

# INVARIANT DOMAIN PRESERVING CENTRAL SCHEMES FOR NONLINEAR HYPERBOLIC SYSTEMS\*

BOJAN POPOV<sup>†</sup> AND YUCHEN HUA<sup>‡</sup>

**Abstract.** We propose a central scheme framework for the approximation of hyperbolic systems of conservation laws in any space dimension. The new central schemes are defined so that any convex invariant set containing the initial data can be an invariant domain for the numerical method. The underlying first-order central scheme is the analog of the guaranteed maximum speed method of [J.-L. Guermond and B. Popov, *SIAM J. Anal.*, 54(4):2466–2489, 2016] adjusted to the finite volume framework. There are three novelties in this work. The first one is that any classical second-order central scheme can be modified to satisfy an invariant domain property of the first-order scheme via a process which we call *convex limiting*. This is done by using convex flux limiting along the lines of [J.-L. Guermond, B. Popov and I. Tomas, *Comput. Meth. Appl. Mech. Engrg.*, 347:143–175, 2019]. The second novelty is the design of a new second-order method based on slope limiting only. The new local slope reconstruction technique is based on convex limiting so that the cell interface values are corrected to fit into a local invariant domain of the hyperbolic system. This new type of slope limiting depends on the hyperbolic system and to the best of our knowledge is the only one to guarantee local invariant domain preservation. Both schemes, flux and slope limiting based, are shown to be second-order accurate for smooth solutions in the  $L^\infty$ -norm and robust in all test cases. The third novelty is a new second-order method based on the MAPR limiter from [I. Christov and B. Popov, *J. Comput. Phys.*, 227(11):5736–5757, 2008] and adaptive slope limiting in the spirit of [A. Kurganov, G. Petrova and B. Popov, *SIAM J. Sci. Comput.*, 29(6):2381–2401, 2007] but based on an entropy commutator. This new method can be used as an underlying high-order method and combined with convex flux limiting to guarantee a local invariant domain property. The time stepping of all methods is done by using strong stability preserving Runge-Kutta methods and the invariant domain property is proved under a standard CFL condition.

**Keywords.** nonlinear hyperbolic systems; Riemann problem; invariant domain; second-order method; convex limiting; finite volume method; central schemes.

**AMS subject classifications.** 65M08; 65M12; 65M15; 35L65.

## 1. Introduction

We consider Godunov-type approximations of nonlinear hyperbolic systems of conservation laws. There are two types of Godunov-type schemes: upwind and central. In general, the upwind-type schemes are considered to be less diffusive and with sharper resolution than the central schemes. However, they are based on exact or approximate Riemann solvers and some of them can be very expensive, especially in the multidimensional case. Central schemes are more efficient because they require only an estimate on the local speed of propagation at each interface and their overall complexity for achieving given accuracy tends to be smaller, see for example [19, 20]. Moreover, there are no Riemann solvers or characteristic decomposition involved which makes them a universal tool for a wide variety of problems. The goal of this work is to develop new central schemes which are robust in the sense that they preserve the invariant domain properties of the underlying hyperbolic system and at the same time keep the low

---

\*Received: August 04, 2019; Accepted (in revised form): September 21, 2020. Communicated by Siddhartha Mishra.

This material is based upon work supported in part by the National Science Foundation grants DMS 1619892, by the Air Force Office of Scientific Research, USAF, under grant/contract number FA99550-12-0358, by the Army Research Office under grant/contract number W911NF-15-1-0517.

<sup>†</sup>Department of Mathematics, Texas A&M University 3368 TAMU, College Station, TX 77843, USA ([popov@math.tamu.edu](mailto:popov@math.tamu.edu)).

<sup>‡</sup>Department of Mathematics, University of Science and Technology of China, No. 96 Jinzhai Road, Hefei City, Anhui Province, 230026, China ([hyc12908@ustc.edu.cn](mailto:hyc12908@ustc.edu.cn)).

computational complexity of the central framework. The new methods being invariant domain preserving means that they are: (i) robust when the solution is close to the boundary of the admissible set of the hyperbolic system, for example vacuum states for the Euler equations; (ii) naturally reduce oscillations by imposing local convex limiting. Moreover, the methods remain second-order accurate in all numerical tests after the local convex limiting process (flux or slope limiting). A high-order method can be made invariant domain preserving by adapting the convex flux limiting process developed in the finite element framework in [11, 12] to the central scheme setup. This is not a big novelty but it has not been done in the finite volume context. We present it here with all details so that practitioners can use it and make an existing finite volume method more robust. Note that we impose local convex limiting on the numerical solution in contrast with many existing positivity preserving methods. For example, the bound preserving limiting developed by [25] for the Euler system is very different from the convex limiting used here and in [11, 12]. Namely, in contrast to [25] we check and correct the local second-order method to fit into a local convex set in phase space: we limit density, internal energy and specific entropy using local values. In the positivity preserving limitation, see (3.5) and (3.7) in [25], the corrections are only active when the local values are close to a global *numerical vacuum* state (defined to be  $10^{-13}$  for the density and pressure) of the Euler system. Away from this numerical vacuum the limiters in (3.5) and (3.7) in [25] are not active, hence some other limitation, typically limiting local variation of the solution, is used to make the method stable.

We develop two new versions of the Kurganov-Tadmor (KT) scheme, see [19], which are invariant domain preserving. The first new scheme we propose is based on a MAPR style limiter [2], an entropy commutator used for adaptive limiting, see §3.3.2, and convex flux limiting (see §4.2) to make the scheme invariant domain preserving. The MAPR style limiting allows for second-order accuracy for smooth solutions in the  $L^\infty$ -norm, adaptive limiting is essential for problems with nonconvex fluxes [21] where phase transition is present, and convex flux limiting, see Algorithm 1, guarantees robustness of the method. The entropy commutator simplifies the adaptive limiting process from [21] and makes the process universal for any hyperbolic system with an entropy. This scheme seems to have the best performance in all tests we considered. The second scheme that we propose is based on an invariant domain preserving slope limiter, see §4.3. That scheme is formally second-order accurate and there is no need to do any extra limitations, the invariant domain is preserved because of the new convex slope limiting, see Algorithm 2. This is a brand new way to do slope limiting because it depends on the hyperbolic system at hand. For scalar equations, the local invariant domain property reduces to local maximum principle preservation and the resulting methods are not new. However, when dealing with systems, the local limitations are not related in any way to maximum principle. The goal of the limiting process is to put the interface states into a local convex set in phase space and this is the only tool used to define and reduce oscillations. All local sets used for limitations are extracted from the so-called *bar* states, see (3.7), which are averages of exact solutions of “fake” Riemann problems. This is similar to the finite element approach in [11, 12]. We are unaware of any other slope limiting methods which guarantee local invariant domain property for any hyperbolic system under a standard CFL and keeps the second-order accuracy of the method. Moreover, this is the most cost efficient method we propose. The paper is organized as follows. The problem is formulated in §2. The central scheme settings are introduced in §3. The quasiconcave limitation is presented in §4. The new invariant domain preserving schemes are also detailed in this section. The main results

of this section are Lemma 4.2, Theorem 4.1, and Lemma 4.4. The performance of the proposed methods is illustrated in §5. In both cases, flux or slope convex limiting, we observe in all numerical tests that the extra cost of convex limiting is not dominating and the convex limiting process does not deteriorate second-order accuracy for smooth solutions. Moreover, both schemes are simple to implement, parameter free, robust, and can be applied without any changes for various hyperbolic systems, for example when phase transitions are present and near vacuum states. In principle, this methodology can be extended for more general second-order finite volume type schemes (which only evolve cell averages in time) on unstructured meshes. However, a general approach for arbitrary unstructured meshes will overcomplicate the presentation of the methods, the description of the local limitations and make the resulting schemes difficult to use in practice.

**2. Preliminaries**

The objective of this section is to introduce the notions of nonlinear hyperbolic systems, Riemann problem and the concept of stability for systems, i.e., the invariant domain property. The reader who is familiar with the notions of invariant domains and Riemann problems may skip this section and go directly to §3. Our definitions of invariant sets and domains are identical to the ones in [6].

**2.1. Nonlinear hyperbolic systems.** Let  $d$  and  $m$  be the space dimension and we consider the following hyperbolic systems in conservation form

$$\begin{cases} \partial_t \mathbf{u} + \nabla \cdot \mathbf{f}(\mathbf{u}) = \mathbf{0}, & \text{for } (\mathbf{x}, t) \in \mathbb{R}^d \times \mathbb{R}_+ \\ \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}) \end{cases} \tag{2.1}$$

where the variable  $\mathbf{u}$  takes value in  $\mathbb{R}^m$  and the flux  $f$  takes value in  $(\mathbb{R}^m)^d$ , and we consider  $\mathbf{u}$  as a column vector  $\mathbf{u} = (u_1, \dots, u_m)^\top$ . The flux is a matrix with entries  $f_{ij}(\mathbf{u}), 1 \leq i \leq m, 1 \leq j \leq d$  and  $\nabla \cdot \mathbf{f}$  is a column vector with entries  $(\nabla \cdot \mathbf{f})_i = \sum_{1 \leq j \leq d} \partial_{(x_j)} f_{ij}$ . For any  $\mathbf{n} = (n_1, \dots, n_d)^\top \in \mathbb{R}^d$ , we denote  $\mathbf{f}(\mathbf{u}) \cdot \mathbf{n}$  the column vector with entries  $\sum_{1 \leq j \leq d} n_j f_{ij}(\mathbf{u})$ , where  $i \in \{1 : m\}$ .

To simplify questions regarding boundary conditions, we assume that either periodic boundary conditions are enforced, or the initial data is compactly supported or constant outside a compact set. In both cases we denote by  $D$  the spatial domain where the approximation is constructed.

We assume that (2.1) is such that there is a clear notion for the solution of the one-dimensional Riemann problem. Namely, we assume that there exists an admissible set  $\mathcal{A} \subset \mathbb{R}^m$  such that the following one-dimensional Riemann problem is (uniquely) solvable

$$\partial_t \mathbf{v} + \partial_x (\mathbf{f}(\mathbf{v}) \cdot \mathbf{n}) = 0, \quad (x, t) \in \mathbb{R} \times \mathbb{R}_+, \quad \mathbf{v}(x, 0) = \begin{cases} \mathbf{v}_L & \text{if } x < 0 \\ \mathbf{v}_R & \text{if } x > 0 \end{cases} \tag{2.2}$$

for any unit vector  $\mathbf{n} \in \mathbb{R}^d$  and any Riemann pair  $(\mathbf{v}_L, \mathbf{v}_R)$  in  $\mathcal{A}^2$ . An important property of the Riemann solution is that it has finite speed of propagation  $\lambda_{\max}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R)$ . Namely, for any  $t \geq 0$  we have

$$\mathbf{v}(x, t) = \begin{cases} \mathbf{v}_L, & \text{if } x \leq -t \lambda_{\max}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R) \\ \mathbf{v}_R, & \text{if } x \geq t \lambda_{\max}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R). \end{cases} \tag{2.3}$$

We assume also that there exists a convex subset  $A$  of  $\mathcal{A}$ , which we call invariant set, such that for any Riemann pair in  $A$ , the average of the Riemann solution over the

Riemann fan  $\int_{-t\lambda_{\max}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R)}^{t\lambda_{\max}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R)} \mathbf{v}(x, t) dx$  is also in  $A$  for all  $(x, t) \in \mathbb{R} \times \mathbb{R}_+$ . The existence of such a set has been established by [3] on a very large class of hyperbolic systems.

The following elementary result is a well-known consequence of (2.3), see e.g., [6, Lem. 2.1].

LEMMA 2.1. *Let  $\mathbf{u}_L, \mathbf{u}_R \in \mathcal{A}$ , let  $\mathbf{u}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R)$  be the Riemann solution to (2.2), let  $\bar{\mathbf{u}}(t, \mathbf{n}, \mathbf{u}_L, \mathbf{u}_R) := \int_{-\frac{1}{2}}^{\frac{1}{2}} \mathbf{u}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R)(x, t) dx$  and assume that  $t\lambda_{\max}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R) \leq \frac{1}{2}$ , then*

$$\bar{\mathbf{u}}(t, \mathbf{n}, \mathbf{u}_L, \mathbf{u}_R) = \frac{1}{2}(\mathbf{u}_L + \mathbf{u}_R) - t(\mathbf{f}(\mathbf{u}_R) \cdot \mathbf{n} - \mathbf{f}(\mathbf{u}_L) \cdot \mathbf{n}). \quad (2.4)$$

**2.2. Invariant sets and domain.** Following [6], we introduce the notions of invariant sets and invariant domains. We associate invariant sets only with the solutions of Riemann problems and define invariant domains only for an approximation process.

DEFINITION 2.1 (Invariant Set). *We say that a set  $A \subset \mathcal{A} \subset \mathbb{R}^m$  is invariant for (2.1) if for any pair  $(\mathbf{u}_L, \mathbf{u}_R) \in A \times A$ , any unit vector  $\mathbf{n} \in \mathcal{S}^{d-1}(\mathbf{0}, 1)$ , and any  $t > 0$ , the average of the entropy solution of the Riemann problem (2.2) over the Riemann fan, say,  $\frac{1}{t(\lambda_m^+ - \lambda_1^-)} \int_{\lambda_1^- t}^{\lambda_m^+ t} \mathbf{u}(\mathbf{n}, \mathbf{u}_L, \mathbf{u}_R)(x, t) dx$ , remains in  $A$ .*

We now introduce the notion of invariant domain for an approximation process. Let  $\mathbf{X}_h \subset L^1(\mathbb{R}^d, \mathbb{R}^m)$  be a finite-dimensional approximation space and let  $S_h : \mathbf{X}_h \ni u_h \mapsto S_h(\mathbf{u}_h) \in \mathbf{X}_h$  be a discrete process over  $\mathbf{X}_h$ . Henceforth we abuse the language by saying that a member of  $\mathbf{X}_h$ , say  $\mathbf{u}_h$ , is in the set  $A \subset \mathbb{R}^m$  when actually we mean that  $\{\mathbf{u}_h(\mathbf{x}) | \mathbf{x} \in \mathbb{R}^d\} \subset A$ .

DEFINITION 2.2 (Invariant Domain). *A convex invariant set  $A \subset \mathcal{A} \subset \mathbb{R}^m$  is said to be an invariant domain for the process  $S_h$  if and only if for any state  $\mathbf{u}_h$  in  $A$ , the state  $S_h(\mathbf{u}_h)$  is also in  $A$ .*

For scalar conservation equations the notions of invariant sets and invariant domains are closely related to the maximum principle, see §2.2.1. In the case of nonlinear systems, the notion of maximum principle does not apply and must be replaced by the notion of invariant domain. For example, the invariant domain theory when  $m=2$  and  $d=1$  relies on the existence of global Riemann invariants, the best known examples are the hyperbolic systems of isentropic gas dynamics in Eulerian and Lagrangian form, see [6] for details. For results on general hyperbolic systems, we refer the reader to [3, 4, 15].

Here, we will illustrate the abstract notions of invariant sets and invariant domains with some examples which are used in this paper.

**2.2.1. Example 1: scalar equations.** Assume  $m=1$  and  $d$  is arbitrary, i.e., (2.1) is a scalar equation. Provided  $\mathbf{f} \in \text{Lip}(\mathbb{R}; \mathbb{R}^d)$ , any bounded interval is an admissible set for (2.1). For any Riemann data  $u_L, u_R$ , the maximum speed of propagation in (2.3) is bounded by  $\lambda_{\max}(u_L, u_R) := \|\mathbf{f} \cdot \mathbf{n}\|_{\text{Lip}(u_{\min}, u_{\max})}$  where  $u_{\min} = \min(u_L, u_R)$ ,  $u_{\max} = \max(u_L, u_R)$ . If  $\mathbf{f}$  is convex and is of class  $C^1$ , we have  $\lambda_{\max}(u_L, u_R) = \max(|\mathbf{n} \cdot \mathbf{f}'(u_L)|, |\mathbf{n} \cdot \mathbf{f}'(u_R)|)$  if  $\mathbf{n} \cdot \mathbf{f}(u_L) \leq \mathbf{n} \cdot \mathbf{f}(u_R)$  and  $\lambda_{\max}(u_L, u_R) = \mathbf{n} \cdot (\mathbf{f}(u_L) - \mathbf{f}(u_R)) / (u_L - u_R)$  otherwise. Any interval  $[a, b] \subset \mathbb{R}$  is admissible and is an invariant set for (2.1), i.e., if  $u_R, u_L \in [a, b]$ , then  $a \leq u(\mathbf{n}, u_L, u_R) \leq b$  for all times, i.e., any interval  $[a, b]$  is an invariant domain for any numerical scheme which satisfies a local maximum principle property.

**2.2.2. Example 2: p-system.** The one-dimensional motion of an isentropic gas is modeled by the so-called p-system. In Lagrangian coordinates the system is written as follows:

$$\begin{cases} \partial_t v + \partial_x u = 0, \\ \partial_t u + \partial_x p(v) = 0, \end{cases} \text{ for } (x, t) \in \mathbb{R} \times \mathbb{R}_+. \tag{2.5}$$

The dependent variables are the velocity  $u$  and the specific volume  $v$ , i.e., the reciprocal of density. The mapping  $v \mapsto p(v)$  is the pressure and is assumed to be of class  $C^2(\mathbb{R}^+; \mathbb{R})$  and to satisfy

$$p' < 0, \quad 0 < p''. \tag{2.6}$$

A typical example is the so-called gamma-law,  $p(v) = rv^\gamma$ , where  $r > 0$  and  $\gamma \geq 1$ . Using the notation  $\mathbf{u} = (v, u)^\top$ , any set  $\mathcal{A} \subset (0, \infty) \times \mathbb{R}$  is admissible. Using the notation  $d\mu := \sqrt{-p'(s)} ds$ , and assuming  $\int_1^\infty < \infty$ , the system has two families of global Riemann invariants:

$$w_1(\mathbf{u}) = u + \int_v^\infty d\mu, \quad \text{and} \quad w_2(\mathbf{u}) = u - \int_v^\infty d\mu. \tag{2.7}$$

Note that  $\int_1^\infty d\mu < \infty$  if  $\gamma > 1$ . Let  $a, b \in \mathbb{R}$ , then it can be shown that any set  $A_{ab} \in \mathbb{R}_+ \times \mathbb{R}$  of the form

$$A_{ab} := \{\mathbf{u} \in \mathbb{R}_+ \times \mathbb{R} \mid a \leq w_2(\mathbf{u}), w_1(\mathbf{u}) \leq b\} \tag{2.8}$$

is an invariant set for the system (2.5), see [15, Exp. 3.5, p. 597]. Moreover,  $A_{ab}$  is an invariant domain for the LxF scheme, [14, Thm. 2.1] and [15, Thm. 4.1], and the guaranteed maximum speed method of [6].

**2.2.3. Example 3: Euler equations.** Consider the compressible Euler equations:

$$\partial_t \mathbf{c} + \nabla \cdot (\mathbf{f}(\mathbf{c})) = \mathbf{0}, \quad \mathbf{c} = \begin{pmatrix} \rho \\ \mathbf{m} \\ E \end{pmatrix}, \quad \mathbf{f}(\mathbf{c}) = \begin{pmatrix} \mathbf{m} \\ \mathbf{m} \otimes \frac{\mathbf{m}}{\rho} + p\mathbb{I} \\ \frac{\mathbf{m}}{\rho}(E + p) \end{pmatrix}. \tag{2.9}$$

where the independent variables are the density  $\rho$ , the momentum vector field  $\mathbf{m}$  and the total energy  $E$ . The velocity vector field  $\mathbf{u}$  is defined by  $\mathbf{u} := \frac{\mathbf{m}}{\rho}$  and the internal energy density  $e$  by  $e := \rho^{-1}E - \frac{1}{2}\|\mathbf{u}\|_{\ell^2}^2$ . The quantity  $p$  is the pressure. The symbol  $\mathbb{I}$  denotes the identity matrix in  $\mathbb{R}^d$ . Let  $s := s(\rho, e)$  be the specific entropy of the system, and assume that  $-s$  is strictly convex as a function of  $\tau := 1/\rho$  and  $e$ . It is known that

$$A_r := \{(\rho, \mathbf{m}, E) \mid \rho > 0, e > 0, s \geq r\} \tag{2.10}$$

is an invariant set for the Euler system for any  $r \in \mathbb{R}$ . For example, it is shown in [4, Thm. 7 and 8] that the  $A_r$  is convex and is an invariant domain for the staggered Lax-Friedrichs scheme. The non-staggered Lax-Friedrichs scheme and the guaranteed maximum speed (GMS) first-order scheme from [6] are also invariant domain preserving.

**3. Central schemes**

In this section, we introduce the analog of the first-order invariant domain preserving scheme from [6] adopted to the central framework and the corresponding second-order invariant domain preserving central scheme. Both of these schemes are derived from the semi-discrete Kurganov-Tadmor scheme, see [19]. For more details on central schemes we refer the reader to [19,20]. Following [19], we restrict our presentation to the case of one and two space dimensions. However, it is not difficult to derive the analogous results in arbitrary space dimension on rectangular meshes. We start with the one-dimensional setup.

**3.1. First-order central scheme.** We first consider the case of one space dimension.

$$\begin{cases} \partial_t \mathbf{u} + \partial_x \mathbf{f}(\mathbf{u}) = 0, & \text{for } (x, t) \in \mathbb{R} \times \mathbb{R}_+ \\ \mathbf{u}(x, 0) = \mathbf{u}_0(x) \end{cases} \tag{3.1}$$

We assume that the space discretization is uniform. That is, we set the cell centers to be  $x_j := j\Delta x$ ,  $j \in \mathbb{Z}$ , and assume that the approximate solution  $\tilde{\mathbf{u}}^n(x) \approx \mathbf{u}(x, t^n)$  at time  $t^n$  is a piecewise constant function

$$\tilde{\mathbf{u}}^n(x) := \sum_j \mathbf{u}_j^n \mathbf{1}_{[x_{j-1/2}, x_{j+1/2}]}, \quad x_{j\pm 1/2} := x_j \pm \frac{\Delta x}{2}. \tag{3.2}$$

The values  $\mathbf{u}_j^n$ ,  $j \in \mathbb{Z}$ , are the cell averages of the approximate solution at time  $t^n$ . The time step  $\Delta t_n := t^{n+1} - t^n$  is generic, determined by a CFL condition for each  $n \geq 0$ , and we will denote  $t^{n+1} = t^n + \Delta t$  where we abuse the notation and drop the dependence on  $n$  in the time step. The invariant domain first-order scheme from [6] can be written as follows

$$\frac{\mathbf{u}_j^{n+1} - \mathbf{u}_j^n}{\Delta t} = - \frac{\mathbf{f}(\mathbf{u}_{j+1}^n) - \mathbf{f}(\mathbf{u}_{j-1}^n)}{2\Delta x} + \frac{\lambda_{i+1/2}^n}{2\Delta x} (\mathbf{u}_{j+1}^n - \mathbf{u}_j^n) - \frac{\lambda_{i-1/2}^n}{2\Delta x} (\mathbf{u}_j^n - \mathbf{u}_{j-1}^n), \tag{3.3}$$

where the quantity  $\lambda_{j+1/2}^n := \lambda_{\max}(\mathbf{u}_j^n, \mathbf{u}_{j+1}^n, \mathbf{f})$  denotes the maximum speed of propagation of the Riemann problem with left state  $\mathbf{u}_j^n$ , right state  $\mathbf{u}_{j+1}^n$  and flux  $\mathbf{f}$ , see (2.2) and (2.3) in section §2.1. As proved in [6, Thm. 4.1] the above scheme is invariant domain preserving if the following CFL condition holds for all  $n \geq 0$

$$\frac{\Delta t \max_j \lambda_{j+1/2}^n}{\Delta x} \leq \frac{1}{2}. \tag{3.4}$$

However, it is easy to verify that using forward Euler time stepping for the first-order semi-discrete central scheme in [19, Eqn. (4.8)] will result in the same discrete method. Therefore, in the one-dimensional case, the fully discrete first-order central scheme from [19] coincides with the invariant domain preserving method from [6] when the maximum speed is defined as above and Euler time stepping is used.

We now rewrite the fully discrete scheme (3.3) in flux form using the notation  $\mathbf{u}_j^{L,n+1}$  to indicate that this is the first-order method:

$$\frac{\mathbf{u}_j^{L,n+1} - \mathbf{u}_j^n}{\Delta t} = - \frac{L_{j+1/2}^n - L_{j-1/2}^n}{\Delta x} \tag{3.5}$$

where  $L_{j+1/2}^n$  is the first-order interface numerical flux

$$L_{j+1/2}^n = \frac{1}{2} (\mathbf{f}(\mathbf{u}_{j+1}^n) + \mathbf{f}(\mathbf{u}_j^n)) - \frac{1}{2} \lambda_{j+1/2}^n (\mathbf{u}_{j+1}^n - \mathbf{u}_j^n). \tag{3.6}$$

REMARK 3.1. The Euler time stepping in (3.5) can be upgraded to any strong stability preserving (SSP) Runge-Kutta (RK) scheme and the method will still be invariant domain preserving under a standard CFL-condition, for more details see the discussion in [10, §4.5] for more details and references on SSP-RK schemes.

For completeness, we recall the following result established in [6] in this setup.

THEOREM 3.1. Let  $A \subset A$  be an invariant set for (3.1) in the sense of Definition (2.2). Assume that  $A$  is convex and that for any admissible states  $\mathbf{u}_L, \mathbf{u}_R$ , the maximum speed of propagation  $\lambda_{\max}(\mathbf{u}_L, \mathbf{u}_R, \mathbf{f})$  is finite. Assume that  $\mathbf{u}_h^0 \in A$  and the CFL condition (3.4) holds. Then we have:

- (1)  $A$  is an invariant domain for the process  $\mathbf{u}_h^n \mapsto \mathbf{u}_h^{n+1}$  where  $\mathbf{u}_h^{n+1}$  is computed with the scheme (3.3) for all  $n \geq 0$ .
- (2) Given  $n \geq 0$  and  $j \in \{1: I\}$ , let  $B \subset A$  be a convex invariant set such that  $\mathbf{u}_j^n \in B$  and  $\mathbf{u}_{j\pm 1}^n \in B$ , then  $\mathbf{u}_j^{n+1} \in B$ .

REMARK 3.2. The proof of Theorem 3.1 relies on the introduction of auxiliary states, called *bar* states, given by

$$\bar{\mathbf{u}}_{j+1/2}^{n+1} := \frac{1}{2}(\mathbf{u}_j^n + \mathbf{u}_{j+1}^n) - \frac{1}{2\lambda_{j+1/2}^n}(\mathbf{f}(\mathbf{u}_{j+1}^n) - \mathbf{f}(\mathbf{u}_j^n)), \tag{3.7}$$

which under the CFL-condition (3.4) are averages of the exact solution of the Riemann problem with a left state  $\mathbf{u}_j^n$ , a right state  $\mathbf{u}_{j+1}^n$  and a flux  $\mathbf{f}$ , see Lemma 2.1. These states are naturally in the local invariant set of the problem, see [6], and are essential for the convex limiting process which we use in this paper.

Now we consider the case of two space dimensions in (2.1) with  $\mathbf{x} := (x, y)$

$$\begin{cases} \partial_t \mathbf{u} + \partial_x \mathbf{f}(\mathbf{u}) + \partial_y \mathbf{g}(\mathbf{u}) = 0, & \text{for } (\mathbf{x}, t) \in \mathbb{R}^2 \times \mathbb{R}_+, \\ \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}). \end{cases} \tag{3.8}$$

We use uniform rectangular mesh with cell centers  $(x_j, y_k) := (j\Delta x, k\Delta y)$ ,  $j, k \in \mathbb{Z}$ , and take the approximate solution  $\tilde{\mathbf{u}}^n(x, y) \approx \mathbf{u}(x, y, t^n)$  at time  $t^n$  to be a piecewise-constant function

$$\tilde{\mathbf{u}}^n(x, y) := \sum_{j,k} \mathbf{u}_{j,k}^n \mathbf{1}_{[x_{j-1/2}, x_{j+1/2}] \times [y_{k-1/2}, y_{k+1/2}]}, \tag{3.9}$$

where the values  $\mathbf{u}_{j,k}^n$ ,  $j, k \in \mathbb{Z}$ , are the cell averages of the approximate solution at time  $t^n$ . The time step is determined in the same way as in the one dimensional case. Namely, we denote  $t^{n+1} = t^n + \Delta t$  where we abuse the notation and drop the dependence on  $n$  in the time step. The fully discrete first-order central scheme is given by

$$\frac{\mathbf{u}_{j,k}^{L,n+1} - \mathbf{u}_{j,k}^{L,n}}{\Delta t} = - \frac{L_{j+1/2,k}^{n,x} - L_{j-1/2,k}^{n,x}}{\Delta x} - \frac{L_{j,k+1/2}^{n,y} - L_{j,k-1/2}^{n,y}}{\Delta y}, \tag{3.10}$$

where the first-order interface fluxes are defined as follows:

$$L_{j+1/2,k}^{n,x} = \frac{1}{2}(\mathbf{f}(\mathbf{u}_{j+1,k}^n) + \mathbf{f}(\mathbf{u}_{j,k}^n)) - \frac{1}{2}\lambda_{j+1/2,k}^{n,x}(\mathbf{u}_{j+1,k}^n - \mathbf{u}_{j,k}^n), \tag{3.11a}$$

$$L_{j,k+1/2}^{n,y} = \frac{1}{2}(\mathbf{g}(\mathbf{u}_{j,k+1}^n) + \mathbf{g}(\mathbf{u}_{j,k}^n)) - \frac{1}{2}\lambda_{j,k+1/2}^{n,y}(\mathbf{u}_{j,k+1}^n - \mathbf{u}_{j,k}^n). \tag{3.11b}$$

Similar to the one-dimensional case we have an invariant domain property for this first-order method,  $\mathbf{u}^{n+1} := \mathbf{u}^{L,n+1}$ , if the local speeds are given by  $\lambda_{j+1/2,k}^{n,x} = \lambda_{\max}(\mathbf{u}_{j,k}^n, \mathbf{u}_{j+1,k}^n; \mathbf{f})$ ,  $\lambda_{j,k+1/2}^{n,y} = \lambda_{\max}(\mathbf{u}_{j,k}^n, \mathbf{u}_{j,k+1}^n; \mathbf{g})$  and the following CFL condition holds for all  $n \geq 0$

$$\max_{j,k} \left( \frac{\Delta t \lambda_{j+1/2,k}^{n,x}}{\Delta x} + \frac{\Delta t \lambda_{j,k+1/2}^{n,y}}{\Delta y} \right) \leq \frac{1}{2} \tag{3.12}$$

Namely, Theorem 3.1 holds under the CFL condition (3.12).

REMARK 3.3. Similar to the one-dimensional case, the proof of Theorem 3.1 relies on the introduction of the bar states

$$\bar{\mathbf{u}}_{j+1/2,k}^{n+1} = \frac{1}{2}(\mathbf{u}_{j,k}^n + \mathbf{u}_{j+1,k}^n) - \frac{1}{2\lambda_{j+1/2,k}}(\mathbf{f}(\mathbf{u}_{j+1,k}^n) - \mathbf{f}(\mathbf{u}_{j,k}^n)), \tag{3.13}$$

and

$$\bar{\mathbf{u}}_{j,k+1/2}^{n+1} = \frac{1}{2}(\mathbf{u}_{j,k}^n + \mathbf{u}_{j,k+1}^n) - \frac{1}{2\lambda_{j,k+1/2}}(\mathbf{f}(\mathbf{u}_{j,k+1}^n) - \mathbf{f}(\mathbf{u}_{j,k}^n)), \tag{3.14}$$

which under the CFL-condition (3.12) are averages of exact solutions of Riemann problems, therefore, these states are naturally in the local invariant set of the problem. In the general multidimensional case ( $d \geq 2$ ), Theorem 3.1 holds when the constant in the CFL condition (3.12) is  $\frac{1}{2d}$  and we call such methods guaranteed maximum speed (GMS) schemes, see [7] for more details on GMS schemes.

**3.2. Second-order central scheme.** Here we recall the second-order central scheme from [19], which we call the KT-scheme in the rest of the paper. In one space dimension case, we assume the same setup of space and time discretization as for first-order GMS-scheme, and assume the approximate solution  $\tilde{\mathbf{u}}^n = \mathbf{u}(x, t^n)$  at time  $t = t^n$  to be piecewise linear

$$\tilde{\mathbf{u}}^n := \sum_j [\mathbf{u}_j^n + (\mathbf{u}_x)_j^n (x - x_j)] \mathbf{1}_{[x_{j-1/2}, x_{j+1/2}]}, \quad x_{j\pm 1/2} := x_j \pm \frac{\Delta x}{2}, \tag{3.15}$$

where values  $\mathbf{u}_j^n$  are cell averages of approximate solutions and  $(\mathbf{u}_x)_j^n$  are approximations of exact derivatives  $\mathbf{u}_x(x_j, t^n)$ . The semi-discrete KT-scheme is given by

$$\begin{aligned} \frac{d}{dt} \mathbf{u}_j(t) = & - \frac{\mathbf{f}(\mathbf{u}_{j+1/2}^+(t)) + \mathbf{f}(\mathbf{u}_{j+1/2}^-(t)) - \mathbf{f}(\mathbf{u}_{j-1/2}^+(t)) - \mathbf{f}(\mathbf{u}_{j-1/2}^-(t))}{2\Delta x} \\ & + \frac{a_{j+1/2}(t)(\mathbf{u}_{j+1/2}^+(t) - \mathbf{u}_{j+1/2}^-(t)) - a_{j-1/2}(t)(\mathbf{u}_{j-1/2}^+(t) - \mathbf{u}_{j-1/2}^-(t))}{2\Delta x} \end{aligned} \tag{3.16}$$

where  $\mathbf{u}_{j+1/2}^+ := \mathbf{u}_{j+1}(t) - \frac{\Delta x}{2}(\mathbf{u}_x)_{j+1}(t)$ ,  $\mathbf{u}_{j+1/2}^- := \mathbf{u}_j(t) + \frac{\Delta x}{2}(\mathbf{u}_x)_j(t)$  are the interface values and  $a_{j+1/2}(t) = \lambda_{\max}(\mathbf{u}_{j+1/2}^-(t), \mathbf{u}_{j+1/2}^+(t), \mathbf{f})$  denotes maximum speed. By setting the second-order numerical flux to be

$$H_{j+1/2}^n := \frac{\mathbf{f}(\mathbf{u}_{j+1/2}^{n,+}) + \mathbf{f}(\mathbf{u}_{j+1/2}^{n,-})}{2} - \frac{a_{j+1/2}^n}{2}(\mathbf{u}_{j+1/2}^{n,+} - \mathbf{u}_{j+1/2}^{n,-}). \tag{3.17}$$

and using a forward Euler in time we obtain the fully discrete KT-scheme

$$\frac{\mathbf{u}_j^{H,n+1} - \mathbf{u}_j^n}{\Delta t} = - \frac{H_{j+1/2}^n - H_{j-1/2}^n}{\Delta x}. \tag{3.18}$$



REMARK 3.4. The KT-scheme (3.16) will reduce to its first-order form (3.3) if we set the slopes  $(\mathbf{u}_x)_j^n$  to zero. As in the first-order case, high-order time stepping is done by using SSP-RK schemes.

In the case of two space dimensions, we use the same rectangular cell and the same time discretization as for the first-order scheme. The approximate solution  $\tilde{\mathbf{u}}^n = \mathbf{u}(x, y, t^n)$  is a piecewise linear function given by

$$\tilde{\mathbf{u}}^n(x, y) := \sum_{j,k} [\mathbf{u}_{j,k}^n + (\mathbf{u}_x)_{j,k}^n(x - x_j) + (\mathbf{u}_y)_{j,k}^n(y - y_k)] \mathbf{1}_{[x_{j-1/2}, x_{j+1/2}] \times [y_{k-1/2}, y_{k+1/2}]}, \tag{3.19}$$

where  $x_{j\pm 1/2} := x_j \pm \frac{\Delta x}{2}$ ,  $y_{k\pm 1/2} := y_k \pm \frac{\Delta y}{2}$ . The values  $\mathbf{u}_{j,k}^n$  are the cell averages of the approximate solutions and  $((\mathbf{u}_x)_{j,k}^n, (\mathbf{u}_y)_{j,k}^n)$  is the approximate gradient on the cell  $[x_{j-1/2}, x_{j+1/2}] \times [y_{k-1/2}, y_{k+1/2}]$  at time  $t = t^n$ . Following [19], we set the numerical fluxes to be

$$H_{j+1/2,k}^{n,x} := \frac{f(\mathbf{u}_{j+1/2,k}^{n,+}) + f(\mathbf{u}_{j+1/2,k}^{n,-})}{2} - \frac{a_{j+1/2,k}^{n,x}}{2} (\mathbf{u}_{j+1/2,k}^{n,+} - \mathbf{u}_{j+1/2,k}^{n,-}), \tag{3.20a}$$

$$H_{j,k+1/2}^{n,y} := \frac{g(\mathbf{u}_{j,k+1/2}^{n,+}) + g(\mathbf{u}_{j,k+1/2}^{n,-})}{2} - \frac{a_{j,k+1/2}^{n,y}}{2} (\mathbf{u}_{j,k+1/2}^{n,+} - \mathbf{u}_{j,k+1/2}^{n,-}), \tag{3.20b}$$

where  $\mathbf{u}_{j+1/2,k}^{n,+} := \mathbf{u}_{j+1,k}^n - \frac{\Delta x}{2} (\mathbf{u}_x)_{j+1,k}^n$ ,  $\mathbf{u}_{j+1/2,k}^{n,-} := \mathbf{u}_{j,k}^n + \frac{\Delta x}{2} (\mathbf{u}_x)_{j,k}^n$ ,  $\mathbf{u}_{j,k+1/2}^{n,+} = \mathbf{u}_{j,k+1}^n - \frac{\Delta y}{2} (\mathbf{u}_y)_{j,k+1}^n$ ,  $\mathbf{u}_{j,k+1/2}^{n,-} = \mathbf{u}_{j,k}^n + \frac{\Delta y}{2} (\mathbf{u}_y)_{j,k}^n$  separately, and the local speeds are  $a_{j+1/2,k}^{n,x} = \lambda_{\max}(\mathbf{u}_{j+1/2,k}^{n,-}(t), \mathbf{u}_{j+1/2,k}^{n,+}(t), \mathbf{f})$ ,  $a_{j,k+1/2}^{n,y} = \lambda_{\max}(\mathbf{u}_{j,k+1/2}^{n,-}(t), \mathbf{u}_{j,k+1/2}^{n,+}(t), \mathbf{g})$ . Then, a forward Euler time step of the semi-discrete KT-scheme can be written as follows

$$\frac{\mathbf{u}_{j,k}^{h,n+1} - \mathbf{u}_{j,k}^n}{\Delta t} = - \frac{H_{j+1/2,k}^{n,x} - H_{j-1/2,k}^{n,x}}{\Delta x} - \frac{H_{j,k+1/2}^{n,y} - H_{j,k-1/2}^{n,y}}{\Delta y}. \tag{3.21}$$

**3.3. Piecewise linear reconstruction.** In order to completely describe the KT-scheme (3.16), we need to define the slope reconstructions in (3.15) and (3.19). It is well known that a nonlinear slope reconstruction is needed. A common approach is to use a nonlinear limiter and we recall here some of the widely used choices below. We present the reconstruction process in one space dimension. The multidimensional case is handled by splitting the gradient reconstruction into onedimensional steps.

**3.3.1. Examples of slope reconstructions.** We start with a classical slope reconstruction based on the so-called minmod limiter. The minmod slope limiter is given by

$$\sigma_j^m(\mathbf{u}) = (\mathbf{u}_x)_j := m\left(\frac{\mathbf{u}_{j+1} - \mathbf{u}_j}{\Delta x}, \frac{\mathbf{u}_j - \mathbf{u}_{j-1}}{\Delta x}\right), \tag{3.22}$$

where the minmod operator is defined as follows

$$m(x_1, x_2, \dots, x_n) = \begin{cases} \min_{1 \leq j \leq n} x_j, & \text{if } x_j > 0 \ \forall j, \\ \max_{1 \leq j \leq n} x_j, & \text{if } x_j < 0 \ \forall j, \\ 0, & \text{otherwise.} \end{cases} \tag{3.23}$$

Another minmod-type  $\theta$ -dependent family with  $1 \leq \theta \leq 2$  is given by

$$\sigma_j^{m,\theta}(\mathbf{u}) = (\mathbf{u}_x)_j := m\left(\theta \frac{\mathbf{u}_{j+1} - \mathbf{u}_j}{\Delta x}, \frac{\mathbf{u}_{j+1} - \mathbf{u}_{j-1}}{2\Delta x}, \theta \frac{\mathbf{u}_j - \mathbf{u}_{j-1}}{\Delta x}\right), \tag{3.24}$$

see for example [16, p.1900]. The range of  $\theta$  ( $1 \leq \theta \leq 2$ ) guarantees a local maximum principle for scalar equations, see [19, Cor. 5.1].

There are many other second-order reconstructions, most notably the so-called uniformly non-oscillatory (UNO) reconstruction, introduced by [13]. We refer the reader to [21] for more nonlinear reconstructions and when it is appropriate to use them.

Another gradient reconstruction, appropriate for unstructured meshes and avoiding the clipping phenomenon of minmod, was introduced in [2]. It is based on the so-called MAPR limiter given by

$$\text{mapr}(x_1, x_2, \dots, x_n) = \{x_i \mid \text{where } |x_i| = \min_{1 \leq j \leq n} |x_j|\} \tag{3.25}$$

$$\sigma_j^{\text{mapr}, \theta}(\mathbf{u}) = (\mathbf{u}_x)_j := \text{mapr}\left(\theta_j \frac{\mathbf{u}_{j+1} - \mathbf{u}_j}{\Delta x}, \frac{\mathbf{u}_{j+1} - \mathbf{u}_{j-1}}{2\Delta x}, \theta_j \frac{\mathbf{u}_j - \mathbf{u}_{j-1}}{\Delta x}\right), \tag{3.26}$$

where  $1 \leq \theta_j$  is a number specified by the user. Using  $\theta_j = 1$  gives the MAPR reconstruction from [2] and the generalized MAPR ( $1 \leq \theta_j \leq 4$ ) is the natural analog of the minmod- $\theta$  limiter. We use a variable theta in this paper,  $\theta_j^n = 2 - R_j^n$  with  $R_j^n$  the entropy commutator defined in (3.30), and apply the MAPR limiter (3.25) in the regions of non-smooth flow, where non-smooth is defined by  $\theta_j^n \leq 1.5$ . Note that, in the regions of smooth flow the entropy commutator is almost zero, see Section 3.3.2, so  $\theta_j^n \leq 1.5$  is a good cutoff for the limiter. If  $\theta_j^n > 1.5$  we use the central slope, i.e., we define  $\sigma_j^{\text{mapr}, \theta} = \frac{\mathbf{u}_{j+1} - \mathbf{u}_{j-1}}{2\Delta x}$ .

We can also define an abstract limiter which guarantees that the second-order reconstruction does not violate a local invariant domain property. In the scalar case this is the well known local maximum principle and minmod- $\theta$  will do the job but for nonlinear systems the invariant domain depends on the systems at hand and minmod- $\theta$  is not an appropriate choice. We call such a limiter an invariant domain limiter  $\sigma_j^{\text{inv}}$ . Namely, we define the invariant slope  $\sigma_j^{\text{inv}}$  to be such that the values  $\mathbf{u}_{j+1/2}^{n,-} := \mathbf{u}_j^n + \frac{\sigma_j^{\text{inv}} \Delta x}{2}$  and  $\mathbf{u}_{j-1/2}^{n,+} := \mathbf{u}_j^n - \frac{\sigma_j^{\text{inv}} \Delta x}{2}$  are in local invariant sets defined by the user. More precisely, we have the following local abstract limitation.

**DEFINITION 3.1.** *Let  $A_{j-1/2}$  be an invariant set of (2.1) containing the states  $\mathbf{u}_{j-1}^n$  and  $\mathbf{u}_j^n$  and  $A_{j+1/2}$  be an invariant set of (2.1) containing  $\mathbf{u}_j^n$  and  $\mathbf{u}_{j+1}^n$ . Then the invariant slope  $\sigma_j^{\text{inv}}$  corresponding to the invariant sets  $A_{j-1/2}$  and  $A_{j+1/2}$  is defined as  $\sigma_j^{\text{inv}} = \ell \frac{\mathbf{u}_{j+1} - \mathbf{u}_{j-1}}{2\Delta x}$  where  $\ell$  is the largest number in  $[0, 1]$  such that  $\mathbf{u}_{j-1/2}^{n,+} \in A_{j-1/2}$  and  $\mathbf{u}_{j+1/2}^{n,-} \in A_{j+1/2}$ .*

The actual computation of such limiters and a precise definition of the local sets  $A_{j-1/2}$  which allows the method to be second order will be given later, see Algorithm 2 and the numerical tests but for the time being we stay in this abstract setting. The key difference between minmod-type slope limiting (or any other classical limiting) and the invariant domain slope limiting is that one tries to impose a local maximum principle (or reduce oscillations in physical space) and the invariant domain slope limiting only limits the slopes so that the invariant domain property is imposed in phase space at cell interfaces. With the above notations we are now ready to describe the two new methods we propose in this paper:

- (1) Method 1 is based on convex flux limiting. It uses the MAPR slope, see (3.25), when  $\theta_j^n \leq 1.5$  and the central slope if  $\theta_j^n > 1.5$ . The value of  $\theta$  is computed using  $\theta_j^n = 2 - R_j^n$  with  $R_j^n$  the entropy commutator defined in (3.30). The method is made

invariant domain preserving using convex flux limiting, see Algorithm 4.2.1. In the paper we refer to this method as the MAPR-EV-CL method, where EV stands for the use of entropy viscosity commutator used to determine  $\theta_j^n$  and CL refers to the convex flux limiting used to enforce the local invariant domain property.

- (2) Method 2 is based on convex slope limiting only. It uses the invariant domain slope limiter, see Definition 3.1, and the resulting scheme is invariant domain preserving under standard CFL, see Theorem 4.1 and Theorem 4.2; in the paper we refer to this method as the INV-CL method, where INV is for invariant domain preserving and CL refers to the convex limiting used to generate this invariant domain preserving slope limiter.

Any other method used in the simulations will be identified by the limiter used in the KT-scheme, for example the method based on the minmod limiter (referred to minmod in the numerical section) will be shown as a standard comparison in all numerical examples. We now continue with the exact construction of  $\theta_j$  required in the definition of the MAPR limiter.

**3.3.2. Entropy based smoothness indicator.** In this section, we are going to consider a different approach to create a new limited reconstruction. Namely, we would like to use the central unlimited slope in smooth regions and a nonlinear minmod-type limited slope in the regions of discontinuities. Moreover, the change between the two reconstructions should happen when a physical discontinuity forms. Similar to [8], the approach we take to detect a discontinuity is to measure an entropy production. Our objective is to construct a second-order method that is entropy consistent and at the same time close to being invariant domain preserving. However, we do not want to rely on the yet to be explained limiting process to enforce entropy consistency. We refer the reader to Lemma 3.2, Lemma 4.6 and §6.1 in [7] and [6, §5.1] for counter-examples of methods that are invariant domain preserving but entropy violating. Invariant domain limitation should be understood as a light post-processing applied to a method that is already entropy consistent and almost invariant domain preserving. In [8, 9], a high-order graph viscosity that is entropy consistent was introduced. However, we do not want the time discretization to interfere with the estimation of the residual, and we follow the entropy viscosity commutator approach proposed in [11]. For simplicity, we present the entropy viscosity commutator in the one dimensional case. Let  $(\eta(\mathbf{u}), \mathbf{F}(\mathbf{u}))$  be the entropy pair of system (3.1), that is,  $\eta$  is a convex function of the vector of conserved variables  $\mathbf{u}$ , and the entropy flux  $\mathbf{F}$  satisfies  $D\mathbf{F}(\mathbf{u}) = \eta'(\mathbf{u})^\top D\mathbf{f}(\mathbf{u})$ . Following [11, §3.4] and [12, §6.4] we measure the discrepancy in the chain rule as follows

$$\Delta_j^n = \mathbf{F}(\mathbf{u}_{j+1}^n) - \mathbf{F}(\mathbf{u}_{j-1}^n) - \eta'(\mathbf{u}_j^n)^\top (\mathbf{f}(\mathbf{u}_{j+1}^n) - \mathbf{f}(\mathbf{u}_{j-1}^n)). \tag{3.27}$$

We set

$$C_j^n = |\mathbf{F}(\mathbf{u}_{j+1}^n) - \mathbf{F}(\mathbf{u}_{j-1}^n)| + |\eta'(\mathbf{u}_j^n)^\top| \cdot |\mathbf{f}(\mathbf{u}_{j+1}^n) - \mathbf{f}(\mathbf{u}_{j-1}^n)|, \tag{3.28}$$

to be a corresponding normalizing coefficient, where for a vector function  $\mathbf{g}$  we denote  $|\mathbf{g}| := \|\mathbf{g}\|_{\ell_2}$ . Notice that in smooth regions  $C_j^n$  could be very close to or even zero. Thus, to avoid division by zero, we set

$$\alpha_j^n = \max(|\mathbf{F}(\mathbf{u}_{j+1}^n)|, |\mathbf{F}(\mathbf{u}_j^n)|, |\mathbf{F}(\mathbf{u}_{j-1}^n)|), \tag{3.29a}$$

$$\beta_j^n = |\eta'(\mathbf{u}_j^n)^\top| \cdot \lambda_j^{\max, n} \cdot (|\mathbf{u}_{j+1}^n - \mathbf{u}_j^n| + |\mathbf{u}_j^n - \mathbf{u}_{j-1}^n|), \tag{3.29b}$$

where  $\lambda_j^{\max,n} := \max(\lambda_{j+1/2}^n, \lambda_{j-1/2}^n)$  is the global maximum speed of propagation at time  $t^n$  and define the normalized entropy viscosity commutator as

$$R_j^n = \frac{|\Delta_j^n|}{\max(C_j^n, \epsilon \alpha_j^n, \epsilon \beta_j^n)}, \quad (3.30)$$

where  $\epsilon$  is a small number, for example  $\epsilon := 10^{-8}$ . By definition, we have that  $R_j^n \in (0, 1]$  because  $|\Delta_j^n| \leq C_j^n$ . Note that (3.27) is a discrete version of  $\mathbf{F}' - \eta' \cdot \mathbf{f}'$  and because on continuous level  $\mathbf{F}' - \eta' \cdot \mathbf{f}' = 0$ , using Taylor's expansion one could prove that  $R_j^n = \mathcal{O}(\Delta x)$  in smooth regions and  $R_j^n \sim 1$  in the regions of shocks. For more details on discrete entropy viscosity commutators, we refer the reader to [11, §3.4]. We now define the local weights  $\theta_j^n$  needed for the MAPR limiter (3.25) as follows

$$1 \leq \theta_j^n := 2 - R_j^n \leq 2. \quad (3.31)$$

In smooth regions, we have  $\theta_j^n = 2 - \mathcal{O}(\Delta x)$  and near shocks we have  $\theta_j^n \sim 1$ .

#### 4. Quasiconcavity based limitation

In this section, we introduce a technique which will modify an existing second-order method, for example the original KT-scheme, and make it local invariant domain preserving. We present two limiting techniques which can do that, both based on the so-called *convex limiting* first introduced in [11]. Both of these limitations will upgrade the KT-scheme to satisfy an invariant domain property and numerically preserve the second-order accuracy of the method. The first approach is called convex flux limiting, see Algorithm 4.2.1, and will make any KT-scheme based on any slope limiting invariant domain preserving. This is the approach we use when applying the MAPR limiter, i.e., Method 1. The second approach is based on convex slope limiting and guarantees that after the slope limitation the linear function is in the local invariant domain, see Algorithm 2. This is the approach we use when applying Method 2, no flux limiting is needed.

Invariant domains are convex sets in phase space. In general, second-order finite volume methods, and in particular the KT-scheme, may violate the invariant domain property if the cell averages are on the boundary of the invariant set. For example, assume that the state  $\mathbf{u}_j^n$  is on  $\partial S$ , where  $S$  is a local invariant set,  $\mathbf{u}_{j-1}^n, \mathbf{u}_j^n, \mathbf{u}_{j+1}^n \in S$ . Then the piecewise linear reconstruction using any nonzero  $\sigma_j$  will create interface values  $\mathbf{u}_{j+1/2}^{n,+}$  and  $\mathbf{u}_{j+1/2}^{n,-}$ . Because  $\mathbf{u}_j^n$  is the midpoint of the segment connecting  $\mathbf{u}_{j+1/2}^{n,+}$  and  $\mathbf{u}_{j+1/2}^{n,-}$ , we conclude that at least one of the points  $\mathbf{u}_{j+1/2}^{n,+}$  or  $\mathbf{u}_{j+1/2}^{n,-}$  is outside  $S$ , thus violates the invariant domain property. Once an interface value moves outside the invariant region, we may have that  $\mathbf{u}_j^{n+1}$  is outside the set  $S$  under a standard CFL condition. We have observed that to be true numerically for the KT-scheme and other second-order finite volume schemes in many numerical tests where the invariant domain boundary is smooth but not affine. Examples of systems with non-affine local invariant sets are the p-system (see §2.2.2) and the Euler equations (see §2.2.3).

**4.1. Invariant domains via quasiconcave constraints.** In order to unify into a single framework all the bounds that we want to enforce on the second-order solution, we are going to rely on the notion of quasiconcavity, which we now recall.

**DEFINITION 4.1 (Quasiconcavity).** *Given a convex set  $\mathcal{A} \subset \mathbb{R}^m$ , we say that a function  $\Psi: \mathcal{A} \rightarrow \mathbb{R}$  is quasiconcave if every upper level set of  $\Psi$  is convex; that is, the set  $L_\lambda(\Psi) := \{\mathbf{u} \in \mathcal{A} | \Psi(\mathbf{u}) \geq \lambda\}$  is convex for any  $\lambda \in \mathbb{R}$  in the range of  $\Psi$ .*

Note that concavity implies quasiconcavity. In all hyperbolic problems we consider, we assume that the invariant domain can be described as an intersection of quasiconcave constraints of the type  $\Psi(\mathbf{u}) \geq 0$  and we will enforce such quasiconcave constraints via the convex limiting procedure introduced in [11, §4.2]. Moreover, in practice we will modify all quasiconcave constraints to be concave constraints because the limitation process is much simpler in the concave case. We now describe the set quasiconcave constraints for all cases considered in the paper.

**4.1.1. Scalar equations.** In Example 1 in §2.2.1, the invariant domain is an interval and we enforce it by imposing a local maximum principle. To be precise, we set

$$u_j^{\min} := \min(u_j^n, u_{j\pm 1}^n, \bar{u}_{j\pm 1/2}^{n+1}), \quad u_j^{\max} := \max(u_j^n, u_{j\pm 1}^n, \bar{u}_{j\pm 1/2}^{n+1}). \tag{4.1}$$

Theorem 3.1 guarantees that the first-order method satisfies the local maximum principle  $u_j^{\min} \leq u_j^{n+1,L} \leq u_j^{\max}$ . So, the convex limiting must enforce  $u_j^{\min} \leq u_j^{n+1} \leq u_j^{\max}$  to guarantee an invariant domain property. By setting  $\Psi_j^1(u) = u - u_j^{\min}$  and  $\Psi_j^2(u) = u_j^{\max} - u$ , we transform imposing the local maximum principle to imposing two linear (therefore quasiconcave) constraints:  $\Psi_j^1(u), \Psi_j^2(u) \geq 0$ .

**4.1.2. The p-system.** In the case of p-system (2.5), see §2.2.2, we use the Riemann invariants (2.7) to set  $a = w_{2,j}^{\min}$  and  $b = w_{1,j}^{\max}$ , with the definition

$$w_{1,j}^{\max} := \max(w_{1,j}^n, w_{1,j\pm 1}^n, \bar{w}_{1,j\pm 1/2}^{n+1}), \quad w_{2,j}^{\min} := \min(w_{2,j}^n, w_{2,j\pm 1}^n, \bar{w}_{2,j\pm 1/2}^{n+1}), \tag{4.2}$$

where  $w_{1,j}^n := w_1(\mathbf{u}_j^n)$ ,  $w_{2,j}^n := w_2(\mathbf{u}_j^n)$ ,  $\bar{w}_{1,j\pm 1/2}^{n+1} := w_1(\bar{\mathbf{u}}_{j\pm 1/2}^{n+1})$ , and  $\bar{w}_{2,j\pm 1/2}^{n+1} := w_2(\bar{\mathbf{u}}_{j\pm 1/2}^{n+1})$ . Theorem 3.1 guarantees that  $w_{2,i}^{\min} \leq w_{2,i}^{L,n+1} \leq w_{1,i}^{L,n+1} \leq w_{1,i}^{\max}$ . Therefore, the local invariant set to be enforced is an intersection of two local concave constraints:  $w_{2,j}^{\min} \leq w_{2,j}^{n+1}$  and  $w_{1,j}^{n+1} \leq w_{1,j}^{\max}$ . By setting  $\Psi_j^1(\mathbf{u}) = w_{1,j}^{\max} - w_1(\mathbf{u})$  and  $\Psi_j^2(\mathbf{u}) = w_2(\mathbf{u}) - w_{2,j}^{\min}$ , we have that the concave constraints we are going to enforce in this case are:  $\Psi_j^1(\mathbf{u}), \Psi_j^2(\mathbf{u}) \geq 0$ .

**4.1.3. Euler equations.** In the case of the Euler system (2.9), it is known that the specific entropy is a quasiconcave function of the conserved variables, that is  $\Phi(\mathbf{u}) := s(\rho, e)$  is quasiconcave. The first-order solution  $\mathbf{u}^{L,n+1}$  satisfies

$$\rho_j^{\max} \geq \rho_j^{L,n+1}, \quad \rho_j^{L,n+1} \geq \rho_j^{\min}, \quad e_j^{L,n+1} \geq 0, \quad s_j^{L,n+1} \geq s_j^{\min}, \tag{4.3}$$

where we set

$$\begin{aligned} \rho_j^{n,\min} &:= \min(\rho_j^n, \rho_{j\pm 1}^n, \bar{\rho}_{j\pm 1/2}^{n+1}), & \rho_j^{n,\max} &:= \max(\rho_j^n, \rho_{j\pm 1}^n, \bar{\rho}_{j\pm 1/2}^{n+1}), \\ e_j^{n,\min} &:= \min(e_j^n, e_{j\pm 1}^n, \bar{e}_{j\pm 1/2}^{n+1}), \\ s_j^{n,\min} &:= \min(\Phi(\mathbf{u}_j^n), \Phi(\mathbf{u}_{j\pm 1}^n), \Phi(\bar{\mathbf{u}}_{j\pm 1/2}^{n+1})). \end{aligned} \tag{4.4}$$

Then, the invariant set (2.10) can be imposed by enforcing that the high-order solution  $\mathbf{u}^{H,n+1}$  be in the intersection of the following four quasiconcave constraints:

$$\rho_j^{n,\max} \geq \rho, \quad \rho \geq \rho_j^{n,\min}, \quad e \geq e_j^{n,\min}, \quad s \geq s_j^{n,\min}. \tag{4.5}$$

Two of the above four constraints are not concave:  $e \geq e_j^{n,\min}$  and  $s - s_j^{n,\min} \geq 0$ . However, one can modify the constraints, assuming that the density is already positive, and

make them concave. For example, the mathematical entropy  $\rho s$  and the total internal energy  $\rho e$  are concave functions of the conserved variables. Therefore, the modified constraints  $\rho e - \rho e_j^{n,\min} \geq 0$  and  $\rho s - \rho s_j^{n,\min} \geq 0$  are concave.

In the case of the  $\gamma$ -law equation of state,  $s = \log(e^{\frac{1}{\gamma-1}} \rho^{-1})$ , one can impose the invariant set (2.10) by enforcing

$$\rho_j^{n,\max} - \rho \geq 0, \quad \rho - \rho_j^{n,\min} \geq 0, \quad \rho e - c_j^{n,\min} \rho^\gamma \geq 0. \tag{4.6}$$

where  $c_j^{n,\min} = \exp((\gamma - 1)s_j^{n,\min})$ .

Note that, the constraints  $\rho > 0$  and  $\rho e - c_j^{n,\min}(\rho_j^{n+1})^\gamma \geq 0$  imply that the internal energy is positive. Thus, by setting  $\Psi_j^1(\mathbf{u}) = \rho_j^{n,\max} - \rho$ ,  $\Psi_j^2(\mathbf{u}) = \rho - \rho_j^{n,\min}$  and  $\Psi_j^3(\mathbf{u}) = \rho e - c_j^{n,\min} \rho^\gamma$ , we will enforce that  $\Psi_j^1(\mathbf{u}), \Psi_j^2(\mathbf{u}), \Psi_j^3(\mathbf{u}) \geq 0$  to guarantee the invariant domain property (2.10).

**4.2. Invariant domain via flux limiting.** In this section we develop a novel limiting technique for enforcing quasiconcave constraints. We adopt the methodology of [11] to the central scheme framework. Simple linear constraints like  $\rho_j^{\min} \leq \rho_{j+1}^n \leq \rho_j^{\max}$  can be easily enforced by using the flux corrected transport (FCT) technique, see for example [24] and [1]. However, the FCT approach is designed for box-like limitation and cannot be easily modified to enforce general convex constraints without losing second-order accuracy. Moreover, the use of the bar states (3.7), (3.13) and (3.14) is critical in the definition of the local constraints, see settings of the invariant bounds (4.1), (4.2) and (4.4) for examples.

**4.2.1. Flux limiting algorithm.** We subtract the first-order update (3.5) from the high-order update (3.18) and get:

$$\mathbf{u}_j^{H,n+1} = \mathbf{u}_j^{L,n+1} - \frac{\Delta t}{\Delta x} (H_{j+1/2}^n - H_{j-1/2}^n - L_{j+1/2}^n + L_{j-1/2}^n). \tag{4.7}$$

By setting  $G_{j+1/2}^n := \frac{2\Delta t}{\Delta x} (H_{j+1/2}^n - L_{j+1/2}^n)$  to be the high/low order flux difference, we rewrite (4.7) as the following convex splitting form

$$\mathbf{u}_j^{H,n+1} = \frac{1}{2}(\mathbf{u}_j^{L,n+1} - G_{j+1/2}^n) + \frac{1}{2}(\mathbf{u}_j^{L,n+1} + G_{j-1/2}^n). \tag{4.8}$$

Following [11, §4.2], we introduce a pair of scalar limiting parameters  $(l_j^+, l_j^-)$  to create a limited second-order update

$$\mathbf{u}_j^{n+1}(l_j^+, l_j^-) := \frac{1}{2}(\mathbf{u}_j^{L,n+1} - l_j^+ G_{j+1/2}^n) + \frac{1}{2}(\mathbf{u}_j^{L,n+1} + l_j^- G_{j-1/2}^n). \tag{4.9}$$

which should satisfy the invariant domain property. Similar to the FCT approach, we recover the first-order solution if  $l_j^+ = l_j^- = 0$  and the second-order solution if  $l_j^+ = l_j^- = 1$ .

Let  $\Psi_j^z$  be a quasiconcave function where  $z$  is one of the constraints describing the local invariant set at cell  $j$ . We denote with  $A_j^z$  the zero level set of  $\Psi_j^z$ . For example, we set  $\Psi_j^{\rho_{\max}} = \rho_j^{\max} - \rho$  and  $A_j^{\rho_{\max}} = \{\mathbf{u} \mid \rho_j^{\max} - \rho \geq 0\}$  for the Euler system. By definition,  $A_j^z$  is a convex set. The goal is to find the largest positive numbers  $l_j^\pm \leq 1$  such that  $\mathbf{u}_j^{n+1}(l_j^+, l_j^-)$  is in  $A_j^z$  i.e.,  $\Psi_j^z(\mathbf{u}_j^{n+1}(l_j^+, l_j^-)) \geq 0$  for any  $0 \leq l_j^+ \leq l_j^+$  and  $0 \leq l_j^- \leq l_j^-$ . In order to simplify the notations, for any  $l \in \mathbb{R}$  we denote  $\mathbf{u}_j^+(l) := \mathbf{u}_j^{L,n+1} - lG_{j+1/2}^n$  and  $\mathbf{u}_j^-(l) := \mathbf{u}_j^{L,n+1} + lG_{j-1/2}^n$ . The following two lemmas describe the flux limiting process for a given constraint  $z$ .

LEMMA 4.1. Let  $\Psi_j^z: \mathcal{A} \rightarrow \mathbb{R}$  be the quasiconcave function mentioned above. Assume that  $\ell_j^{z,\pm} \in [0, 1]$  are such that  $\Psi_j^z(\mathbf{u}_j^+(\ell_j^{z,+})) \geq 0$  and  $\Psi_j^z(\mathbf{u}_j^-(\ell_j^{z,-})) \geq 0$ , then we have that  $\Psi_j^z(\mathbf{u}_j^{n+1}(\ell_j^{z,+}, \ell_j^{z,-})) \geq 0$ .

LEMMA 4.2. Let's define  $\ell_j^{z,\pm}$  to be

$$\ell_j^{z,+} = \begin{cases} 1 & \text{if } \Psi_j^z(\mathbf{u}_j^+(1)) \geq 0, \\ \max\{\ell \in [0, 1] \mid \Psi_j^z(\mathbf{u}_j^+(\ell)) \geq 0\} & \text{otherwise.} \end{cases} \quad (4.10)$$

$$\ell_j^{z,-} = \begin{cases} 1 & \text{if } \Psi_j^z(\mathbf{u}_j^-(1)) \geq 0, \\ \max\{\ell \in [0, 1] \mid \Psi_j^z(\mathbf{u}_j^-(\ell)) \geq 0\} & \text{otherwise.} \end{cases} \quad (4.11)$$

Now we set  $\ell_{j+1/2}^z = \min(\ell_j^{z,+}, \ell_{j+1}^{z,-})$ , we have that for all  $\ell_{j+1/2}^z \in [0, \ell_{j+1/2}^z]$ , it holds that  $\Psi_j^z(\mathbf{u}_j^{n+1}(\ell_{j+1/2}^z, \ell_{j-1/2}^z)) \geq 0$ .

REMARK 4.1. We refer the reader to [11] because the proofs of Lemma 4.1 and Lemma 4.2 are analogous to the proofs of Lemma 4.3 and Lemma 4.4 in [11].

The second-order update  $\mathbf{u}_j^{n+1}$  is a convex combination of  $\mathbf{u}_j^+$  and  $\mathbf{u}_j^-$ , see (4.9). Therefore, using the above lemmas, the limited update

$$\mathbf{u}_j^{z,n+1} = \frac{1}{2} \mathbf{u}_j^+(\ell_{j+1/2}^z) + \frac{1}{2} \mathbf{u}_j^-(\ell_{j-1/2}^z). \quad (4.12)$$

satisfies the constraint  $z$ . We now describe the full flux limiting process for all local constraints in the following algorithm.

---

**Algorithm 1** Convex flux limiting

---

**Input:**  $\mathbf{u}_j^{L,n+1}$ ,  $G_{j+1/2}^n$ ,  $k^{\max}$ ,  $z_1, \dots, z_q$ .

**Output:**  $\mathbf{u}_j^{n+1}$

- 1: **for**  $i = 1$  **to**  $k^{\max}$  **do**
  - 2:   **for**  $z = 1$  **to**  $z_q$  **do**
  - 3:     Compute limiting parameters  $\ell_{j+1/2}^z$  via Lemma 4.1 and Lemma 4.2.
  - 4:   **end for**
  - 5:   Set  $\ell_{j+1/2} := \min_{z \in \{z_1, \dots, z_q\}} \ell_{j+1/2}^z$ .
  - 6:   Update  $\mathbf{u}_j^{n+1} = \mathbf{u}_j^{n+1}(\ell_{j-1/2}, \ell_{j+1/2})$  via (4.9).
  - 7:   Update  $G_{j+1/2}^{n+1} = \frac{2\Delta t}{\Delta x} (H_{j+1/2}^{n+1} - L_{j+1/2}^{n+1})$
  - 8: **end for**
  - 9: **Return**  $\mathbf{u}_j^{n+1}$ .
- 

REMARK 4.2. In the numerical experiments reported at the end of this paper, we take  $k^{\max} = 1$ .

REMARK 4.3. The computational cost of finding  $\ell_j^{z,\pm}$  for a given  $j$  can be reduced by setting  $\ell_j^+ = \ell_j^- := \ell_j$  in (4.9) and denote  $\mathbf{u}_j^{n+1}(\ell) := \mathbf{u}_j^{n+1}(\ell, \ell)$ . Then  $\ell_j$  is computed with one line search

$$\ell_j = \begin{cases} 1 & \text{if } \Psi_j^z(\mathbf{u}_j^{n+1}(1)) \geq 0, \\ \max\{\ell \in [0, 1] \mid \Psi_j^z(\mathbf{u}_j^{n+1}(\ell)) \geq 0\} & \text{otherwise.} \end{cases} \quad (4.13)$$

If  $\Psi_j^z(\mathbf{u}_j^\pm(\ell_j)) \geq 0$ , we set  $\ell_{j+1/2} = \min(\ell_j, \ell_{j+1})$  and skip **step 3** in the algorithm. In practice, the single search is successful most of the time and therefore we make one line search instead of  $2d$  line searches where  $d$  is the space dimension.

In the case of two space dimensions, subtracting (3.10) from (3.21) and setting  $G_{j+1/2,k}^{n,x} := \frac{4\Delta t}{\Delta x}(H_{j+1/2,k}^{n,x} - L_{j+1/2,k}^{n,x})$ ,  $G_{j,k+1/2}^{n,y} := \frac{4\Delta t}{\Delta y}(H_{j,k+1/2}^{n,y} - L_{j,k+1/2}^{n,y})$ , we obtain the following convex splitting form for the high-order solution

$$\begin{aligned} \mathbf{u}_{j,k}^{H,n+1} &= \frac{1}{4}(\mathbf{u}_{j,k}^{L,n+1} - G_{j+1/2,k}^{n,x}) + \frac{1}{4}(\mathbf{u}_{j,k}^{L,n+1} + G_{j-1/2,k}^{n,x}) \\ &\quad + \frac{1}{4}(\mathbf{u}_{j,k}^{L,n+1} - G_{j,k+1/2}^{n,y}) + \frac{1}{4}(\mathbf{u}_{j,k}^{L,n+1} + G_{j,k-1/2}^{n,y}). \end{aligned} \quad (4.14)$$

The corresponding limited second-order update is given by

$$\begin{aligned} \mathbf{u}_{j,k}^{n+1}(l_{j,k}^{x,\pm}, l_{j,k}^{y,\pm}) &:= \frac{1}{4}(\mathbf{u}_{j,k}^{L,n+1} - l_{j,k}^{x,+} G_{j+1/2,k}^{n,x}) + \frac{1}{4}(\mathbf{u}_{j,k}^{L,n+1} + l_{j,k}^{x,-} G_{j-1/2,k}^{n,x}) \\ &\quad + \frac{1}{4}(\mathbf{u}_{j,k}^{L,n+1} - l_{j,k}^{y,+} G_{j,k+1/2}^{n,y}) + \frac{1}{4}(\mathbf{u}_{j,k}^{L,n+1} + l_{j,k}^{y,-} G_{j,k-1/2}^{n,y}) \\ &:= \frac{1}{4}\mathbf{u}_{j,k}^{x,+}(l_{j,k}^{x,+}) + \frac{1}{4}\mathbf{u}_{j,k}^{x,-}(l_{j,k}^{x,-}) + \frac{1}{4}\mathbf{u}_{j,k}^{y,+}(l_{j,k}^{y,+}) + \frac{1}{4}\mathbf{u}_{j,k}^{y,-}(l_{j,k}^{y,-}). \end{aligned} \quad (4.15)$$

where the four scalar limiting parameters  $l_{j,k}^{x,\pm}$  and  $l_{j,k}^{y,\pm}$  are computed such that each of the four states  $\mathbf{u}_{j,k}^{x,+}(l_{j,k}^{x,+})$ ,  $\mathbf{u}_{j,k}^{x,-}(l_{j,k}^{x,-})$ ,  $\mathbf{u}_{j,k}^{y,+}(l_{j,k}^{y,+})$ ,  $\mathbf{u}_{j,k}^{y,-}(l_{j,k}^{y,-})$  satisfies the invariant domain property. Similar to the one dimensional case, given a quasiconcave function  $\Psi_{j,k}^z$  which describes a local constraint  $z$  at a cell  $(j,k)$  with a zero level set  $A_{j,k}^z$ , we find the largest positive numbers  $\ell_{j,k}^{x,\pm}, \ell_{j,k}^{y,\pm} \leq 1$  such that the above mentioned four states are in  $A_{j,k}^z$  for any  $0 \leq l_{j,k}^{x,+} \leq \ell_{j,k}^{x,+}$ ,  $0 \leq l_{j,k}^{x,-} \leq \ell_{j,k}^{x,-}$  and  $0 \leq l_{j,k}^{y,+} \leq \ell_{j,k}^{y,+}$ ,  $0 \leq l_{j,k}^{y,-} \leq \ell_{j,k}^{y,-}$ . Analogous to the one dimensional case, these limiters are computed via line searches and one can use a single line search instead of four most of the time, see Remark 4.3.

**4.3. Invariant domain via slope limiting.** It is well known in the literature that one can reduce oscillations via either flux limiting or slope limiting. Both limitations are different but give similar numerical results. In this section we describe a convex limiting procedure using slope limiting. The key difference is that the local invariant sets to be enforced are now located at interfaces and are different from the local invariant sets at cell centers used in flux limiting, see §4.1 and §4.2.

We start with the one dimensional case. Instead of enforcing  $\mathbf{u}_j^{n+1}$  to be in the invariant set, we will limit the interface values given by the local linear reconstructions:  $\mathbf{u}_{j-1/2}^{n,\pm}$  and  $\mathbf{u}_{j+1/2}^{n,\pm}$ . In the fully discrete KT-scheme (3.17)-(3.18) we will change the notation and use  $\lambda_{j+1/2}^n$  instead of  $a_{j+1/2}^n$  to denote the local speed in (3.17). The original second-order KT-scheme is not invariant domain preserving in general. However, a modification of the KT-scheme is going to be invariant domain preserving under a new CFL-condition.

**THEOREM 4.1.** *Let  $A$  be a convex invariant set of (2.1),  $n \geq 0$  and  $j \in \mathbb{Z}$  be such that  $\mathbf{u}_j^n$  and all interface values  $\mathbf{u}_{j\pm 1/2}^{n,\pm}$  are in  $A$ . Assume that the second-order solution  $\mathbf{u}_j^{n+1}$  is computed with the KT-scheme (3.17)-(3.18), where  $\lambda_{j-1/2}^n := a_{j-1/2}^n = \lambda_{\max}(\mathbf{u}_{j-1/2}^{n,-}, \mathbf{u}_{j-1/2}^{n,+}, \mathbf{f})$  and  $\lambda_{j+1/2}^n := a_{j+1/2}^n = \lambda_{\max}(\mathbf{u}_{j+1/2}^{n,-}, \mathbf{u}_{j+1/2}^{n,+}, \mathbf{f})$ . Let  $\lambda_{j,\pm}^n := \lambda_{\max}(\mathbf{u}_{j-1/2}^{n,+}, \mathbf{u}_{j+1/2}^{n,-}, \mathbf{f})$  be the in cell local speed and define the maximum local*



speed by  $\lambda_j^{\max} := \max(\lambda_{j-1/2}^n, \lambda_{j+1/2}^n, \lambda_{j,\pm}^n)$ . Then under the CFL condition  $\frac{\Delta t}{\Delta x} \lambda_j^{\max} \leq \frac{1}{4}$ , we have that  $\mathbf{u}_j^{n+1} \in A$ .

*Proof.* Using the definition of  $\lambda_{j+1/2}^n$  we have that the bar state

$$\bar{\mathbf{u}}_{j+1/2,\pm}^{n+1} := \frac{\mathbf{u}_{j+1/2}^{n,+} + \mathbf{u}_{j+1/2}^{n,-}}{2} - \frac{\mathbf{f}(\mathbf{u}_{j+1/2}^{n,+}) - \mathbf{f}(\mathbf{u}_{j+1/2}^{n,-})}{2\lambda_{j+1/2}^n} \tag{4.16}$$

is in the invariant domain  $A$ , see Lemma 2.1. Similarly, we have that the bar state

$$\bar{\mathbf{u}}_{j-1/2,\pm}^{n+1} := \frac{\mathbf{u}_{j-1/2}^{n,+} + \mathbf{u}_{j-1/2}^{n,-}}{2} - \frac{\mathbf{f}(\mathbf{u}_{j-1/2}^{n,+}) - \mathbf{f}(\mathbf{u}_{j-1/2}^{n,-})}{2\lambda_{j-1/2}^n} \tag{4.17}$$

is also in  $A$ . With this notation, we represent the KT-scheme update as follows

$$\begin{aligned} \mathbf{u}_j^{n+1} = & \mathbf{u}_j^n + \frac{\Delta t}{\Delta x} \lambda_{j+1/2}^n \bar{\mathbf{u}}_{j+1/2,\pm}^{n+1} + \frac{\Delta t}{\Delta x} \lambda_{j-1/2}^n \bar{\mathbf{u}}_{j-1/2,\pm}^{n+1} \\ & + \frac{2\Delta t}{\Delta x} \left( -\frac{\lambda_{j+1/2}^n \mathbf{u}_{j+1/2}^{n,-}}{2} - \frac{\lambda_{j-1/2}^n \mathbf{u}_{j-1/2}^{n,+}}{2} - \frac{\mathbf{f}(\mathbf{u}_{j+1/2}^{n,-}) - \mathbf{f}(\mathbf{u}_{j-1/2}^{n,+})}{2} \right). \end{aligned} \tag{4.18}$$

We now define another bar state

$$\bar{\mathbf{u}}_{j,\pm}^{n+1} := \frac{\mathbf{u}_{j-1/2}^{n,+} + \mathbf{u}_{j+1/2}^{n,-}}{2} - \frac{\mathbf{f}(\mathbf{u}_{j+1/2}^{n,-}) - \mathbf{f}(\mathbf{u}_{j-1/2}^{n,+})}{2\lambda_j^{\max}} \tag{4.19}$$

which is also in the invariant domain  $A$  because  $\lambda_j^{\max} \geq \lambda_{\max}(\mathbf{u}_{j-1/2}^{n,+}, \mathbf{u}_{j+1/2}^{n,-}, \mathbf{f})$ , see Lemma 2.1. Using (4.19) in (4.18) and the fact that  $\frac{\mathbf{u}_{j-1/2}^{n,+} + \mathbf{u}_{j+1/2}^{n,-}}{2} = \mathbf{u}_j^n$ , we obtain

$$\begin{aligned} \mathbf{u}_j^{n+1} = & \left(1 - 4\frac{\Delta t}{\Delta x} \lambda_j^{\max}\right) \mathbf{u}_j^n \\ & + \frac{\Delta t}{\Delta x} \lambda_{j+1/2}^n \bar{\mathbf{u}}_{j+1/2,\pm}^{n+1} + \frac{\Delta t}{\Delta x} \lambda_{j-1/2}^n \bar{\mathbf{u}}_{j-1/2,\pm}^{n+1} + \frac{2\Delta t}{\Delta x} \lambda_j^{\max} \bar{\mathbf{u}}_{j,\pm}^{n+1} \\ & + \frac{\Delta t}{\Delta x} (\lambda_j^{\max} - \lambda_{j-1/2}^n) \mathbf{u}_{j-1/2}^{n,+} + \frac{\Delta t}{\Delta x} (\lambda_j^{\max} - \lambda_{j+1/2}^n) \mathbf{u}_{j-1/2}^{n,-}. \end{aligned} \tag{4.20}$$

Under the CFL-condition  $\frac{\Delta t}{\Delta x} \lambda_j^{\max} \leq \frac{1}{4}$ , we have that  $\mathbf{u}_j^{n+1}$  is a convex combination of  $\mathbf{u}_j^n$ , the bar states  $\bar{\mathbf{u}}_{j-1/2,\pm}^{n+1}$ ,  $\bar{\mathbf{u}}_{j+1/2,\pm}^{n+1}$ ,  $\bar{\mathbf{u}}_{j,\pm}^{n+1}$ , and the interface states  $\mathbf{u}_{j-1/2}^{n,+}$ ,  $\mathbf{u}_{j-1/2}^{n,-}$ . Since the interface states  $\mathbf{u}_{j-1/2}^{n,+}$ ,  $\mathbf{u}_{j-1/2}^{n,-}$  are assumed to be in the invariant domain  $A$ , and because of the definition of  $\lambda_j^{\max}$  all bar states are also in  $A$ , then it follows by convexity that  $\mathbf{u}_j^{n+1} \in A$ .  $\square$

REMARK 4.4. Note that, in order to have the invariant domain property we need to design a limited piecewise linear reconstruction so that the interface values  $\mathbf{u}_{j\pm 1/2}^{n,\pm}$  are in the local invariant set  $A$ . If the local slope is set to be zero we recover the first-order result, see Theorem 3.1 and Remark 3.2.

We now describe a convex slope limiting process which guarantees that the interface states  $\mathbf{u}_{j-1/2}^{n,-}$  and  $\mathbf{u}_{j-1/2}^{n,+}$  and can be modified to be in a given local invariant set  $A$ . As before, we are going to impose a finite set of quasiconcave constraints  $\Psi_{j-1/2}^z$ ,  $z \in \{z_1, \dots, z_q\}$ , with the assumption that enforcing these guarantees that the interface

states  $\mathbf{u}_{j-1/2}^{n,-}$  and  $\mathbf{u}_{j-1/2}^{n,+}$  are in  $A$ . Similar to the flux limiting process, see §4.2.1, we denote by  $A_{j-1/2}^z$  the zero level set of  $\Psi_{j-1/2}^z$ , thus enforcing  $\mathbf{u} \in A_{j-1/2}^z$  is equivalent to enforcing  $\Psi_{j-1/2}^z(\mathbf{u}) \geq 0$ . By definition we have  $\mathbf{u}_{j-1/2}^{n,+} = \mathbf{u}_j^n - \frac{\Delta x}{2}(\mathbf{u}_x^n)_j$  and  $\mathbf{u}_{j-1/2}^{n,-} = \mathbf{u}_{j-1}^n + \frac{\Delta x}{2}(\mathbf{u}_x^n)_{j-1}$ . Setting  $(\mathbf{u}_x^n)_j = l_j \sigma_j^a$  for any  $j \in \mathbb{Z}$ , where  $\sigma_j^a := \frac{\mathbf{u}_{j+1}^n - \mathbf{u}_{j-1}^n}{2\Delta x}$  is the central slope and  $l_j \in [0, 1]$  is a slope limiter, we define the limited interface values to be

$$\mathbf{u}_{j-1/2}^{n,+}(l_j) = \mathbf{u}_j^n - \frac{\Delta x}{2} l_j \sigma_j^a, \quad \mathbf{u}_{j-1/2}^{n,-}(l_{j-1}) = \mathbf{u}_{j-1}^n + \frac{\Delta x}{2} l_{j-1} \sigma_{j-1}^a. \quad (4.21)$$

Similar to the flux limiting limitation, we need to find the largest values  $l_j^{z,-}$  and  $l_{j-1}^{z,+}$  such that for a given  $j$  both  $\mathbf{u}_{j-1/2}^{n,+}(l_j^{z,-})$  and  $\mathbf{u}_{j-1/2}^{n,-}(l_{j-1}^{z,+})$  are in  $A_{j-1/2}^z$ . This is described in the following lemma.

LEMMA 4.3. *Let's define  $l_j^{z,-}$  and  $l_{j-1}^{z,+}$  to be*

$$l_j^{z,-} = \begin{cases} 1 & \text{if } \Psi_{j-1/2}^z(\mathbf{u}_{j-1/2}^{n,+}(1)) \geq 0, \\ \max\{l_j \in [0, 1] \mid \Psi_{j-1/2}^z(\mathbf{u}_{j-1/2}^{n,+}(l_j)) \geq 0\} & \text{otherwise.} \end{cases} \quad (4.22)$$

$$l_{j-1}^{z,+} = \begin{cases} 1 & \text{if } \Psi_{j-1/2}^z(\mathbf{u}_{j-1/2}^{n,-}(1)) \geq 0, \\ \max\{l_{j-1} \in [0, 1] \mid \Psi_{j-1/2}^z(\mathbf{u}_{j-1/2}^{n,-}(l_{j-1})) \geq 0\} & \text{otherwise.} \end{cases} \quad (4.23)$$

Then for all  $l_j^{z,-} \in [0, l_j^{z,-}]$  and  $l_{j-1}^{z,+} \in [0, l_{j-1}^{z,+}]$ , it holds that  $\Psi_{j-1/2}^z(\mathbf{u}_{j-1/2}^{n,+}(l_j^{z,-})) \geq 0$  and  $\Psi_{j-1/2}^z(\mathbf{u}_{j-1/2}^{n,-}(l_{j-1}^{z,+})) \geq 0$ .

Let's denote  $\mathbf{u}_j^{n+1}(l_{j-1}, l_j, l_{j+1})$  to be the limited second-order update computed with interface values  $\mathbf{u}_{j-1/2}^{n,-}(l_{j-1}), \mathbf{u}_{j-1/2}^{n,+}(l_j), \mathbf{u}_{j+1/2}^{n,-}(l_j), \mathbf{u}_{j+1/2}^{n,+}(l_{j+1})$ . Note that we recover the first-order solution if  $l_{j-1} = l_j = l_{j+1} = 0$  and the second-order solution if  $l_{j-1} = l_j = l_{j+1} = 1$ . The goal is to find a set of local limiters which preserve the invariant domain property. A straightforward application of Theorem 4.1 and Lemma 4.3 gives the following.

LEMMA 4.4. *Let  $l_{j-1}^{z,+}, l_j^{z,-}$  be the slope limiters computed via Lemma 4.3 for any  $j \in \mathbb{Z}$  and  $n \geq 0$ . If we set  $l_j^z = \min(l_j^{z,-}, l_j^{z,+})$  for  $j \in \mathbb{Z}$ , then we have that  $\mathbf{u}_j^{n+1}(l_{j-1}^z, l_j^z, l_{j+1}^z) \in A_{j-1/2}^z \cup A_{j+1/2}^z$ .*

REMARK 4.5. Note that, the underlying assumption is that both sets  $A_{j-1/2}^z$  and  $A_{j+1/2}^z$  are in a local invariant set  $A$ . We now describe the slope limiting algorithm for all local constraints.

---

**Algorithm 2** Convex slope limiting

---

**Input:**  $\mathbf{u}_j^n, \mathbf{u}_{j-1/2}^{n,+}, \mathbf{u}_{j+1/2}^{n,-}, z_1, \dots, z_q$ .

**Output:**  $\mathbf{u}_j^{n+1}$

- 1: **for**  $z=1$  **to**  $z_q$  **do**
  - 2:     Compute limiting parameters  $l_j^z$  via Lemma 4.3 and Lemma 4.4.
  - 3: **end for**
  - 4: Set  $l_j := \min_{z \in \{z_1, \dots, z_q\}} l_j^z$  for all  $j \in \mathbb{Z}$ .
  - 5: Update  $\mathbf{u}_{j-1/2}^{n+1,+}$  and  $\mathbf{u}_{j+1/2}^{n+1,-}$ .
  - 6: Update  $\mathbf{u}_j^{n+1} = \mathbf{u}_j^{n+1}(l_{j-1}, l_j, l_{j+1})$ .
  - 7: **Return**  $\mathbf{u}_j^{n+1}$ .
-

In the case of two space dimensions, we limit the second-order solution via a similar approach. Consider the KT-scheme (3.20)-(3.21), where  $\mathbf{u}_{j+1/2,k}^{n,\pm}$  and  $\mathbf{u}_{j,k+1/2}^{n,\pm}$  are the interface values given by the local linear reconstructions. Using  $\lambda_{j+1/2,k}^{n,x}$  and  $\lambda_{j,k+1/2}^{n,y}$  instead of  $a_{j+1/2,k}^{n,x}$  and  $a_{j,k+1/2}^{n,y}$  respectively in (3.20), we have the following result.

**THEOREM 4.2.** *Let  $A$  be a convex invariant set of (2.1),  $n \geq 0$  and  $j, k \in \mathbb{Z}$  be such that  $\mathbf{u}_{j,k}^n$  and all interface values  $\mathbf{u}_{j\pm 1/2,k}^{n,\pm}$  and  $\mathbf{u}_{j,k\pm 1/2}^{n,\pm}$  are in  $A$ . Assume that the second-order solution  $\mathbf{u}_{j,k}^{n+1}$  is computed with the KT-scheme (3.20)-(3.21), where  $\lambda_{j-1/2,k}^{n,x} := a_{j-1/2,k}^{n,x} = \lambda_{\max}(\mathbf{u}_{j-1/2,k}^{n,-}, \mathbf{u}_{j-1/2,k}^{n,+}, \mathbf{f})$ ,  $\lambda_{j+1/2,k}^{n,x} := a_{j+1/2,k}^{n,x} = \lambda_{\max}(\mathbf{u}_{j+1/2,k}^{n,-}, \mathbf{u}_{j+1/2,k}^{n,+}, \mathbf{f})$ ,  $\lambda_{j,k-1/2}^{n,y} := a_{j,k-1/2}^{n,y} = \lambda_{\max}(\mathbf{u}_{j,k-1/2}^{n,-}, \mathbf{u}_{j,k-1/2}^{n,+}, \mathbf{g})$  and  $\lambda_{j,k+1/2}^{n,y} := a_{j,k+1/2}^{n,y} = \lambda_{\max}(\mathbf{u}_{j,k+1/2}^{n,-}, \mathbf{u}_{j,k+1/2}^{n,+}, \mathbf{g})$ . We set the in cell local speeds to be  $\lambda_{j,k,\pm}^{n,x} := \lambda_{\max}(\mathbf{u}_{j-1/2,k}^{n,+}, \mathbf{u}_{j+1/2,k}^{n,-}, \mathbf{f})$  and  $\lambda_{j,k,\pm}^{n,y} := \lambda_{\max}(\mathbf{u}_{j,k-1/2}^{n,+}, \mathbf{u}_{j,k+1/2}^{n,-}, \mathbf{g})$  and define the maximum local speeds in the  $x$ - and the  $y$ -direction respectively by*

$$\lambda_{j,k}^{\max,x} := \max(\lambda_{j-1/2,k}^{n,x}, \lambda_{j+1/2,k}^{n,x}, \lambda_{j,k,\pm}^{n,x}), \tag{4.24a}$$

$$\lambda_{j,k}^{\max,y} := \max(\lambda_{j,k-1/2}^{n,y}, \lambda_{j,k+1/2}^{n,y}, \lambda_{j,k,\pm}^{n,y}). \tag{4.24b}$$

Then under the CFL condition  $\frac{\Delta t}{\Delta x} \lambda_{j,k}^{\max,x} + \frac{\Delta t}{\Delta y} \lambda_{j,k}^{\max,y} \leq \frac{1}{4}$ , we have that  $\mathbf{u}_{j,k}^{n+1} \in A$ .

**REMARK 4.6.** The proof for Theorem 4.2 is analogous to the proof of Theorem 4.1 and we omit it. This result generalizes the scalar maximum principle proved in [19, Thm. 5.1] to an invariant domain property for an arbitrary hyperbolic system under the CFL condition

$$\max_{j,k} \left( \frac{\Delta t}{\Delta x} \lambda_{j,k}^{\max,x}, \frac{\Delta t}{\Delta y} \lambda_{j,k}^{\max,y} \right) \leq \frac{1}{8}. \tag{4.25}$$

Similar to the one dimensional case, we define the limited interface values by

$$\begin{aligned} \mathbf{u}_{j-1/2,k}^{n,+}(l_{j,k}^x) &= \mathbf{u}_{j,k}^n - \frac{\Delta x}{2} l_{j,k}^x \sigma_{j,k}^{a,x}, & \mathbf{u}_{j-1/2,k}^{n,-}(l_{j-1,k}^x) &= \mathbf{u}_{j-1,k}^n + \frac{\Delta x}{2} l_{j-1,k}^x \sigma_{j-1,k}^{a,x}, \\ \mathbf{u}_{j,k-1/2}^{n,+}(l_{j,k}^y) &= \mathbf{u}_{j,k}^n - \frac{\Delta y}{2} l_{j,k}^y \sigma_{j,k}^{a,y}, & \mathbf{u}_{j,k-1/2}^{n,-}(l_{j,k-1}^y) &= \mathbf{u}_{j,k-1}^n + \frac{\Delta y}{2} l_{j,k-1}^y \sigma_{j,k-1}^{a,y}. \end{aligned} \tag{4.26}$$

where  $\sigma_{j,k}^{a,x}$  and  $\sigma_{j,k}^{a,y}$  are two-dimensional central slopes defined by  $\sigma_{j,k}^{a,x} := \frac{\mathbf{u}_{j+1,k}^n - \mathbf{u}_{j-1,k}^n}{2\Delta x}$  and  $\sigma_{j,k}^{a,y} := \frac{\mathbf{u}_{j,k+1}^n - \mathbf{u}_{j,k-1}^n}{2\Delta y}$  and  $l_{j,k}^x, l_{j,k}^y \in [0, 1]$  are the to be computed slope limiters. For a given constraint  $z$ , let  $A_{j\pm 1/2,k}^{z,x}$  and  $A_{j,k\pm 1/2}^{z,y}$  be the local invariant sets at the interfaces of cell  $[x_{j-1/2}, x_{j+1/2}] \times [y_{k-1/2}, y_{k+1/2}]$  such that  $A_{j\pm 1/2,k}^{z,x}$  and  $A_{j,k\pm 1/2}^{z,y}$  are all in  $A$ . Using Lemma 4.3 and Lemma 4.4, we find largest positive  $l_{j,k}^x, l_{j,k}^y \in [0, 1]$  such that  $\mathbf{u}_{j-1/2,k}^{n,+}(l_{j,k}^x), \mathbf{u}_{j-1/2,k}^{n,-}(l_{j-1,k}^x) \in A_{j-1/2,k}^{z,x}$  and  $\mathbf{u}_{j,k-1/2}^{n,+}(l_{j,k}^y), \mathbf{u}_{j,k-1/2}^{n,-}(l_{j,k-1}^y) \in A_{j,k-1/2}^{z,y}$  for all  $j, k \in \mathbb{Z}$ .

**4.4. Application to the Euler system.** In this section, we will illustrate how to apply the two types of convex limiting processes stated in §4 for the Euler system of gas dynamics. More specifically, we explain how to apply the quasiconcave limitation from §4.1. The general approach we follow in all convex limitations is the one from [11, §4.1]. That is, we enforce the same type of limitations on the density and specific entropy. In all limitations, to reduce the computational cost, we apply one diagonal search first as described in Remark 4.3. We assume that the equation of state is a gamma-law, i.e., the specific entropy is  $s = \log(e^{\frac{1}{\gamma-1}} \rho^{-1})$ , and we impose the local invariant sets by enforcing the constraints (4.6).

**4.4.1. Flux limiting.** Using the same notation as in §4.1.3 and defining  $\Delta G_j^n = (\Delta G_j^{\rho,n}, \Delta G_j^{m,n}, \Delta G_j^{E,n})^\top := (G_{j+1/2}^n - G_{j-1/2}^n)$ , we perform the flux limitation process as follows:

- (1) We limit the density by setting  $\psi_j^1(l) := \Psi_j^1(\mathbf{u}_j^{n+1}(l)) = \rho_j^{n,\max} - \rho_j^{L,n+1} + l\Delta G_j^{\rho,n}$  and  $\psi_j^2(l) := \Psi_j^2(\mathbf{u}_j^{n+1}(l)) = \rho_j^{L,n+1} - l\Delta G_j^{\rho,n} - \rho_j^{n,\min}$ . We compute the limiters on density by

$$l_j^\rho = \begin{cases} \min\left(\frac{|\rho_j^{n,\max} - \rho_j^{L,n+1}|}{|\Delta G_j^{\rho,n}| + \epsilon_\rho}, 1\right), & \text{if } \psi_j^1(1) < 0, \\ 1 & \text{if } \psi_j^1(1) \geq 0 \text{ \& } \psi_j^2(1) \geq 0, \\ \min\left(\frac{|\rho_j^{L,n+1} - \rho_j^{n,\min}|}{|\Delta G_j^{\rho,n}| + \epsilon_\rho}, 1\right), & \text{if } \psi_j^2(1) < 0, \end{cases} \quad (4.27)$$

where we take  $\epsilon_\rho = \epsilon \rho_j^{n,\max}$ , where  $\epsilon = 10^{-16}$  to avoid division by zero.

- (2) For limitation on the specific entropy, we use  $\psi_j^3(l) := \Psi_j^3(l) = (\rho_j^{L,n+1} - l\Delta G_j^{\rho,n})(e_j^{L,n+1} - l\Delta G_j^{e,n}) - c_j^{n,\min}(\rho_j^{L,n+1} - l\Delta G_j^{\rho,n})^\gamma$  which is a concave down function of  $l$ . We define  $l_j^s$  as follows. If  $\psi_j^3(\min(1, l_j^\rho)) \geq 0$ , we take  $l_j^s = \min(1, l_j^\rho)$ ; if  $\psi_j^3(0) > 0$  and  $\psi_j^3(\min(1, l_j^\rho)) < 0$ , we define  $l_j^s$  to be the unique positive root of  $\psi_j^3(l) = 0$ ; if  $\psi_j^3(0) = 0$  and  $\psi_j^3(\min(1, l_j^\rho)) < 0$ , then  $\psi_j^3(l) = 0$  has exactly two roots and we take  $l_j^s$  to be the largest nonnegative root of  $\psi_j^3(l) = 0$ .

Let  $l_j = l_j^s$ , then it follows by Lemma 4.8 and Lemma 4.13 in [11] that the limited solution  $\mathbf{u}_j^{n+1}(l_j) := \mathbf{u}_j^{L,n+1} - l_j \lambda \Delta G_j^n$  satisfies the invariant domain property described in §4.1.3.

**4.4.2. Slope Limiting.** We start by constructing the local invariant constraints at cell interface  $x_{j-1/2}$ . We set

$$\begin{aligned} \rho_{j-1/2}^{n,\min} &:= \min(\rho_j^n, \rho_{j-1}^n, \bar{\rho}_{j-1/2}^{n+1}), & \rho_{j-1/2}^{n,\max} &:= \max(\rho_j^n, \rho_{j-1}^n, \bar{\rho}_{j-1/2}^{n+1}). \\ s_{j-1/2}^{n,\min} &:= \min(\Phi(\mathbf{u}_j^n), \Phi(\mathbf{u}_{j-1}^n), \Phi(\bar{\mathbf{u}}_{j-1/2}^{n+1})). \end{aligned} \quad (4.28)$$

Analogous to the flux limiting case, we set  $\Psi_{j-1/2}^1(\mathbf{u}) = \rho_{j-1/2}^{n,\max} - \rho$ ,  $\Psi_{j-1/2}^2(\mathbf{u}) = \rho - \rho_{j-1/2}^{n,\min}$ ,  $\Psi_{j-1/2}^3(\mathbf{u}) = \rho e - c_{j-1/2}^{n,\min} \rho^\gamma$ , where  $c_{j-1/2}^{n,\min} = \exp((\gamma-1)s_{j-1/2}^{n,\min})$ . We impose the invariant domain property by enforcing  $\Psi_{j-1/2}^1(\mathbf{u}), \Psi_{j-1/2}^2(\mathbf{u}), \Psi_{j-1/2}^3(\mathbf{u}) \geq 0$  on each interface. Using the same notation as in §4.3, we denote  $\sigma_j^a := (\sigma_j^{a,\rho}, \sigma_j^{a,m}, \sigma_j^{a,E})^\top$  to be the central slope and define

$$\mathbf{u}_{j+1/2}^-(l) := \mathbf{u}_j^n + \frac{\Delta x}{2} l \sigma_j^a, \quad \mathbf{u}_{j-1/2}^+(l) := \mathbf{u}_j^n - \frac{\Delta x}{2} l \sigma_j^a. \quad (4.29)$$

Depending on the sign of  $\sigma_j^{a,\rho}$ , we limit the density as follows:

- (1) If  $\sigma_j^{a,\rho} > 0$ , we set

$$l_j^\rho = \min\left(\frac{|\rho_{j+1/2}^{n,\max} - \rho_j^n|}{|\frac{\Delta x}{2} \sigma_j^{a,\rho}| + \epsilon_\rho^{n,+}}, \frac{|\rho_j^n - \rho_{j-1/2}^{n,\min}|}{|\frac{\Delta x}{2} \sigma_j^{a,\rho}| + \epsilon_\rho^{n,-}}, 1\right), \quad (4.30)$$

where  $\epsilon_\rho^{n,+} = \epsilon \rho_{j+1/2}^{n,\max}$  and  $\epsilon_\rho^{n,-} = \epsilon \rho_{j-1/2}^{n,\max}$ , with  $\epsilon = 10^{-16}$  to avoid division by zero.

- (2) If  $\sigma_j^{a,\rho} < 0$ , we set

$$l_j^\rho = \min\left(\frac{|\rho_j^n - \rho_{j+1/2}^{n,\min}|}{|\frac{\Delta x}{2} \sigma_j^{a,\rho}| + \epsilon_\rho^{n,+}}, \frac{|\rho_{j-1/2}^{n,\max} - \rho_j^n|}{|\frac{\Delta x}{2} \sigma_j^{a,\rho}| + \epsilon_\rho^{n,-}}, 1\right). \quad (4.31)$$

(3) If  $\sigma_j^{a,\rho} = 0$ , we take  $l_j^\rho = 1$ .

For limitation on the specific entropy, we enforce  $\Psi_{j+1/2}^3(\mathbf{u}_{j+1/2}^-(l)) \geq 0$  and  $\Psi_{j-1/2}^3(\mathbf{u}_{j-1/2}^+(l)) \geq 0$ . Let us denote  $\psi_j^+(l) = \Psi_{j+1/2}^3(\mathbf{u}_{j+1/2}^-(l))$  and  $\psi_j^-(l) = \Psi_{j-1/2}^3(\mathbf{u}_{j-1/2}^+(l))$ . Note that both of these functions are concave-down functions of  $l$ . Therefore, following the same approach as in §4.4.1, we compute  $l_j^{s,+}$  and  $l_j^{s,-}$ . Then, the slope limiter for the specific entropy is defined as  $l_j^s = \min(l_j^{s,+}, l_j^{s,-})$ . After all limitation, the following result holds.

LEMMA 4.5. *Let  $l_j = l_j^s$  for all  $j \in \mathbb{Z}$ , then for any  $l \in [0, l_j]$ , we have that  $\Psi_{j-1/2}^z(\mathbf{u}_{j-1/2}^+(l)) \geq 0$  and  $\Psi_{j+1/2}^z(\mathbf{u}_{j+1/2}^-(l)) \geq 0$ ,  $z = 1, 2, 3$ .*

**4.5. Local relaxations.** It has been observed in many instances that enforcing strict local bounds in the limitation process may reduce the accuracy of the method. In the scalar case it is well known that enforcing strict local maximum principle results in the so-called clipping phenomenon and the rate of error in  $L^\infty$  is reduced. In the case of systems, such effects can be observed when the local states are close to the boundary of the invariant domain. We refer the reader to [17] and [11] for further discussion on the relaxation bounds for the Euler system, in particular for relaxation of the minimum principle on the specific entropy. We follow that approach from [11, §4.7] which was originally proposed for the Euler system but could be easily applied in the general setting. We observe in all numerical tests that after relaxation the limited method keeps the accuracy of the unlimited method.

For simplicity we restrict ourselves to the case of one space dimension with the two dimensional case being analogous. Let  $\Omega$  be the computational domain in space and let  $h$  be the mesh size. That is, we take  $h := \frac{\Delta x}{|\Omega|}$  with  $|\Omega|$  being the diameter of the set. Let  $z$  denote the quantity to be limited. For example,  $z$  could be  $-w_1, w_2$  for P-system, see §4.1.2, or  $z$  could be  $\rho, -\rho, s$  for Euler system, see §4.1.3. We give two types of relaxations as follows.

(1) Limitation on a constraint  $z$  describing a smooth curved part of the boundary. For example,  $z = -w_1, w_2$  in the p-system, or  $z = s$  in the Euler system, let  $x_{ij} = \frac{1}{2}(x_i + x_j)$ , we define  $\Delta z_j^n := \max_{j \neq i \in \{j, j \pm 1\}} (z^n(x_{ij}) - z_j^{\min})$  and set

$$\overline{z_j^{\min,1}} := z_j^{\min} - \min(r_h |z_j^{\min}|, |\overline{\Delta z_j^n}|). \tag{4.32}$$

Then we will use  $\overline{z_j^{\min,1}}$  instead of  $z_j^{\min}$  to be the bound of the local invariant sets.

(2) Limitation on constraint  $z$  describing linear part on the boundary. For example,  $z = u$  or  $z = -u$  in scalar equations, or  $z = \pm \rho$  in the Euler system. Setting  $\Delta^2 z_j^n = z_{j-1}^n - 2z_j^n + z_{j+1}^n$ , we define  $\overline{\Delta^2 z_j^n} := \frac{1}{6} \sum_{j \neq i \in \{j, j \pm 1\}} (\frac{1}{2} \Delta^2 z_i^n + \frac{1}{2} \Delta^2 z_j^n)$  and  $\widetilde{\Delta^2 z_j^n} := \text{m}\{\frac{1}{2} \Delta^2 z_i^n | i \in \{j, j \pm 1\}\}$ , where  $\text{m}$  is the minmod operator defined in §3.3.1. The relaxed bound of local invariant set is defined as

$$\overline{z_j^{\min,2}} := z_j^{\min} - \min(r_h |z_j^{\min}|, |\overline{\Delta^2 z_j^n}|), \tag{4.33a}$$

$$\widetilde{z_j^{\min,2}} := z_j^{\min} - \min(r_h |z_j^{\min}|, |\widetilde{\Delta^2 z_j^n}|). \tag{4.33b}$$

We will use  $\overline{z_j^{\min,2}}$  or  $\widetilde{z_j^{\min,2}}$  instead of  $z_j^{\min}$  to be the bounds of local invariant sets. It is observed in the numerical tests that both of these two bounds defined in (4.33) are robust and give similar results.

REMARK 4.7. The relaxation technique on the bounds was first introduced in [11, §4.7] for the Euler system and later on generalized for general hyperbolic systems in [12, §7.6]. We simply apply this idea here in (4.32) and (4.33), with  $r_h = \min(1, h^{1.5})$  to restrict the relaxation order to  $\mathcal{O}(h^{1.5})$ . The exponent 1.5 seems to give good results and enforces the original invariant domain in the limit  $h \rightarrow 0$ , see [12, §7.6] for more details.

**5. Numerical illustrations**

In this section, we report numerical test to illustrate the performance of the two limiting techniques mentioned above. Our code is constructed using finite volume method on uniform cells of size  $\Delta x = h$  (in the 1d case) or  $\Delta x = \Delta y = h$  (in the 2d case). Time stepping is done by using SSP-RK3 methods (three stages, third-order), see [18, 22]. The time step is computed by using formula  $\tau = CFL \times \frac{h}{\lambda^{\max, n}}$ , where  $\lambda^{\max, n}$  is the global maximum speed of the method used at time  $t_n$ . For  $p \in [1, \infty]$ , we introduce a consolidated error indicator at time  $t$  by adding the relative error in  $L^p$ -norm of all conserved variables:

$$\delta_p^\alpha(t) = \sum_i \frac{\|\mathbf{u}_h^i(t) - \mathbf{u}^i(t)\|_{L^p(D)}}{\|\mathbf{u}^i(t)\|_{L^p(D)}}, \quad \mathbf{u} = (u^1, u^2, \dots, u^m), \tag{5.1}$$

where  $\alpha = f$  or  $s$ , corresponding to flux limiting or slope limiting error. For the Euler system, in all examples we use a  $\gamma$ -law equation of state, i.e.,  $p = (\gamma - 1)\rho e$ . We test the two new methods, MAPR-EV-CL and INV-CL, against the classical Minmod method (KT-scheme based on the minmod slope limiter given in (3.22)). Note that the MAPR-EV limiter (3.25) is only applied when  $\theta_j^n \leq 1.5$  and in the regions of smooth flow ( $\theta_j^n > 1.5$ ) we define  $\sigma^{\text{mapr}, \theta} = \frac{\mathbf{u}_{j+1} - \mathbf{u}_{j-1}}{2\Delta x}$ .

**5.1. Linear transport equation.** To illustrate the second-order accuracy of the newly developed schemes, we start with the one-dimensional linear transport equation,  $u_t + u_x = 0$ ,  $0 \leq x \leq 2\pi$ . The initial condition is  $u(x, 0) = \sin(x)$  and the exact solution is given by  $u(x, t) = \sin(x - t)$  assuming periodic boundary conditions. We run the test for  $0 \leq t \leq 0.5$  and the results are given in Table 5.1 and Table 5.2. For  $L^1$ -convergence, we observe that all three methods presented in the paper reach a second-order accuracy and relatively smaller errors are obtained using MAPR-EV-CL and INV-CL. For  $L^\infty$ -convergence, the convergence rate is around 2 if we use MAPR-EV-CL or INV-CL. However, if we take the Minmod limiter only, the rate seems to be  $4/3$  which is expected for the classical minmod-type limitation.

# of cells	Minmod limiter		MAPR-EV-CL		INV-CL	
	$\delta_1^f(t)$	rate	$\delta_1^f(t)$	rate	$\delta_1^s(t)$	rate
100	2.90E-03		1.05E-04		1.05E-04	
200	7.79E-04	1.90	2.62E-05	2.00	2.62E-05	2.00
400	2.04E-04	1.93	6.55E-06	2.00	6.56E-06	2.00
800	5.35E-05	1.93	1.64E-06	2.00	1.64E-06	2.00
1600	1.41E-05	1.93	4.09E-07	2.00	4.11E-07	2.00
3200	3.63E-06	1.95	1.02E-07	2.00	1.03E-07	2.00

TABLE 5.1. Linear equation,  $L^1$ -Convergence tests with  $CFL = 0.25$ .

**5.2. Burgers' equation.** We consider the one dimensional inviscid Burgers' equation,  $u_t + (\frac{u^2}{2})_x = 0$ , with an exact solution with limited smoothness,  $\partial_x u(\cdot, t)$  is with

# of cells	Minmod limiter		MAPR-EV-CL		INV-CL	
	$\delta_\infty^f(t)$	rate	$\delta_\infty^f(t)$	rate	$\delta_\infty^s(t)$	rate
100	1.27E-02		1.65E-04		1.65E-04	
200	5.55E-03	1.19	4.11E-05	2.00	4.11E-05	2.00
400	2.29E-03	1.27	1.03E-05	2.00	1.03E-05	2.00
800	9.33E-04	1.30	2.57E-06	2.00	2.57E-06	2.00
1600	3.83E-04	1.29	6.43E-07	2.00	6.43E-07	2.00
3200	1.55E-04	1.30	1.61E-07	2.00	1.61E-07	2.00

TABLE 5.2. Linear equation,  $L^\infty$ -Convergence tests with CFL=0.25.

# of cells	Minmod limiter		MAPR-EV-CL limiter		Slope limiter	
	$\delta_1^f(t)$	rate	$\delta_1^f(t)$	rate	$\delta_1^s(t)$	rate
100	1.17E-03		5.43E-04		6.47E-04	
200	4.24E-04	1.47	1.65E-04	1.72	2.30E-04	1.49
400	1.56E-04	1.44	5.33E-05	1.63	8.47E-05	1.44
800	5.93E-05	1.40	1.84E-05	1.53	3.28E-05	1.37
1600	2.28E-05	1.38	6.68E-06	1.46	1.28E-05	1.35
3200	8.89E-06	1.36	2.51E-06	1.41	5.10E-06	1.33
6400	3.48E-06	1.35	9.63E-07	1.38	2.03E-06	1.33

TABLE 5.3. 1D Burgers' equation, Convergence tests with CFL=0.25.

bounded variation:  $u(x,t)=0$  if  $x < 0.25$ ;  $u(x,t) = \frac{4x-1}{4t+1}$  if  $0.25 \leq x < 0.5+t$ ;  $u(x,t) = 1$  if  $0.5+t \leq x \leq 1$ . The computation is done for  $0 \leq t \leq 0.5$  and the results are reported for  $t=0.4$  in Table 5.3. We observe that using the method based on the MAPR limiter gives the optimal rate in  $L^1$ . This is a super-convergence effect that we observe for scalar equations. However, the methods based on the minmod and the invariant slope limiter have a convergence rate around  $4/3$  which is expected for a method based on mass lumping, see [5] for details. Moreover, the convex flux limiting process doesn't affect the convergence rate of the unlimited MAPR-EV method.

**5.3. KPP test case.** We consider the so-called KPP-test, a two dimensional scalar conservation equation with a non-convex flux, see [21, §5.3] for more details. This test checks if the high-order method has enough viscosity to resolve correctly the composite wave structure of the unique entropy solution, see Figure 5.1.

# of cells	Minmod limiter		MAPR-EV-CL		INV-CL	
	$\delta_1^f(t)$	rate	$\delta_1^f(t)$	rate	$\delta_1^s(t)$	rate
100	4.61E-03		2.26E-03		2.27E-03	
200	2.02E-03	1.19	9.59E-04	1.24	9.53E-04	1.25
400	8.72E-04	1.21	3.88E-04	1.31	3.91E-04	1.29
800	3.54E-04	1.30	1.51E-04	1.36	1.53E-04	1.35
1600	1.44E-04	1.30	6.00E-05	1.33	6.05E-05	1.34
3200	5.80E-05	1.31	2.42E-05	1.31	2.43E-05	1.32

TABLE 5.4. The p-system, expansion wave, Convergence tests with CFL=0.25.

$$\partial_t u + \partial_x \sin u + \partial_y \cos u = 0, u(x, y, 0) = \begin{cases} \frac{14\pi}{4}, & \text{if } \sqrt{x^2 + y^2} \leq 1, \\ \frac{\pi}{4}, & \text{otherwise.} \end{cases} \tag{5.2}$$

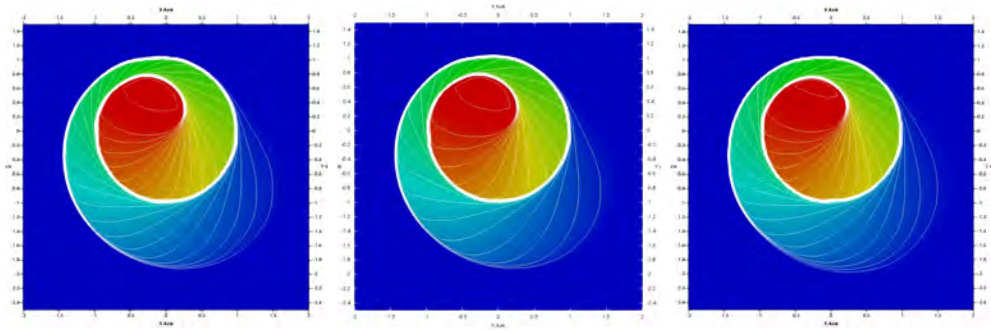


FIG. 5.1. *KPP-wave: CFL=0.25, t=1, 40000 cells. Left: Minmod limiter; Center: MAPR-EV-CL; Right: INV-CL.*

All schemes are able to resolve correctly the composite wave structure. Note that if we use the MAPR limiter with  $\theta=2$  the method will fail to converge to the correct solution, see [21] for details.

**5.4. The p-system.** We assume the pressure is given by  $p(v) = rv^{-\gamma}$  for the p-system (2.5). In the numerical example we take  $\gamma=3$ , and compute the Riemann problem with initial data  $(v_l, u_l) = (1, 0)$  and  $(v_r, u_r) = (2^{\frac{2}{\gamma-1}}, \frac{1}{\gamma-1})$ . The exact solution is a single rarefaction wave, see [11, §5] for more details on this test. The convergence rate for all methods is around  $\frac{4}{3}$ , see Table 5.4. The convex limiting process does not affect the rate of the unlimited MAPR-EV method.

**5.5. The Euler system, 1D smooth wave.** We start with a one-dimensional test whose purpose is to estimate the convergence rate of the methods on a very smooth solution. We set  $v(x, t) = 1$ ,  $p(x, t) = 1$  and

$$\rho(x, t) = \begin{cases} 1 + 2^6(x_1 - x_0)^{-6}(x - t - x_0)^3(x_1 - x + t)^3, & \text{if } x_0 \leq x - t < x_1, \\ 1 & \text{otherwise,} \end{cases} \quad (5.3)$$

where  $x_0 = 0.1$ ,  $x_1 = 0.3$  and  $\gamma = \frac{7}{5}$ . This is an exact solution for Euler, see [11, §5] for more details. The numerical solution is computed from  $t=0$  to  $t=0.1$ . The results are shown in Table 5.5.

# of cells	Minmod limiter		MAPR-EV-CL		INV-CL	
	$\delta_\infty^f(t)$	rate	$\delta_\infty^f(t)$	rate	$\delta_\infty^s(t)$	rate
100	1.53E-01		3.40E-02		2.75E-02	
200	6.64E-02	1.21	8.09E-03	2.07	6.68E-03	2.04
400	2.83E-02	1.23	2.45E-03	1.72	3.32E-03	1.01
800	1.17E-02	1.27	6.55E-04	1.90	1.15E-03	1.54
1600	4.78E-03	1.29	1.70E-04	1.94	3.44E-04	1.74
3200	1.93E-03	1.31	4.35E-05	1.97	9.20E-05	1.90

TABLE 5.5. *1D smooth wave, Convergence tests with CFL=0.25.*

**5.6. The Euler system, 1-rarefaction wave.** We consider the Riemann problem with the following initial data:  $(\rho_L, v_L, p_L) = (3, c_L, 1)$ ,  $(\rho_R, v_R, p_R) = (\frac{1}{2}, v_L +$



# of cells	Minmod limiter		MAPR-EV-CL		INV-CL	
	$\delta_1^f(t)$	rate	$\delta_1^f(t)$	rate	$\delta_1^s(t)$	rate
100	1.63E-02		1.63E-02		1.45E-02	
200	7.51E-03	1.12	3.96E-03	2.04	4.31E-03	1.75
400	3.15E-03	1.25	1.11E-03	1.83	1.31E-03	1.72
800	1.23E-03	1.36	3.37E-04	1.72	4.12E-04	1.67
1600	4.71E-04	1.38	1.09E-04	1.63	1.37E-04	1.59
3200	1.84E-04	1.36	3.59E-05	1.60	5.02E-05	1.45

TABLE 5.6. 1D Euler, 1-rarefaction wave, Convergence tests with CFL=0.25.

$\frac{2}{\gamma-1}(c_L - c_R), p_L(\frac{\rho_R}{\rho_L})^\gamma)$ , where  $c_L = \sqrt{\gamma p_L / \rho_L}$ ,  $c_R = \sqrt{\gamma p_R / \rho_R}$  and  $\gamma = \frac{7}{5}$ . The exact solution is described in [11, §5]. The numerical solution is computed starting from initial time  $t = \frac{0.2}{v_R - c_R}$  and running to final time  $t = 0.2$ . The results are given in Table 5.6. The  $L_1$ -convergence rate is best for the MAPR-EV-CL method.

**5.7. The Euler system, Leblanc shock tube.** We continue with a Riemann problem that is known in the literature as the Leblanc shocktube. The results are shown in Table 5.7. The performance of all methods is similar.

# of cells	Minmod limiter		MAPR-EV-CL		INV-CL	
	$\delta_\infty^f(t)$	rate	$\delta_\infty^f(t)$	rate	$\delta_\infty^s(t)$	rate
100	1.25E-01		1.24E-01		1.24E-01	
200	8.92E-02	0.49	8.04E-02	0.62	7.98E-02	0.64
400	5.79E-02	0.62	5.18E-02	0.64	4.91E-02	0.70
800	3.27E-02	0.83	2.75E-02	0.91	2.52E-02	0.96
1600	1.84E-02	0.83	1.44E-02	0.93	1.34E-02	0.91
3200	9.30E-03	0.98	7.57E-03	0.93	6.63E-03	1.02

TABLE 5.7. 1D Euler, Leblanc shocktube, Convergence tests with CFL=0.25.

**5.8. The Euler system, blast wave.** We consider the well known Woodward-Collela blast wave. The computations are done on the domain  $D = (0, 1)$  with CFL=0.25. The final time is  $t = 0.038$ . The results are shown in Figure 5.2. The MAPR-EV-CL method has the best resolution of the contact wave.

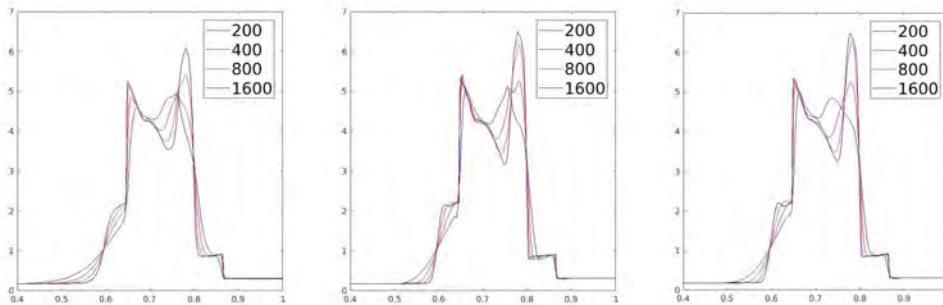


FIG. 5.2. Blast wave,  $t = 0.038$ , CFL=0.25. Left: Minmod limiter; Center: MAPR-EV-CL; Right: INV-CL.

**5.9. The Euler system, Isentropic Vortex.** We consider a two-dimensional problem introduced in [23]. The flow field is isentropic and the solution is smooth. Let  $\rho_\infty = P_\infty = T_\infty = 1$ ,  $\mathbf{u}_\infty = (1, 1)^\top$  be free stream values. We define the following perturbation values for the velocity and the temperature:

$$\delta \mathbf{u}(\mathbf{x}, t) = \frac{\beta}{2\pi} \exp\left(\frac{1-r^2}{2}\right)(-\bar{x}_2, \bar{x}_1), \quad \delta T(\mathbf{x}, t) = \frac{(\gamma-1)\beta^2}{8\gamma\pi^2} \exp(1-r^2), \quad (5.4)$$

where  $\beta = 5$  is a constant defining the vortex strength,  $\gamma = \frac{7}{5}$ ,  $\bar{\mathbf{x}} = (x_1 - x_1^c(t), x_2 - x_2^c(t))$ , where  $\mathbf{x}^c(t) = (x_1^0 + t, x_2^0)$  is the position of the vortex, and  $r^2 = \|\bar{\mathbf{x}}\|_{\ell^2}^2$ . The exact solution is a passive convection of the vortex with the mean velocity  $\mathbf{u}_\infty$ :

$$\rho(\mathbf{x}, t) = (T_\infty + \delta T)^{1/(\gamma-1)}, \quad \mathbf{u}(\mathbf{x}, t) = \mathbf{u}_\infty + \delta \mathbf{u}, \quad p(\mathbf{x}, t) = \rho^\gamma. \quad (5.5)$$

We perform the numerical computations in the rectangle  $D = (0, 20) \times (0, 20)$  from  $t=0$  until  $t=2$ , and we take  $x_1^0 = x_2^0 = 10$ . The results are shown in Table 5.8 and Figure 5.3. In this test it is critical to use local relaxation in the convex limitation process, see Section 4.5, to achieve the optimal convergence order. Both the MAPR-EV-CL and the INV-CL methods are optimal in this case.

# of cells	Minmod limiter		MAPR-EV-CL		INV-CL	
	$\delta_\infty^f(t)$	rate	$\delta_\infty^f(t)$	rate	$\delta_\infty^s(t)$	rate
2500	2.66E-01		1.25E-01		8.66E-02	
10000	1.19E-01	1.17	1.85E-02	2.75	1.85E-02	2.22
40000	5.82E-02	1.03	3.57E-03	2.38	3.57E-03	2.38
160000	2.94E-02	0.99	7.08E-04	2.34	7.08E-04	2.34

TABLE 5.8. *Isentropic vortex test case, Convergence tests with CFL=0.25.*

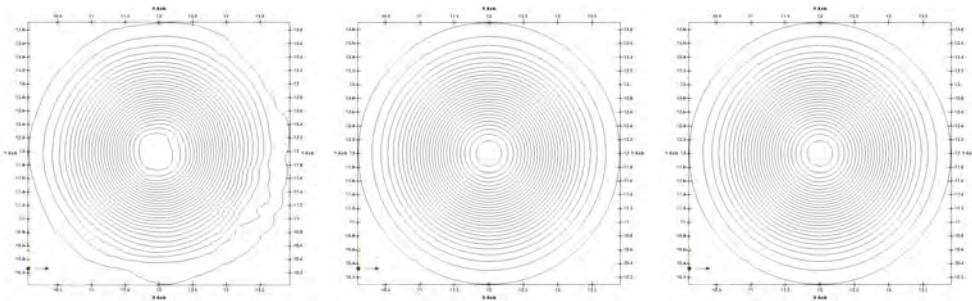


FIG. 5.3. *Isentropic vortex at t=2, CFL=0.25. Left: Mimmod limiter; Center: MAPR-EV-CL; Right: INV-CL.*

**5.10. The Euler system, Mach 3 Test.** Now we consider the classical Mach 3 flow in a wind tunnel with a forward facing step. The computational domain is  $D = (0, 3) \times (0, 1) \setminus (0.6, 3) \times (0, 0.2)$ ; the geometry of the domain is shown in Figure 5.4. The initial data is  $\rho = 1.4$ ,  $p = 1$ ,  $\mathbf{v} = (3, 0)^\top$ . The inflow boundary conditions are  $\rho|_{\{x=0\}} = 1.4$ ,  $p|_{\{x=0\}} = 1$ ,  $\mathbf{v}|_{\{x=0\}} = (3, 0)^\top$  and at the outflow boundary,  $\{x=3\}$ , we do nothing. On the top and bottom boundaries of the channel we enforce  $\mathbf{v} \cdot \mathbf{n} = 0$ . The computational results at  $t=4$  are shown in Figure 5.4. The MAPR-EV scheme (center)

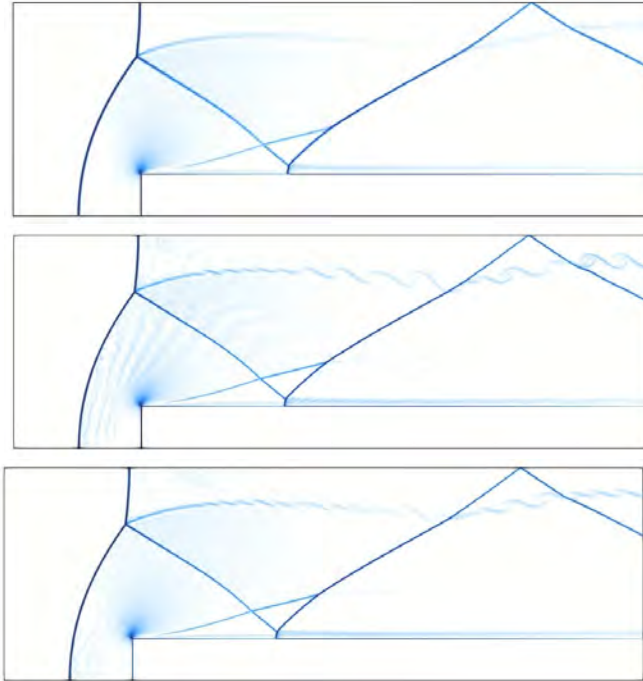


FIG. 5.4. Mach 3 step, density at  $t=4$ ,  $CFL=0.25$ . Top: Mimmod limiter; Center: MAPR-EV-CL; Bottom: INV-CL;

requires convex limiting to run and the results are superior in the region of the contact wave. The INV-CL scheme (bottom) has some instability in the contact but it is less pronounced at this mesh size. The minmod method (top) is the most diffusive scheme in this case, see Figure 5.4.

#### REFERENCES

- [1] J.P. Boris and D.L. Book, *Flux-corrected transport. I. SHASTA, a fluid transport algorithm that works*, J. Comput. Phys., **11(1)**:38–69, 1973. 4.2
- [2] I. Christov and B. Popov, *New non-oscillatory central schemes on unstructured triangulations for hyperbolic systems of conservation laws*, J. Comput. Phys., **227(11)**:5736–5757, 2008. 1, 3.3.1, 3.3.1
- [3] K.N. Chueh, C.C. Conley, and J.A. Smoller, *Positively invariant regions for systems of nonlinear diffusion equations*, Indiana Univ. Math. J., **26(2)**:373–392, 1977. 2.1, 2.2
- [4] H. Frid, *Maps of convex sets and invariant regions for finite-difference systems of conservation laws*, Arch. Ration. Mech. Anal., **160(3)**:245–269, 2001. 2.2, 2.2.3
- [5] J.-L. Guermond and R. Pasquetti, *A correction technique for the dispersive effects of mass lumping for transport problems*, Comput. Meth. Appl. Mech. Engrg., **253**:186–198, 2013. 5.2
- [6] 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. 2, 2.1, 2.2, 2.2, 2.2.2, 2.2.3, 3, 3.1, 3.1, 3.1, 3.1, 3.2, 3.3.2
- [7] J.-L. Guermond and B. Popov, *Invariant domains and second-order continuous finite element approximation for scalar conservation equations*, SIAM J. Numer. Anal., **55(6)**:3120–3146, 2017. 3.1, 3.3.2
- [8] J.-L. Guermond, R. Pasquetti, and B. Popov, *Entropy viscosity method for nonlinear conservation laws*, J. Comput. Phys., **230(11)**:4248–4267, 2011. 3.3.2

- [9] J.-L. Guermond, M. Nazarov, B. Popov, and Y. Yang, *A second-order maximum principle preserving Lagrange finite element technique for nonlinear scalar conservation equations*, SIAM J. Numer. Anal., **52(4):2163–2182**, 2014. [3.3.2](#)
- [10] J.-L. Guermond, B. Popov, L. Saavedra, and Y. Yang, *Invariant domains preserving arbitrary Lagrangian Eulerian approximation of hyperbolic systems with continuous finite elements*, SIAM J. Sci. Comput., **39(2):A385–A414**, 2017. [3.1](#)
- [11] 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. [1](#), [3.3.2](#), [3.3.2](#), [4](#), [4.1](#), [4.2](#), [4.2.1](#), [4.1](#), [4.4](#), [4.4.1](#), [4.5](#), [4.7](#), [5.4](#), [5.5](#), [5.6](#)
- [12] J.-L. Guermond, B. Popov, and I. Tomas, *Invariant domain preserving discretization-independent schemes and convex limiting for hyperbolic systems*, Comput. Meth. Appl. Mech. Engrg., **347:143–175**, 2019. [1](#), [3.3.2](#), [4.7](#)
- [13] A. Harten and S. Osher, *Uniformly high-order accurate nonoscillatory schemes. I*, SIAM J. Numer. Anal., **24(2):279–309**, 1987. [3.3.1](#)
- [14] D. Hoff, *A finite difference scheme for a system of two conservation laws with artificial viscosity*, Math. Comput., **33(148):1171–1193**, 1979. [2.2.2](#)
- [15] D. Hoff, *Invariant regions for systems of conservation laws*, Trans. Amer. Math. Soc., **289(2):591–610**, 1985. [2.2](#), [2.2.2](#)
- [16] G.-S. Jiang and E. Tadmor, *Nonoscillatory central schemes for multidimensional hyperbolic conservation laws*, SIAM J. Sci. Comput., **19(6):1892–1917**, 1998. [3.3.1](#)
- [17] B. Khobalatte and B. Perthame, *Maximum principle on the entropy and second-order kinetic schemes*, Math. Comput., **62(205):119–131**, 1994. [4.5](#)
- [18] J.F.B.M. Kraaijevanger, *Contractivity of Runge-Kutta methods*, BIT Numer. Math., **31(3):482–528**, 1991. [5](#)
- [19] A. Kurganov and E. Tadmor, *New high-resolution central schemes for nonlinear conservation laws and convection-diffusion equations*, J. Comput. Phys., **160(1):241–282**, 2000. [1](#), [3](#), [3.1](#), [3.2](#), [3.2](#), [3.3.1](#), [4.6](#)
- [20] A. Kurganov, S. Noelle, and G. Petrova, *Semidiscrete central-upwind schemes for hyperbolic conservation laws and Hamilton-Jacobi equations*, SIAM J. Sci. Comput., **23(3):707–740 (electronic)**, 2001. [1](#), [3](#)
- [21] A. Kurganov, G. Petrova, and B. Popov, *Adaptive semidiscrete central-upwind schemes for nonconvex hyperbolic conservation laws*, SIAM J. Sci. Comput., **29(6):2381–2401**, 2007. [1](#), [3.3.1](#), [5.3](#), [5.3](#)
- [22] C.-W. Shu and S. Osher, *Efficient implementation of essentially non-oscillatory shock-capturing schemes*, J. Comput. Phys., **77(2):439–471**, 1988. [5](#)
- [23] H.C. Yee, N.D. Sandham, and M.J. Djomehri, *Low-dissipative high-order shock-capturing methods using characteristic-based filters*, J. Comput. Phys., **150(1):199–238**, 1999. [5.9](#)
- [24] S.T. Zalesak, *Fully multidimensional flux-corrected transport algorithms for fluids*, J. Comput. Phys., **31(3):335–362**, 1979. [4.2](#)
- [25] X. Zhang and C.-W. Shu, *Maximum-principle-satisfying and positivity-preserving high-order schemes for conservation laws: survey and new developments*, Proc. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci., **467(2134):2752–2776**, 2011. [1](#)