notice that $(1, \infty) \ni \gamma_z \mapsto c(\gamma_Z)$ is continuous and $c(\gamma_Z) \in (\frac{1}{2}, 1]$.

688  LEMMA A.3. *Let $p_Z > 0$, $\rho_Z$ be such that $0 < 1 - b\rho_Z < 1$, and $\gamma_Z \in (1, \infty)$. Then*
689  $c(\gamma_Z)f_Z^R(p_Z) = f_Z^S(p_Z) = 0$ *and* $c(\gamma_Z)f_Z^R(p) < f_Z^S(p)$ *for all $p \in (p_Z, \infty)$.*

690  *Proof.* The proof of the assertion is in the supplementary material. $\square$

691  **References.**

692  [1] R. Abgrall and S. Karni. Computations of compressible multifluids. *J. Comput. Phys.*, 169(2):
693      594–623, 2001.
694  [2] H. Callen. *Thermodynamics and an introduction to thermostatistics, 2nd ed.* John Wiley &
695      Sons, New York, 1985. Second edition.
696  [3] B. Clayton, J.-L. Guermond, and B. Popov. Upper bound on the maximum wave speed in
697      Riemann problems for the Euler equations with tabulated equation of state, apr 2021. URL
698      https://doi.org/10.5281/zenodo.4685868.
699  [4] B. Clayton, J.-L. Guermond, and B. Popov. Solution to the double sonic shock Riemann
700      problem with the van der Waals equation of state, Apr. 2021. URL https://doi.org/10.5281/
701      zenodo.4685958.
702  [5] P. Colella and H. M. Glaz. Efficient solution algorithms for the Riemann problem for real
703      gases. *J. Comput. Phys.*, 59(2):264–289, 1985.
704  [6] M. S. Cramer and R. Sen. Exact solutions for sonic shocks in van der waals gases. *The Physics
705      of Fluids*, 30(2):377–385, 1987.
706  [7] J. K. Dukowicz. A general, noniterative Riemann solver for Godunov's method. *J. Comput.
707      Phys.*, 61(1):119–137, 1985.
708  [8] M. Dumbser and V. Casulli. A conservative, weakly nonlinear semi-implicit finite volume
709      scheme for the compressible Navier-Stokes equations with general equation of state. *Appl.
710      Math. Comput.*, 272(part 2):479–497, 2016.
711  [9] M. Fossati and L. Quartapelle. The riemann problem for hyperbolic equations under a non-
712      convex flux with two inflection points, 2014. URL https://arxiv.org/abs/1402.5906.
713  [10] E. Godlewski and P.-A. Raviart. *Numerical approximation of hyperbolic systems of conserva-
714      tion laws*, volume 118 of *Applied Mathematical Sciences*. Springer-Verlag, New York, 1996.

[11] J.-L. Guermond and B. Popov. Fast estimation from above of the maximum wave speed in the Riemann problem for the Euler equations. *J. Comput. Phys.*, 321:908–926, 2016.

[12] J.-L. Guermond and B. Popov. Invariant domains and first-order continuous finite element approximation for hyperbolic systems. *SIAM J. Numer. Anal.*, 54(4):2466–2489, 2016.

[13] J.-L. Guermond, M. Nazarov, B. Popov, and I. Tomas. Second-order invariant domain preserving approximation of the Euler equations using convex limiting. *SIAM J. Sci. Comput.*, 40(5):A3211–A3239, 2018.

[14] J.-L. Guermond, B. Popov, and I. Tomas. Invariant domain preserving discretization-independent schemes and convex limiting for hyperbolic systems. *Comput. Methods Appl. Mech. Engrg.*, 347:143–175, 2019.

[15] H. Holden and N. H. Risebro. *Front tracking for hyperbolic conservation laws*, volume 152 of *Applied Mathematical Sciences*. Springer, Heidelberg, second edition, 2015.

[16] M. J. Ivings, D. M. Causon, and E. F. Toro. On Riemann solvers for compressible liquids. *Internat. J. Numer. Methods Fluids*, 28(3):395–418, 1998.

[17] G. Lai. Interactions of composite waves of the two-dimensional full Euler equations for van der Waals gases. *SIAM J. Math. Anal.*, 50(4):3535–3597, 2018.

[18] P. D. Lax. Weak solutions of nonlinear hyperbolic equations and their numerical computation. *Comm. Pure Appl. Math.*, 7:159–193, 1954.

[19] P. D. Lax. Hyperbolic systems of conservation laws. II. *Comm. Pure Appl. Math.*, 10:537–566, 1957.

[20] M. Maier and M. Kronbichler. Efficient parallel 3d computation of the compressible euler equations with an invariant-domain preserving second-order finite-element scheme, 2021.

[21] M. Maier and I. Tomas. tamiko/step-69: step-69 v20200305, Mar. 2020. URL https://doi.org/10.5281/zenodo.3698223.

[22] C. Pantano, R. Saurel, and T. Schmitt. An oscillation free shock-capturing method for compressible van der Waals supercritical fluid flows. *J. Comput. Phys.*, 335:780–811, 2017.

[23] J. Pike. Riemann solvers for perfect and near-perfect gases. *AIAA Journal*, 31(10):1801–1808, 1993.

[24] L. Quartapelle, L. Castelletti, A. Guardone, and G. Quaranta. Solution of the Riemann problem of classical gasdynamics. *J. Comput. Phys.*, 190(1):118–140, 2003.

[25] O. Redlich and J. N. S. Kwong. On the thermodynamics of solutions. v. an equation of state. fugacities of gaseous solutions. *Chemical Reviews*, 44(1):233–244, 1949.

[26] P. L. Roe and J. Pike. Efficient construction and utilisation of approximate riemann solutions. In *Proc. of the Sixth Int'l. Symposium on Computing Methods in Applied Sciences and Engineering, VI*, pages 499–518, NLD, 1985. North-Holland Publishing Co.

[27] E. F. Toro. *Riemann solvers and numerical methods for fluid dynamics*. Springer-Verlag, Berlin, third edition, 2009.

[28] E. F. Toro, C. E. Castro, and B. J. Lee. A novel numerical flux for the 3D Euler equations with general equation of state. *J. Comput. Phys.*, 303:80–94, 2015.

[29] J. O. Valderrama. The state of the cubic equations of state. *Industrial & Engineering Chemistry Research*, 42(8):1603–1618, 2003.