



**HAL**  
open science

## Discrete entropy inequalities via an optimization process

Nina Aguillon, Emmanuel Audusse, Vivien Desveaux, Julien Salomon

► **To cite this version:**

Nina Aguillon, Emmanuel Audusse, Vivien Desveaux, Julien Salomon. Discrete entropy inequalities via an optimization process. *ESAIM: Mathematical Modelling and Numerical Analysis*, 2024, 58 (1), pp.363-391. 10.1051/m2an/2023098 . hal-04482948

**HAL Id: hal-04482948**

**<https://hal.science/hal-04482948>**

Submitted on 28 Feb 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## DISCRETE ENTROPY INEQUALITIES VIA AN OPTIMIZATION PROCESS

NINA AGUILLON<sup>1,\*</sup>, EMMANUEL AUDUSSE<sup>2</sup>, VIVIEN DESVEAUX<sup>3</sup> AND JULIEN SALOMON<sup>4</sup>

**Abstract.** The solutions of hyperbolic systems may contain discontinuities. These weak solutions verify not only the original PDEs, but also an entropy inequality that acts as a selection criterion determining whether a discontinuity is physical or not. Obtaining a discrete version of these entropy inequalities when approximating the solutions numerically is crucial to avoid convergence to unphysical solutions or even instability. However such a task is difficult in general, if not impossible for schemes of order 2 or more. In this paper, we introduce an optimization framework that enables us to quantify a posteriori the decrease or increase of entropy of a given scheme, locally in space and time. We use it to obtain maps of numerical diffusion and to prove that some schemes do not have a discrete entropy inequality. A special attention is devoted to the widely used second order MUSCL scheme for which almost no theoretical results are known.

**Mathematics Subject Classification.** 35L03, 65M08, 76M12, 35L40, 76-10.

Received May 19, 2023. Accepted November 30, 2023.

### 1. INTRODUCTION

Many physical phenomena can be described by means of a hyperbolic system, also called a system of conservation laws. Some famous hyperbolic systems include the Lighthill–Whitham–Richards or Aw–Rascle models for traffic flow, the shallow water equations of Barr de Saint-Venant, the Euler equation for fluid dynamics and the inviscid magnetohydrodynamics equation.

This class of partial differential equations (PDEs) does not contain any regularization terms such as diffusion or dispersion. Their solutions typically develop discontinuities in finite time. These discontinuities are observed in traffic jams, during floods caused by dam breaks, at hydraulic jumps, or in aeronautics. The PDE should be understood in the weak sense to allow such discontinuous solutions. Doing so makes it possible to construct infinitely many discontinuous solutions for the same initial data. An additional criterion should consequently be included to select only the physical weak solution. It generally takes the form of an entropy (or energy) inequality and is related to the second law of thermodynamics.

---

*Keywords and phrases.* Numerical diffusion, finite volume methods, discrete entropy inequality, MUSCL scheme.

<sup>1</sup> Sorbonne-Université and INRIA Paris, CNRS, Université de Paris, Laboratoire Jacques-Louis Lions (LJLL), 75005 Paris, France.

<sup>2</sup> Université Sorbonne Paris Nord and INRIA Paris, CNRS, Laboratoire Analyse, Géométrie et Applications (LAGA), 99 av. J.-B. Clément, 93430 Villetaneuse, France.

<sup>3</sup> Université de Picardie Jules Verne, CNRS, LAMFA, 33 rue Saint-Leu, 80039 Amiens Cedex 1, France.

<sup>4</sup> INRIA Paris, ANGE Project-Team, 75589 Paris Cedex 12, France and Laboratoire Jacques-Louis Lions, Sorbonne Université, CNRS, 75005 Paris, France.

\*Corresponding author: [nina.aguillon@sorbonne-universite.fr](mailto:nina.aguillon@sorbonne-universite.fr)

Discretizing the PDE to obtain a numerical approximation of the solution can be done in several ways. In this paper we focus on finite volume schemes which are well adapted to the low regularity of the solution and built around the idea of conservation laws. In the design of such schemes, it seems important that the entropy also decreases at the numerical level. This condition ensures that the scheme will not converge towards a nonphysical solution of the PDE. The loss of entropy in each cell during one time step is called the numerical diffusion. Discrete entropy inequalities have mainly been obtained for first order schemes in space and time. Realistic codes use high-order discretizations and splitting techniques and usually incorporate ideas and knowledge from the well-understood first-order framework. For this reason, they probably verify a discrete entropy inequality in most cases. However, no explicit formulas are known in practice.

This paper proposes to quantify the numerical diffusion with an a posteriori minimization technique where the scheme is used as a black box. Our primary goal is to obtain maps of numerical diffusion which quantify in space and time the loss of entropy coming from the choice of discretization. As a practical motivation and a long term objective, let us mention numerical oceanic circulation models where the numerical diffusion is related to artificial changes of salinity, density and temperature between two adjacent distinct water masses, and is usually called spurious mixing. This phenomenon is identified as a major issue in numerical cores for climate application [13]. We refer the reader to [5, 18] for the quantification of spurious mixing in simplified configurations and [17] in a realistic setting.

This paper provides a mathematical insight on the quantification of numerical diffusion in realistic codes but is still far away from oceanic applications. Maps of numerical diffusion are obtained by minimizing a functional which takes into account the consistency of the numerical entropy fluxes and the fact that the entropy should decrease at each time step. This minimization couples all the cells of the mesh, but we also propose a local and cheap quantification that gives qualitatively good results. A different perspective allows us to construct the worst initial data in terms of entropy by minimizing a different but related functional. We apply this procedure to prove that no discrete entropy inequality exists for most of the versions of the widely used MUSCL approach with a 2 steps Runge-Kutta time discretization. This was suspected in [2]. A notable exception is the limitation in the entropy variables with a HLL first order scheme [1].

We now present the mathematical framework on hyperbolic PDEs and their discretization with finite volume schemes, with an emphasis on discrete entropy inequality which is at the core of this work. Our main results and the outline of the paper are described hereafter.

## Fundamentals on discrete entropy inequalities

Consider a hyperbolic system in 1 dimension (1D) in space

$$\partial_t u(x, t) + \partial_x f(u(x, t)) = 0, \quad x \in \mathbf{R}, t \in \mathbf{R}^+, \quad (1)$$

where the vectorial unknown  $u$  belongs to some convex domain  $\Omega \subset \mathbf{R}^d$ . The flux  $f : \Omega \rightarrow \mathbf{R}^d$  is  $C^1$ -regular and its Jacobian matrix  $Df$  is diagonalizable with real eigenvalues. We are only interested in weak solutions of (1) that additionally satisfy the entropy inequality (or energy inequality)

$$\partial_t \eta(u) + \partial_x G(u) \leq 0, \quad (2)$$

where the entropy  $\eta : \Omega \rightarrow \mathbf{R}$  is strictly convex. The entropy flux  $G$  is linked to the entropy  $\eta$  through the relation on their Jacobian matrices  $D\eta Df = DG$ . Such hyperbolic systems arise in particular in the modeling of nonviscous flows. In this paper, we consider the scalar ( $d = 1$ ) Burgers equations related to the Lighthill–Whitham–Richards model for traffic flows and the Euler equations of inviscid gas dynamics for which  $d = 3$ .

We now turn to the numerical discretization of (1) with a finite volume technique. For the sake of simplicity a space interval  $[a, b]$  with periodic boundary conditions is considered throughout the paper. A finite number of cells  $M$  is fixed with size  $\Delta x = \frac{b-a}{M}$  and we note  $x_{j-1/2} = a + (j-1)\Delta x$  for  $j \in \llbracket 1, M \rrbracket$ . The points  $a = x_{1/2}$  and  $b = x_{M+1/2}$  are identified and all the space subscript are considered modulo  $M$ . For the sake of simplicity, every family  $(X)_{j \in \llbracket 1, M \rrbracket}$  of quantities indexed by the set of cells will be denoted by  $(X_j)_j$ .

We also consider a discretization in time  $0 = t^0 < t^1 < \dots t^n < \dots$ . A Courant–Friedrichs–Lewy (CFL) condition is imposed at each time step. It reads, for some CFL  $\alpha \in (0, 1)$  depending on the scheme,

$$(t^{n+1} - t^n) \max_{1 \leq j \leq M} \rho(Df(u_j^n)) = \alpha \Delta x, \tag{3}$$

where  $\rho(DF(u_j^n))$  is the spectral radius of the Jacobian matrix  $DF(u_j^n)$ . The time step varies with the time iteration but for the sake of simplicity, we will denote it independently of  $n$  as  $t^{n+1} - t^n = \Delta t$ .

A finite volume scheme writes

$$u_j^{n+1} = u_j^n - \frac{\Delta t}{\Delta x} (f_{j+1/2}^n - f_{j-1/2}^n), \tag{4}$$

where the vectors  $u_j^n$  and  $f_{j+1/2}^n$  correspond to the numerical approximations of the mean of the exact solution  $u$  and flux  $f(u)$

$$u_j^n \approx \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u(x, t^n) dx, \quad f_{j+1/2}^n \approx \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} f(u(x_{j+1/2}, s)) ds.$$

**Definition 1.** A consistent finite volume scheme with a stencil of  $s_L \in \mathbf{N}$  cells to the left and  $s_R \in \mathbf{N}$  cells to the right is a formula expressing the numerical flux  $f_{j+1/2}^n$  in terms of its  $s_L + s_R$  neighbor cells  $f_{j+1/2}^n = \mathcal{F}(u_{j-s_L+1}^n, \dots, u_{j+s_R}^n)$  such that, for all  $u \in \Omega$ ,  $\mathcal{F}(u_{-s_L+1}, \dots, u_{s_R})$  tends to  $f(u)$  as  $(u_{-s_L+1}, \dots, u_{s_R})$  tends to  $(u, \dots, u)$ .

Equations (1) and (2) are understood in a weak sense to allow discontinuous solutions. Inequality (2) does not hold for every discontinuity but selects only entropy satisfying shocks. The fact that the scheme (4) is stable and computes the entropy solution, with only entropy satisfying shocks, is strongly related to the existence of a numerical counterpart of (2) at the discrete level, called a discrete entropy inequality

$$\forall j \in \llbracket 1, M \rrbracket, \quad \eta(u_j^{n+1}) \leq \eta(u_j^n) - \frac{\Delta t}{\Delta x} (G_{j+1/2}^n - G_{j-1/2}^n). \tag{5}$$

In the spirit of Definition 1, we introduce the notion of entropy satisfying scheme.

**Definition 2.** Consider a consistent finite volume scheme of  $s_L \in \mathbf{N}$  cells to the left and  $s_R \in \mathbf{N}$  to the right. It is a consistent entropy satisfying scheme if there exists a numerical entropy flux function  $\mathcal{G}$  such that the following two conditions are satisfied.

- Inequality (5) holds for all data  $(u_k^n)_k$  with the choice  $G_{j+1/2}^n = \mathcal{G}(u_{j-s_L+1}^n, \dots, u_{j+s_R}^n)$ . The remainder

$$D_j^n = \eta(u_j^{n+1}) - \eta(u_j^n) + \frac{\Delta t}{\Delta x} (G_{j+1/2}^n - G_{j-1/2}^n) \tag{6}$$

is called numerical diffusion and is nonpositive in all cells.

- The numerical entropy flux function  $\mathcal{G}$  is consistent: for all  $u \in \Omega$ ,  $\mathcal{G}(u_{-s_L+1}, \dots, u_{s_R})$  tends to  $G(u)$  as  $(u_{-s_L+1}, \dots, u_{s_R})$  tends to  $(u, \dots, u)$ .

There exist several first order schemes (4) with a stencil  $s_L = s_R = 1$  for which an explicit formula for  $\mathcal{G}$  yielding to nonpositive diffusion  $D_j^n$  is known:

- for scalar equations  $d = 1$ , monotone schemes, see [14];
- for hyperbolic systems  $d \geq 2$  the Godunov and HLL schemes are entropy satisfying; see [14]. In the specific case of gas dynamics, the HLLC scheme and some relaxation or kinetic schemes are also entropy satisfying see [3, 20] and references therein.

For schemes of order larger than 1 the specific form of (5) seems out of reach for hyperbolic systems and the question is still largely open. Some works present results in that direction, mainly for second order schemes.

- In [4, 9] and [1], inequality (5) is slightly modified. The schemes are either difficult to implement or there is no guarantee that they capture entropy solution.
- The local discrete entropy inequality (5) is replaced by the weaker global condition  $\sum_j \eta(u_j^{n+1}) \leq \sum_j \eta(u_j^n)$  in [10] for the multilayer shallow water equations, in [20] and [16] for gas dynamics and in [6] for a conservation law with nonconvex flux.
- A different approach consists in using a second order scheme and to go back to first order if (5) does not hold. The MOOD technique (see [7, 11]) was initially developed to ensure the positiveness of some quantities like the density and the pressure, as well as some discrete maximum principle. This method was later extended in [2] to ensure the discrete entropy inequalities (5) hold. However it is limited to the gas dynamics (23).

On the other hand many schemes of order 2 or more are employed in the applications for their realistic results. They are typically designed to be of high order when the solution is smooth and include corrections to ensure stability, such as limiters or explicit numerical diffusion. Positiveness and lack of spurious oscillations are also often taken into account. However, the existence of discrete entropy inequality often remains an open question, as is the case for the Piecewise Parabolic Method (PPM) [8] or the (Weighted) Essentially Non Oscillatory method (ENO/WENO) [15, 21, 26].

In this work, we are concerned with the *a posteriori* determination of the numerical entropy fluxes  $G_{j+1/2}^n$  and numerical diffusions  $D_j^n$  of Definition 2. We do not attempt to find an analytical choice of  $\mathcal{G}$ . Instead, we follow a different approach and construct an adequate functional  $\mathcal{J}$  that is minimized to obtain

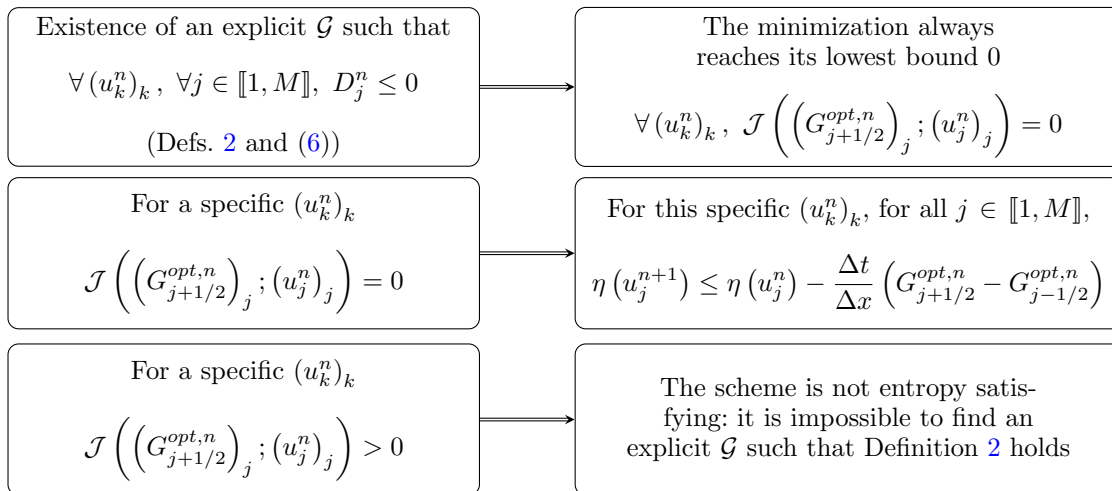
$$\left(G_{j+1/2}^{opt,n}\right)_j = \arg \min_{\gamma = (\gamma_{j+1/2}^n)_j \in \mathbf{R}^M} \left\{ \mathcal{J} \left( \gamma; (u_j^n)_j \right) \right\} \quad (7)$$

and the corresponding a posteriori numerical diffusion

$$D_j^{opt,n} = \eta(u_j^{n+1}) - \eta(u_j^n) + \frac{\Delta t}{\Delta x} \left( G_{j+1/2}^{opt,n} - G_{j-1/2}^{opt,n} \right). \quad (8)$$

The first objective is to be able to quantify the numerical diffusion in settings where finding an explicit formula for  $\mathcal{G}$  is out of reach. The functional is built so that the two constraints of Definition 2, namely the nonpositivity of the diffusion and the consistency of the fluxes, are satisfied “as much as possible” and can be applied to any explicit finite volume scheme satisfying Definition 1.

This a posteriori approach offers a novel method to construct numerical entropy fluxes. There are strong links between Definition 2 of entropy satisfying scheme and the fact that the minimum of the functional  $\mathcal{J}$  is 0. Our main results are summarized in the diagram below.



The paper is organized as follows. The construction of the functional  $\mathcal{J}$  and the minimization procedure are introduced in Section 2. In Section 3, we prove several theoretical results about the minimizers of the functional  $\mathcal{J}$ , including the two first implications of the previous diagram. The first one ensures that the minimization procedure is efficient on the class of entropy satisfying schemes where an expression for the numerical entropy fluxes  $G_{j+1/2}$  in (5) can be found in the literature. The second one is more practical, since it is applied on a specific initial data for a generic scheme and allows to visualize the distribution of numerical diffusion by making the best choice for  $G_{j+1/2}$ . We also prove a Lax-Wendroff theorem adapted to our minimization procedure in the sense that if a numerical scheme converges and we can find numerical entropy fluxes such that  $\mathcal{J}\left(\left(G_{j+1/2}^{opt,n}\right)_j; \left(u_j^n\right)_j\right) = 0$ , then the limit is an entropy solution.

The minimization procedure is illustrated on various testcases in Section 4 both for the Burgers equation and for the Euler equations of gas dynamics. The third implication of the previous diagram is nothing but the contraposition of the first one. We exploit it in Section 5 to build initial data that cannot satisfy any discrete entropy inequality. As an illustration, we prove that most variants of the MUSCL scheme with a second-order Runge-Kutta time discretization for the gas dynamics are not entropy satisfying.

## 2. MINIMIZATION PROCEDURE

The minimization is carried out at some fixed time iteration  $n$  which is referred to as “the initial data”. In other words  $\left(u_j^n\right)_j$  is fixed by the user and  $u_j^{n+1}$  is then obtained with Scheme (4). The superscript  $n$  plays no particular role and one can fix, e.g.,  $n = 0$ .

In this section, the numerical entropy flux is computed by a minimization procedure (7). The function  $\mathcal{J} : \mathbf{R}^M \rightarrow \mathbf{R}$  is defined by the sum of two contributions  $\mathcal{J} = \mathcal{J}^D + \mathcal{J}^C$ . The first part  $\mathcal{J}^D$  gathers the contributions of positive numerical diffusion

$$\mathcal{J}^D(\gamma; (u_j^n)_j) = \sum_{j=1}^M \max\left(0, \eta(u_j^{n+1}) - \eta(u_j^n) + \frac{\Delta t}{\Delta x} (\gamma_{j+1/2}^n - \gamma_{j-1/2}^n)\right)^2.$$

Indeed the quantity  $\eta(u_j^{n+1}) - \eta(u_j^n) + \frac{\Delta t}{\Delta x} (\gamma_{j+1/2}^n - \gamma_{j-1/2}^n)$  corresponds to the numerical diffusion (6) in the  $j$ -th cell if the numerical entropy fluxes are given by  $G_{j+1/2} = \gamma_{j+1/2}$ . A negative value is consistent with the Definition 2 of an entropy satisfying scheme; thus only positive contributions are kept in the function  $\mathcal{J}^D$ . Overall, we have  $\mathcal{J}^D(\gamma; (u_j^n)_j) = 0$  if and only if Inequality (5) holds with the choice  $G = \gamma$ .

It remains to take into account the consistency property of the numerical entropy fluxes of Definition 2, which is the role of the second part of the functional. We make use of *a priori* consistency bounds on the numerical entropy fluxes  $m_{j+1/2}^n \leq G_{j+1/2}^n \leq M_{j+1/2}^n$  depending on  $s_L$  cells on the left and  $s_R$  on the right. These bounds are defined in Lemma 1 and satisfy the following consistency property

$$u_{j-s_L+1}^n = \dots = u_{j+s_R}^n = u \implies m_{j+1/2}^n = M_{j+1/2}^n = G(u).$$

This motivates the choice

$$\mathcal{J}^C(\gamma; (u_j^n)_j) = \left(\frac{\Delta t}{\Delta x}\right)^2 \sum_{j=0}^M \left(\max\left(0, \gamma_{j+1/2}^n - M_{j+1/2}^n\right)^2 + \max\left(0, m_{j+1/2}^n - \gamma_{j+1/2}^n\right)^2\right),$$

which vanishes if and only if  $\forall j \in \{1/2, \dots, M+1/2\}$ ,  $m_{j+1/2}^n \leq \gamma_{j+1/2}^n \leq M_{j+1/2}^n$ . In the particular case where  $u_{j-s_L+1}^n = \dots = u_{j+s_R}^n = u$ , the  $j$ -th term of the sum in  $\mathcal{J}^C(\gamma; (u_j^n)_j)$  vanishes if and only if  $\gamma_{j+1/2}^n = G(u)$ , meaning that  $\gamma$  is consistent. The function  $\mathcal{J}$  is nonnegative,  $\mathcal{C}^1$ -regular and convex with respect to  $\gamma$ , but the set where it vanishes has no reason either to exist or to be reduced to a single point.

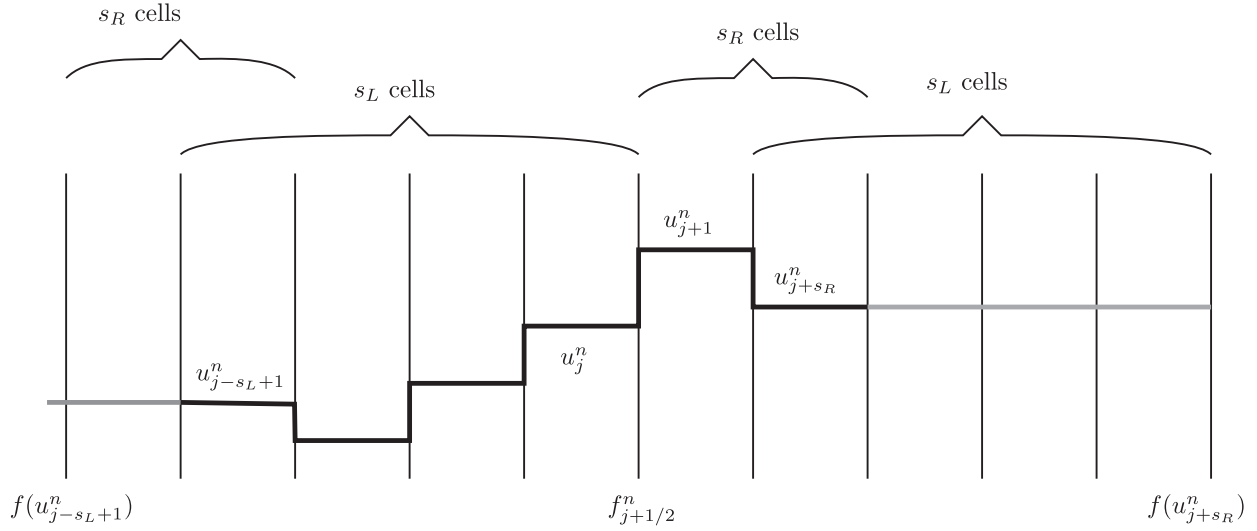


FIGURE 1. Modification of the initial data away from the stencil. Consistency can be used to compute the fluxes at interfaces  $j + 3/2 - s_L - s_R$  and  $j - 1/2 + s_L + s_R$  and the value of  $f_{j+1/2}^n$  remains unchanged.

We now define the bounds  $m_{j+1/2}^n$  and  $M_{j+1/2}^n$ . To do so, we zoom on the interface  $j + 1/2$  and construct a modified initial data  $(\tilde{u}_k^{n,j+1/2})_k$  as illustrated on Figure 1 in such a way that the consistency property can be used as close as possible to the interface  $j + 1/2$  without changing the value of the numerical flux  $f_{j+1/2}^n$ . The easiest way to proceed consists in extending the initial data by constant values outside of the stencil.

Then the finite volume scheme (4) applied to this new initial data gives numerical fluxes  $(\tilde{f}_{k+1/2}^{n,j+1/2})_k$ . They are used to compute the update  $(\hat{u}_k^{n,j+1/2})_k$ . The benefit brought by this modified data is that the evolution of the total entropy on the left and the right of interface  $j + 1/2$  can be expressed only in terms of  $G_{j+1/2}^n$ ,  $(\tilde{u}_k^{n,j+1/2})_k$  and  $(\hat{u}_k^{n,j+1/2})_k$ . It yields bounds on the individual numerical entropy flux  $G_{j+1/2}^n$  that can be constructed from the initial data  $(u_k^n)_k$  and depends on the chosen numerical scheme.

**Lemma 1.** *Consider a consistent entropy satisfying scheme in the sense of Definition 2. Then, for all  $j \in \llbracket 1, M \rrbracket$ , we have  $m_{j+1/2}^n \leq G_{j+1/2}^n \leq M_{j+1/2}^n$ , with  $s$*

$$\begin{cases} M_{j+1/2}^n = G(u_{j-s_L+1}^n) + \frac{\Delta x}{\Delta t} \sum_{k=j-s_L-s_R+2}^j \eta(\tilde{u}_k^{n,j+1/2}) - \eta(\hat{u}_k^{n,j+1/2}), \\ m_{j+1/2}^n = G(u_{j+s_R}^n) + \frac{\Delta x}{\Delta t} \sum_{k=j+1}^{j+s_L+s_R-1} \eta(\tilde{u}_k^{n,j+1/2}) - \eta(\hat{u}_k^{n,j+1/2}), \end{cases} \quad (9)$$

where

$$\tilde{u}_k^{n,j+1/2} = u_{\min(\max(k,j-s_L+1),j+s_R)}^n \quad (10)$$

and for all  $k \in \llbracket 1, M \rrbracket$ ,

$$\hat{u}_k^{n,j+1/2} = \tilde{u}_k^{n,j+1/2} - \frac{\Delta t}{\Delta x} \left( \mathcal{F} \left( \tilde{u}_{k-s_L+1}^{n,j+1/2}, \dots, \tilde{u}_k^{n,j+1/2}, \tilde{u}_{k+1}^{n,j+1/2}, \dots, \tilde{u}_{k+s_R}^{n,j+1/2} \right) - \mathcal{F} \left( \tilde{u}_{k-s_L}^{n,j+1/2}, \dots, \tilde{u}_{k-1}^{n,j+1/2}, \tilde{u}_k^{n,j+1/2}, \dots, \tilde{u}_{k+s_R-1}^{n,j+1/2} \right) \right).$$

*Proof.* The flux at the interface  $j + 1/2$  is independent of the values in cells  $j - s_L$  and smaller, and in cells  $j + s_R + 1$  and larger. Thus, we modify the initial data by extending it by  $u_{j-s_L+1}^n$  on its left and by  $u_{j+s_R}^n$  on its right

$$\tilde{u}_k^{n,j+1/2} = \begin{cases} u_{j-s_L+1}^n & \text{if } k \leq j - s_L, \\ u_k^n & \text{if } j - s_L + 1 \leq k \leq j + s_R, \\ u_{j+s_R}^n & \text{if } j + s_R + 1 \leq k, \end{cases}$$

which corresponds to (10). We update this modified initial data with the chosen finite volume scheme  $\mathcal{F}$  and obtain  $\hat{u}_k^{n,j+1/2}$ .

The interest of the modified initial data  $\left( \tilde{u}_k^{n,j+1/2} \right)_k$  is that the fluxes at interfaces  $j - s_L - s_R + 3/2$  (and before) and  $j + s_R + s_L - 1/2$  (and after) are given by consistency:

$$\begin{aligned} \mathcal{F} \left( \tilde{u}_{j-2s_L-s_R+2}^{n,j+1/2}, \dots, \tilde{u}_{j-s_L-s_R+1}^{n,j+1/2}, \tilde{u}_{j-s_L-s_R+2}^{n,j+1/2}, \dots, \tilde{u}_{j-s_L+1}^{n,j+1/2} \right) &= f \left( u_{j-s_L+1}^n \right), \\ \mathcal{F} \left( \tilde{u}_{j+s_R}^{n,j+1/2}, \dots, \tilde{u}_{j+s_L+s_R-1}^{n,j+1/2}, \tilde{u}_{j+s_L+s_R}^{n,j+1/2}, \dots, \tilde{u}_{j+s_L+2s_R-1}^{n,j+1/2} \right) &= f \left( u_{j+s_R}^n \right). \end{aligned}$$

Now, as the scheme is entropy satisfying

$$\forall k, \eta \left( \hat{u}_k^{n,j+1/2} \right) \leq \eta \left( \tilde{u}_k^{n,j+1/2} \right) - \frac{\Delta t}{\Delta x} \left( \tilde{G}_{k+1/2}^{n,j+1/2} - \tilde{G}_{k-1/2}^{n,j+1/2} \right),$$

where  $\tilde{G}_{k+1/2}^{n,j+1/2} = \mathcal{G} \left( \tilde{u}_{k-s_L+1}^{n,j+1/2}, \dots, \tilde{u}_k^{n,j+1/2}, \tilde{u}_{k+1}^{n,j+1/2}, \dots, \tilde{u}_{k+s_R}^{n,j+1/2} \right)$  for some function  $\mathcal{G}$ , see Definition 2. These inequalities are not of much use because we do not know an explicit formula for  $\mathcal{G}$ . However, we deduced from the construction of  $\left( \tilde{u}_k^{n,j+1/2} \right)_k$  several consequences. The consistency yields at interfaces  $j - s_L - s_R + 3/2$  and  $j + s_L + s_R - 1/2$

$$\tilde{G}_{j-s_L-s_R+3/2}^{n,j+1/2} = G \left( u_{j-s_L+1}^n \right) \quad \text{and} \quad \tilde{G}_{j+s_L+s_R-1/2}^{n,j+1/2} = G \left( u_{j+s_R}^n \right).$$

On the other hand the flux at interface  $j + 1/2$  remains unchanged

$$\tilde{G}_{j+1/2}^{n,j+1/2} = \mathcal{G} \left( \tilde{u}_{j-s_L+1}^{n,j+1/2}, \dots, \tilde{u}_{j+s_R}^{n,j+1/2} \right) = \mathcal{G} \left( u_{j-s_L+1}^n, \dots, u_{j+s_R}^n \right) = G_{j+1/2}^n.$$

We eliminate the other numerical entropy fluxes by summation:

$$\begin{aligned} \sum_{k=j+1}^{j+s_L+s_R-1} \eta \left( \hat{u}_k^{n,j+1/2} \right) &\leq \left( \sum_{k=j+1}^{j+s_L+s_R-1} \eta \left( \tilde{u}_k^{n,j+1/2} \right) \right) - \frac{\Delta t}{\Delta x} \left( G \left( u_{j+s_R}^n \right) - G_{j+1/2}^n \right), \\ \sum_{k=j-s_L-s_R+2}^j \eta \left( \hat{u}_k^{n,j+1/2} \right) &\leq \sum_{k=j-s_L-s_R+2}^j \eta \left( \tilde{u}_k^{n,j+1/2} \right) - \frac{\Delta t}{\Delta x} \left( G_{j+1/2}^n - G \left( u_{j-s_L+1}^n \right) \right). \end{aligned}$$

We can now bound  $G_{j+1/2}^n$  from above and below and conclude. □

**Remark 1.** In the case of two points scheme  $s_L = s_R = 1$ , Lemma 1 correspond to the notion of “interface entropy inequality” in Definition 2.7 of [3]. In this specific case it implies the desired discrete entropy inequality (5) at the cost of a time step  $\Delta t$  twice smaller Proposition 2.9 of [3].



### 3. MAIN RESULTS

In this section, we state some properties of the minimizers of  $\mathcal{J}$  and revisit some standard concepts of finite volume scheme theory from our a posteriori point of view.

#### 3.1. Entropy dissipation and zero minimization

**Proposition 1.** *Consider a finite volume scheme (4) that admits a discrete entropy inequality (5) for some numerical entropy fluxes  $\left(G_{j+1/2}^n = \mathcal{G}(u_{j-s_L+1}^n, \dots, u_{j+s_R}^n)\right)_j$ . Then for all initial data  $(u_j^n)_j$  we have*

$$\mathcal{J}\left(\left(G_{j+1/2}^n\right)_j; (u_j^n)_j\right) = 0.$$

*Proof.* This property follows from the definition of the functional. If (5) holds, then

$$\mathcal{J}^D\left(\left(G_{j+1/2}^n\right)_j; (u_j^n)_j\right) = 0.$$

It remains to prove that the second part also vanishes, which is the case if and only if  $m_{j+1/2}^n \leq G_{j+1/2}^n \leq M_{j+1/2}^n$ . This holds by construction on the bounds  $m_{j+1/2}^n$  and  $M_{j+1/2}^n$ . The scheme is also entropy diminishing on the modified initial data  $\left(\tilde{u}_k^{n,j+1/2}\right)_k$ . The lower (resp. upper) bound follows from the diminution of total entropy in the  $s_L + s_R - 1$  cells on the left (resp. right cells) during the time step. Details are given in the proof of Lemma 1 in the previous section.  $\square$

**Remark 2.** The contraposition of Proposition 1 states that if there exists an initial data  $(u_j^n)_j$  such that for all  $\left(G_{j+1/2}^n\right)_j$ ,  $\mathcal{J}\left(\left(G_{j+1/2}^n\right)_j; (u_j^n)_j\right) > 0$  then the scheme is not entropy satisfying. We further exploit this fact in Section 5 to build another minimization procedure which constructs the worst initial data in terms of entropy and use it to prove that some schemes are non entropy satisfying.

**Proposition 2.** *Fix the initial data  $(u_j^n)_j$ . Suppose that there exists  $\left(G_{j+1/2}^{opt,n}\right)_j$  such that the functional vanishes*

$$\mathcal{J}\left(\left(G_{j+1/2}^{opt,n}\right)_j; (u_j^n)_j\right) = 0. \text{ Then the corresponding scheme verifies}$$

$$\forall j \in \llbracket 1, M \rrbracket, \quad \eta(u_j^{n+1}) \leq \eta(u_j^n) - \frac{\Delta t}{\Delta x} \left(G_{j+1/2}^{opt,n} - G_{j-1/2}^{opt,n}\right).$$

*The scheme is also consistent in the sense that if  $u_{j-s_L+1}^n = \dots = u_j^n = \dots = u_{j+s_R}^n$ , then  $G_{j+1/2}^{opt,n} = G(u_j^n)$ .*

Let us insist that this result only shows that there exist numerical entropy fluxes yielding to nonpositive numerical diffusion for the particular choice of initial data  $(u_j^n)_j$ , which does not mean that the scheme is always entropy satisfying.

*Proof.* Suppose that  $\mathcal{J}\left(\left(G_{j+1/2}^{opt,n}\right)_j; (u_j^n)_j\right) = 0$ . The first contribution  $\mathcal{J}^D$  is zero thus

$$\forall j \in \llbracket 1, M \rrbracket, \quad \eta(u_j^{n+1}) - \eta(u_j^n) + \frac{\Delta t}{\Delta x} \left(G_{j+1/2}^{opt,n} - G_{j-1/2}^{opt,n}\right) \leq 0,$$

which is exactly (5). The term  $\mathcal{J}^C$  is also zero, thus  $m_{j+1/2}^n \leq G_{j+1/2}^{opt,n} \leq M_{j+1/2}^n$ . If  $(u_{j-s_L+1}^n, \dots, u_{j+s_R}^n) = (u_j^n, \dots, u_j^n)$ , then  $\left(\tilde{u}_k^{n,j+1/2}\right)_k$  is constant equal to  $u_j^n$ , as well as  $\left(\hat{u}_k^{n,j+1/2}\right)_k$ , so that  $m_{j+1/2}^n = M_{j+1/2}^n = G(u_j^n)$  and  $G_{j+1/2}^{opt,n} = G(u_j^n)$ .  $\square$

### 3.2. Discrepancy between global minimizers

When the set where the functional  $\mathcal{J}$  vanishes is nonempty, it most likely contains several solutions. We first quantify how close they are from each other.

**Proposition 3.** *Suppose that  $\mathcal{J}\left(\left(G_{j+1/2}^n\right)_j; \left(u_j^n\right)_j\right) = 0$  and  $\mathcal{J}\left(\left(\bar{G}_{j+1/2}^n\right)_j; \left(u_j^n\right)_j\right) = 0$ , and that the numerical flux  $\mathcal{F}$  used in (4) is consistent and Lipschitz regular. Then for all  $j$ ,*

$$\left|G_{j+1/2}^n - \bar{G}_{j+1/2}^n\right| = \sum_{k \in \{j-s_L+1, \dots, j+s_R\}} O\left(\left|u_k^n - u_j^n\right|^2\right).$$

When the scheme is known to satisfy a given discrete entropy inequality and when the solution is smooth, this result shows that the difference between the numerical entropy flux found in the literature and the one returned by the minimization procedure is of order  $\Delta x^2$ . Note that the consistency property alone gives an order  $\Delta x$ .

*Proof.* Suppose that the minimization procedure has two different global minimizers  $\mathcal{J}\left(\left(G_{j+1/2}^n\right)_j; \left(u_j^n\right)_j\right) = 0$  and  $\mathcal{J}\left(\left(\bar{G}_{j+1/2}^n\right)_j; \left(u_j^n\right)_j\right) = 0$ . Then  $m_{j+1/2}^n \leq G_{j+1/2}^n \leq M_{j+1/2}^n$  and  $m_{j+1/2}^n \leq \bar{G}_{j+1/2}^n \leq M_{j+1/2}^n$ , thus  $\left|\bar{G}_{j+1/2}^n - G_{j+1/2}^n\right| \leq M_{j+1/2}^n - m_{j+1/2}^n$ , where

$$\begin{aligned} M_{j+1/2}^n - m_{j+1/2}^n &= G\left(u_{j-s_L+1}^n\right) - G\left(u_{j+s_R}^n\right) \\ &\quad + \frac{\Delta x}{\Delta t} \sum_{k=j-s_L-s_R+2}^{j+s_L+s_R-1} \eta\left(\tilde{u}_k^{n,j+1/2}\right) - \eta\left(\hat{u}_k^{n,j+1/2}\right). \end{aligned} \tag{11}$$

By convexity of the entropy  $\eta$ ,

$$\begin{aligned} \eta\left(\hat{u}_k^{n,j+1/2}\right) &= \eta\left(\tilde{u}_k^{n,j+1/2} - \frac{\Delta t}{\Delta x}\left(\tilde{f}_{k+1/2}^{n,j+1/2} - \tilde{f}_{k-1/2}^{n,j+1/2}\right)\right), \\ &\geq \eta\left(\tilde{u}_k^{n,j+1/2}\right) - \frac{\Delta t}{\Delta x} D\eta\left(\tilde{u}_k^{n,j+1/2}\right)\left(\tilde{f}_{k+1/2}^{n,j+1/2} - \tilde{f}_{k-1/2}^{n,j+1/2}\right). \end{aligned}$$

Combining this inequality with (11), we get

$$M_{j+1/2}^n - m_{j+1/2}^n \leq G\left(u_{j-s_L+1}^n\right) - G\left(u_{j+s_R}^n\right) + \sum_{k=j-s_L-s_R+2}^{j+s_L+s_R-1} D\eta\left(\tilde{u}_k^{n,j+1/2}\right)\left(\tilde{f}_{k+1/2}^{n,j+1/2} - \tilde{f}_{k-1/2}^{n,j+1/2}\right),$$

which rewrites as

$$\begin{aligned} M_{j+1/2}^n - m_{j+1/2}^n &\leq G\left(u_{j-s_L+1}^n\right) - G\left(u_{j+s_R}^n\right) + D\eta\left(u_j^n\right)\left(f\left(u_{j+s_R}^n\right) - f\left(u_{j-s_L+1}^n\right)\right) \\ &\quad + \sum_{k=j-s_L-s_R+2}^{j+s_L+s_R-1} \left(D\eta\left(\tilde{u}_k^{n,j+1/2}\right) - D\eta\left(u_j^n\right)\right)\left(\tilde{f}_{k+1/2}^{n,j+1/2} - \tilde{f}_{k-1/2}^{n,j+1/2}\right). \end{aligned}$$

If the numerical flux is Lipschitz regular, the last sum is  $\sum_{k=j-s_L+1}^{j+s_R} O\left(\left|u_k^n - u_j^n\right|^2\right)$ . A first order expansion at point  $u_j^n$  of the other contributions is

$$\left(DG\left(u_j^n\right) - D\eta\left(u_j^n\right) Df\left(u_j^n\right)\right)\left(u_{j-s_L+1}^n - u_{j+s_R}^n\right) + O\left(\left|u_{j-s_L+1}^n - u_j^n\right|^2\right) + O\left(\left|u_{j+s_R}^n - u_j^n\right|^2\right).$$

The result then follows from  $D\eta Df = DG$ . □

### 3.3. Lax-Wendroff theorem

One of the main theoretical results about numerical schemes for systems of conservation laws is the Lax-Wendroff theorem, which states that if a numerical scheme converges in a certain sense, then the limit is a weak solution. In addition, if the scheme satisfies relevant discrete entropy inequalities, then the limit is an entropy solution.

The latter statement usually requires the numerical entropy flux to be a continuous and consistent function  $\mathcal{G}$  of the neighboring approximations. This is not the case in this work, since the numerical entropy fluxes  $G_{j+1/2}^n$  are obtained through a minimization procedure and we do not use an entropy flux function  $\mathcal{G}$ . However, it is possible to adapt the Lax-Wendroff theorem to our framework. This extension somehow means it is possible to replace the classical Definition 2 of entropy satisfying scheme by the requirement “for all initial data  $(u_j^n)_j$ , there exists  $G^{opt,n}$  such that  $\mathcal{J}\left(G^{opt,n}, (u_j^n)_j\right) = 0$ ”. The latter condition differs on the treatment of the consistency property but leads to an identical Lax-Wendroff theorem. Note that proving such a result has no reason to be easier than finding an explicit function  $\mathcal{G}$  and is possibly more difficult as  $\mathcal{G}$  has  $s_L + s_R + 1$  variables while  $\mathcal{J}$  is defined on  $\mathbf{R}^M$ .

Since the time discretization will be important in this section, the objects related to the optimization problem (7) at time  $t^n$ , *i.e.*, with initial data  $(u_j^n)_j$ , will be noted with a superscript  $n$ . The following Theorem is given on a bounded space domain  $[a, b]$  with periodic boundary conditions since it is the framework of this paper. It can easily be generalized on the unbounded space domain  $\mathbf{R}$ .

For a discretization  $\Delta = (\Delta x, \Delta t)$ , we introduce the piecewise constant function  $u_\Delta$  defined almost everywhere (a.e.) in  $[a, b] \times [0, +\infty)$  by

$$u_\Delta(x, t) = u_j^n, \quad x_{j-1/2} < x < x_{j+1/2}, \quad t^n \leq t < t^{n+1}.$$

**Theorem 1.** *Assume the scheme (4) is consistent in the sense of Definition 1, with  $\mathcal{F}$  continuous. Consider a sequence  $\Delta_k = (\Delta x_k, \Delta t_k)$  which converges to  $(0, 0)$  with the ratio  $\lambda = \frac{\Delta t_k}{\Delta x_k}$  being constant. Assume that*

- *there exists a compact  $K \subset \Omega$  such that  $u_{\Delta_k}(x, t) \in K$  for a.e.  $(x, t) \in [a, b] \times [0, +\infty)$  and for all  $k \in \mathbf{N}$ ;*
- *the sequence  $u_{\Delta_k}$  converges in  $L^1_{loc}([a, b] \times (0, +\infty))$  to a function  $u$ .*

*Then  $u$  is a weak solution of (1).*

*Furthermore, if for a given entropy pair  $(\eta, G)$ , for all discretization  $\Delta_k = (\Delta x_k, \Delta t_k)$  and for all  $n \in \mathbf{N}$ , there exists a family  $\left(G_{j+1/2}^{opt,n}\right)_j$  such that  $\mathcal{J}\left(\left(G_{j+1/2}^{opt,n}\right)_j; (u_j^n)_j\right) = 0$ , then the solution  $u$  satisfies the entropy inequality (2) in the sense of distributions on  $[a, b] \times (0, +\infty)$ .*

*Proof.* For the sake of simplicity, the subscript  $k$  in  $\Delta_k$  will be omitted all along the proof. The proof of the convergence of  $u_\Delta$  toward a weak solution of (1) is exactly the same as in the original Lax-Wendroff theorem. The reader is referred for instance to [14, 19] for a complete proof.

Concerning the entropy inequality, for a fixed discretization  $\Delta$ , let us consider a compactly supported test function  $\varphi \in C_0^1([a, b] \times (0, +\infty))$ , with  $\varphi \geq 0$ . Let us notice that the function  $\varphi$  is assumed to vanish when  $t = 0$ . This choice is made for the sake of simplicity of the proof. However, the extension to the case where  $\varphi$  does not vanish at  $t = 0$  is immediate.

Since  $\mathcal{J}^D\left(\left(G_{j+1/2}^{opt,n}\right)_j; (u_j^n)_j\right) = 0$ , we have for all  $j$  and  $n$

$$\eta(u_j^{n+1}) - \eta(u_j^n) + \lambda \left(G_{j+1/2}^{opt,n} - G_{j-1/2}^{opt,n}\right) \leq 0.$$

Multiplying this inequality by  $\Delta x \varphi(x_j, t^n)$ , summing over  $j$  and  $n$  and performing a summation by parts, we obtain

$$\Delta x \sum_{j,n} \eta(u_j^{n+1}) (\varphi(x_j, t^{n+1}) - \varphi(x_j, t^n)) + \Delta t \sum_{j,n} G_{j+1/2}^{opt,n} (\varphi(x_{j+1}, t^n) - \varphi(x_j, t^n)) \geq 0.$$

Introducing the piecewise constant functions

$$\begin{aligned} \varphi_\Delta(x, t) &= \varphi(x_j, t^n), \quad x_{j-1/2} < x < x_{j+1/2}, \quad t^n \leq t < t^{n+1}, \\ G_\Delta^{opt}(x, t) &= G_{j+1/2}^{opt, n}, \quad x_j < x < x_{j+1}, \quad t^n \leq t < t^{n+1}, \end{aligned}$$

where  $x_j = \frac{x_{j-1/2} + x_{j+1/2}}{2}$ , the last inequality writes

$$\begin{aligned} \int_{[a, b] \times [\Delta t, +\infty)} \eta(u_\Delta(x, t)) \frac{\varphi_\Delta(x, t) - \varphi_\Delta(x, t - \Delta t)}{\Delta t} dx dt \\ + \int_{[a, b] \times \mathbf{R}^+} G_\Delta^{opt}(x, t) \frac{\varphi_\Delta(x + \Delta x/2) - \varphi_\Delta(x - \Delta x/2)}{\Delta x} dx dt \geq 0. \end{aligned} \tag{12}$$

As in the classical Lax-Wendroff theorem, the first integral converges as  $\Delta \rightarrow (0, 0)$  to

$$\int_{[a, b] \times \mathbf{R}^+} \eta(u(x, t)) \partial_t \varphi(x, t) dx dt.$$

The difference with the classical Lax-Wendroff theorem lies in the second integral of (12). First, let us introduce the following piecewise constant functions:

$$\begin{aligned} M_\Delta(x, t) &= M_{j+1/2}^n, \quad x_j < x < x_{j+1}, \quad t^n \leq t < t^{n+1}, \\ m_\Delta(x, t) &= m_{j+1/2}^n, \quad x_j < x < x_j, \quad t^n \leq t < t^{n+1}, \\ \tilde{u}_\Delta^{j+1/2}(x, t) &= \tilde{u}_k^{j+1/2}, \quad x_{k-1/2} < x < x_{k+1/2}, \quad t^n \leq t < t^{n+1}, \\ \tilde{f}_\Delta^{j+1/2}(x, t) &= \tilde{f}_{k+1/2}^{j+1/2}, \quad x_k < x < x_{k+1}, \quad t^n \leq t < t^{n+1}, \\ \hat{u}_\Delta^{j+1/2}(x, t) &= \hat{u}_k^{j+1/2}, \quad x_{k-1/2} < x < x_{k+1/2}, \quad t^n \leq t < t^{n+1}. \end{aligned} \tag{13}$$

$$\tag{14}$$

Since  $\varphi$  is smooth, the term  $\frac{\varphi_\Delta(x + \Delta x/2) - \varphi_\Delta(x - \Delta x/2)}{\Delta x}$  uniformly converges to  $\partial_x \varphi(x, t)$ . Moreover, according to (9) and  $\tilde{u}_k^{n, j+1/2}$  and since the ratio  $\lambda$  is constant,  $M_{j+1/2}^n$  and  $m_{j+1/2}^n$  are continuous functions of a finite number of  $u_k^n$ , that all lie in the compact  $K$ . Therefore the functions  $M_\Delta$  and  $m_\Delta$  are uniformly bounded. Since  $\mathcal{J}^C \left( \left( G_{j+1/2}^{opt, n} \right)_j ; \left( u_j^n \right)_j \right) = 0$ , it follows that the inequalities

$$m_\Delta^n(x, t) \leq G_\Delta^{opt}(x, t) \leq M_\Delta^n(x, t) \tag{15}$$

hold for a.e.  $(x, t) \in [a, b] \times [0, +\infty)$ . As a consequence, the function  $G_\Delta^{opt}$  is also uniformly bounded and therefore

$$\int_{[a, b] \times \mathbf{R}^+} G_\Delta^{opt}(x, t) \left( \frac{\varphi_\Delta(x + \frac{\Delta x}{2}) - \varphi_\Delta(x - \frac{\Delta x}{2})}{\Delta x} - \partial_x \varphi(x, t) \right) dx dt \rightarrow 0. \tag{16}$$

Next we have

$$\begin{aligned} \tilde{u}_\Delta^{j+1/2}(x, t) &= u_\Delta(x + \min(\max(k, j - s_L + 1), j + s_R) \Delta x - k \Delta x, t), \\ \tilde{f}_\Delta^{j+1/2}(x, t) &= \mathcal{F} \left( \tilde{u}_\Delta^{j+1/2}(x - (s_L - 1/2) \Delta x, t), \dots, \tilde{u}_\Delta^{j+1/2}(x + (s_R - 1/2) \Delta x, t) \right), \\ \hat{u}_\Delta^{j+1/2}(x, t) &= \tilde{u}_\Delta^{j+1/2}(x, t) - \lambda \left( \tilde{f}_\Delta^{j+1/2}(x + \Delta x/2, t) - \tilde{f}_\Delta^{j+1/2}(x - \Delta x/2, t) \right). \end{aligned} \tag{17}$$

It follows from (9) that

$$\begin{aligned}
 M_\Delta(x, t) &= G(u_\Delta(x - (s_L - 1/2)\Delta x, t)) \\
 &+ \lambda \sum_{k=j-s_L-s_R+2}^j \left( \eta\left(\tilde{u}_\Delta^{j+1/2}(x - (j + 1/2 - k)\Delta x, t)\right) \right. \\
 &\left. - \eta\left(\hat{u}_\Delta^{j+1/2}(x - (j + 1/2 - k)\Delta x, t)\right) \right). \tag{18}
 \end{aligned}$$

Since  $u_\Delta$  converges in  $L^1_{loc}([a, b] \times (0, +\infty))$  to  $u$ , then up to a subsequence,  $\tilde{u}_\Delta^{j+1/2}$  converges to  $u$  a.e. Thanks to the continuity and the consistency of  $\mathcal{F}$ , we get from (13) that  $\tilde{f}_\Delta^{j+1/2}$  converges a.e. to  $f(u)$ . Then according to (14) and (17), we deduce that  $\hat{u}_\Delta^{j+1/2}$  also converges a.e. to  $u$ . Finally, it follows from (18) that  $M_\Delta$  converges a.e. to  $G(u)$ . A similar process shows that up to a subsequence,  $m_\Delta$  converges a.e. to  $G(u)$ . We deduce from (15) that up to a subsequence,  $G_\Delta^{opt}$  converges a.e. to  $G(u)$ . The dominated convergence theorem then ensures that

$$\int_{[a,b] \times \mathbf{R}^+} G_\Delta^{opt}(x, t) \partial_x \varphi(x, t) \, dx dt \rightarrow \int_{[a,b] \times \mathbf{R}^+} G(u) \partial_x \varphi(x, t) \, dx dt. \tag{19}$$

Summing (16) and (19), we obtain

$$\int_{[a,b] \times \mathbf{R}^+} G_\Delta^{opt}(x, t) \frac{\varphi_\Delta(x + \Delta x/2) - \varphi_\Delta(x - \Delta x/2)}{\Delta x} \, dx dt \rightarrow \int_{[a,b] \times \mathbf{R}^+} G(u) \partial_x \varphi(x, t) \, dx dt.$$

Hence the limit of (12) writes

$$\int_{[a,b] \times \mathbf{R}^+} \eta(u(x, t)) \partial_t \varphi(x, t) \, dx dt + \int_{[a,b] \times \mathbf{R}^+} G(u) \partial_x \varphi(x, t) \, dx dt \geq 0,$$

which concludes the proof. □

#### 4. A POSTERIORI QUANTIFICATION OF THE NUMERICAL DIFFUSION: NUMERICAL RESULTS

We now apply the minimization of the functional  $\mathcal{J}$  to obtain maps of numerical diffusion in different settings<sup>1</sup>. The numerical entropy fluxes  $G_{j\pm 1/2}^{opt,n}$  are the results of the minimization (7) and we define the corresponding a posteriori numerical diffusion according to (6) by

$$D_j^{opt,n} = \eta(u_j^{n+1}) - \eta(u_j^n) + \frac{\Delta t}{\Delta x} \left( G_{j+1/2}^{opt,n} - G_{j-1/2}^{opt,n} \right). \tag{20}$$

In practice the minimization is performed using the blackbox `fminunc` of MATLAB [23], with the initial guess  $\gamma_{j+1/2} = \frac{1}{2} \left( m_{j+1/2}^n + M_{j+1/2}^n \right)$ . The values for  $D_j^{opt,n}$  are often small, in particular on fine meshes or when the solution is smooth. However, these values remain far away from the machine precision estimated at  $10^{-16}$ .

##### 4.1. Continuous solution of the Burgers equation

In this first test case, we consider the scalar ( $d = 1$ ) Burgers equation

$$\partial_t u + \partial_x \left( \frac{u^2}{2} \right) = 0 \quad \text{and} \quad \partial_t (u^2) + \partial_x \left( \frac{2u^3}{3} \right) \leq 0. \tag{21}$$

---

<sup>1</sup>In all the following numerical tests, we use a personal laptop, equipped with a 2.6 GHz Intel Core i7 6-core processor of 16 GB and 2667 MHz DDR4 memory.

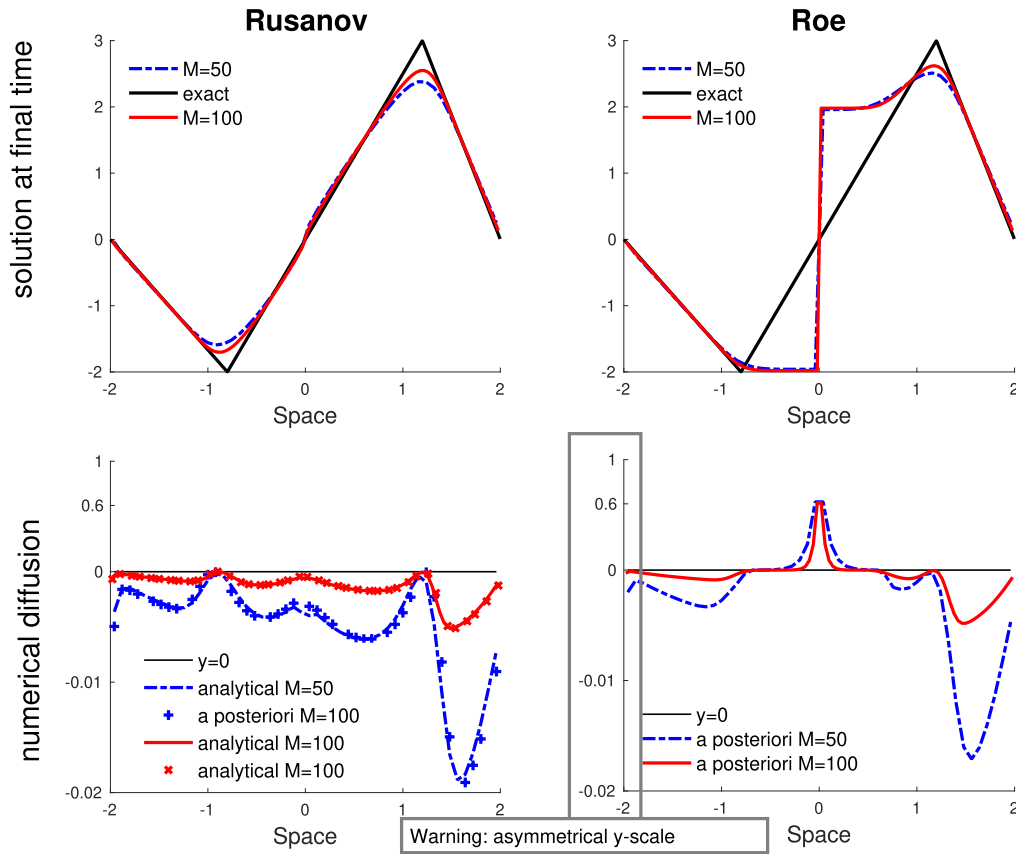


FIGURE 2. Quantification of the numerical diffusion on Testcase (22) for the entropy satisfying Rusanov scheme and the entropy violating Roe scheme.

The choice of  $\eta(u) = u^2$  is arbitrary; for the scalar Burgers equation any convex entropy  $\eta$  could be considered. The initial data is

$$u^0(x) = \begin{cases} -2 - x & \text{if } -2 < x \leq 0, \\ 3 - \frac{3}{2}x & \text{if } 0 < x \leq 2, \end{cases} \quad (22)$$

and we use periodic boundary conditions. The exact solution is piecewise linear, and for  $t < \frac{2}{3}$  is given by

$$u(x, t) = \begin{cases} \frac{-(x+2)}{1-t} & \text{if } -2 \leq x \leq -2t, \\ \frac{x}{t} & \text{if } -2t \leq x \leq 3t, \\ \frac{-3(x-2)}{2-3t} & \text{if } 3t \leq x \leq 2. \end{cases}$$

The space interval  $[-2, 2]$  is discretized with 50 or 100 cells, and the functional  $\mathcal{J}$  is minimized at the last iteration at  $T = 0.4$ . The parameter  $\alpha$  in the Courant–Friedrichs–Levy condition (3) is set to  $\alpha = 0.5$ .

We compared the entropy satisfying Rusanov scheme and the entropy violating Roe scheme with a forward Euler march in time, see the appendix for the definitions of the corresponding fluxes  $\mathcal{F}$ . The stencil is  $s_L = s_R = 1$ . The minimization procedure detects the two previous behaviors:  $D_j^{opt,n}$  remains nonpositive everywhere for

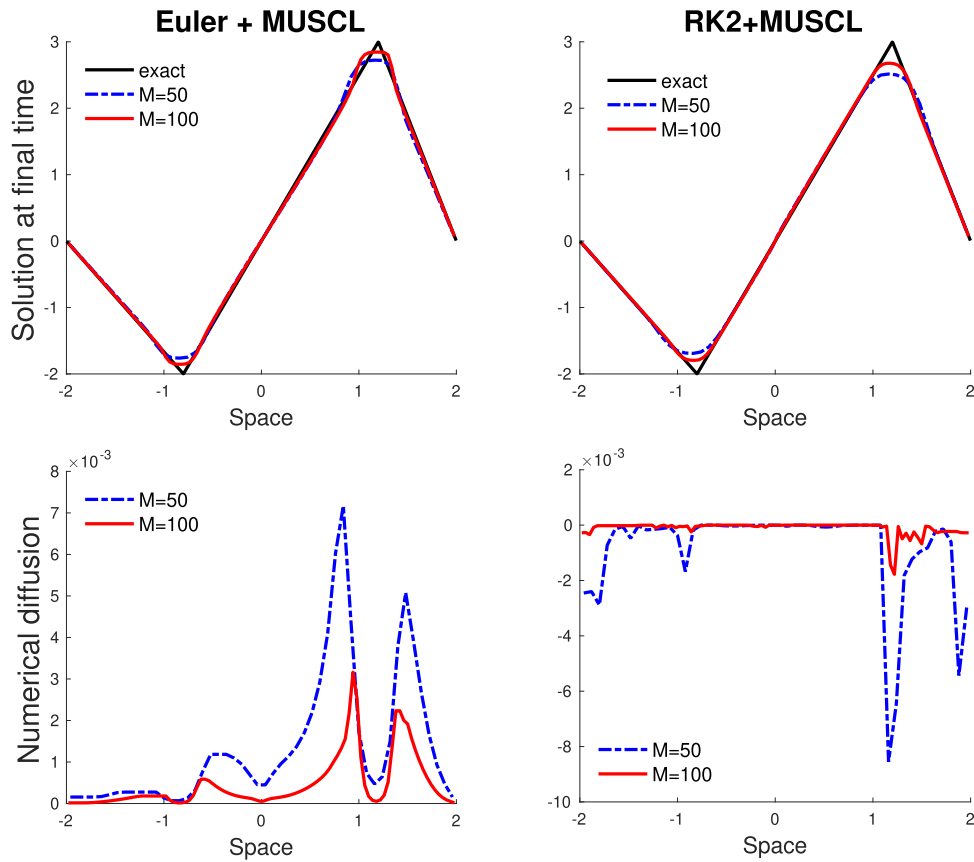


FIGURE 3. Influence of the time discretization on the numerical diffusion for the MUSCL scheme Testcase (22).

the Rusanov scheme, while it has large strictly positive values localized around the stationary entropy creating shock at the sonic point for the Roe schemes.

The two plots of Figure 3 show the results for the MUSCL second order flux in space based on a Rusanov flux and a minmod limiter. We illustrate here the well-known fact that increasing the order in space without increasing the order in time is detrimental in terms of entropy, whereas increasing both orders is efficient. The reader will find more specific details related to the MUSCL scheme in the following section. Together with a forward Euler march in time, (5) is grossly false in the sense that the total entropy increases with time:  $\sum_j \eta(u_j^{n+1}) > \sum_j \eta(u_j^n)$ . Combined with a second order Runge-Kutta time stepping, the a posteriori procedure finds a discrete entropy inequality (5).

To conclude this section, let us illustrate Proposition 3. The Rusanov scheme with an Euler method is entropy satisfying. With the choice of  $\mathcal{G}$  given in the appendix (A.2) the discrete entropy inequality (5) holds. We denote by  $D_j^n$  the associated nonpositive numerical diffusion (6). We compare the total discrepancy  $\Delta x \sum |D_j^{opt,n} - D_j^n|$  between this quantity and the one obtained through minimization (20) and obtained as expected that it decreases one order faster than the total diffusion  $\Delta x \sum_j |D_j^n| = \Delta x \sum_j |\eta(u_j^{n+1}) - \eta(u_j^n)|$  as the mesh becomes finer and finer. Note that this quantity does not depend on the numerical entropy fluxes  $G_{j+1/2}^n$ .

## 4.2. Gas dynamics

We now turn to the Euler equations. The unknown is  $u = (\rho, \rho v, E)$ , where  $\rho$  is the density of the fluid,  $v$  is its velocity and  $E$  its total energy. This system reads

$$\begin{cases} \partial_t \rho + \partial_x (\rho v) = 0 \\ \partial_t (\rho v) + \partial_x (\rho v^2 + p) = 0 \\ \partial_t E + \partial_x (v(E + p)) = 0. \end{cases} \quad (23)$$

The pressure force  $p$  is related to  $\rho$  and  $E$  through an ideal gas equation of state

$$p = (\gamma - 1) \left( E - \frac{\rho v^2}{2} \right) \quad \text{where } \gamma \in (1, 3].$$

Both the density and the pressure should remain nonnegative, thus

$$\Omega = \left\{ (\rho, \rho v, E) \in \mathbf{R}^3 : \rho \geq 0 \text{ and } E \geq \frac{\rho v^2}{2} \right\}.$$

There exists an infinite number of entropy inequalities for these equations, see [1]. We consider the classical inequality on the specific entropy  $s = \frac{p}{\rho^\gamma}$

$$\partial_t (-\rho \ln(s)) + \partial_t (-v \rho \ln(s)) \leq 0. \quad (24)$$

The minimization method for quantifying the numerical diffusion does not depend on the number of unknowns  $d$ , since the quantification of the numerical diffusion only concerns the scalar equation on the entropy evolution (24). We focus on a widely chosen scheme of order two in space and time, namely the Van Leer version of the MUSCL scheme [25].

In this scheme the piecewise constant in space approximation  $(u_j^n)_j$  is replaced by a reconstructed piecewise affine data. For scalar equations  $d = 1$ , the reconstruction procedure heavily relies on the fact that the exact solutions of (1) verify a maximum principle property and are total variation diminishing (TVD). A family of functions called limiters is considered to determine the slopes in each cell [14] and allows to keep these features at the discrete level. Both properties are lost for hyperbolic systems  $d \geq 2$ . Limiters are also used but several choices are possible. We investigate the effects of some of them in this section.

It remains to decide how the evolution in time of the reconstructed initial data is performed. A first possibility is to compute the exact flux of (1) by solving generalized Riemann problems. This is only feasible for some particular hyperbolic systems and very costly from a computational point of view. This path is followed in [4] for scalar equations and in [9] for systems. A more convenient approach is to combine the reconstruction in space with a second order method in time (typically a second order Runge-Kutta method (RK2)) to obtain a precise enough approximation of the fluxes. Each individual substep is based on a first order solver.

In any case the discrete entropy inequalities found in the literature for the MUSCL approach differ from (5). The quantity  $\eta(u_j^n)$  is replaced with a linear approximation in equation (1.9) of [4] for scalar equations or with a nonlinear entropy diminishing projection operator in Theorem 2.9 of [9] for systems. In [1],  $\eta(u_j^n)$  is replaced by a convex combination of three terms that not only depends on  $u_j^n$  but also on  $u_{j-1}^n$  and  $u_{j+1}^n$ , see equations (2.7) and (2.10) of [1]. None of these modified entropy inequality is sufficient to prove a Lax-Wendroff theorem. Several numerical simulations [2] even indicate that the MUSCL+RK2 scheme may converge to incorrect solutions.

We first reproduce the Sod tube testcase of Toro [24]

$$\begin{cases} \rho^0(x) = \mathbf{1}_{x < 0} + 0.125 \times \mathbf{1}_{x \geq 0}, \\ u^0(x) = 0.75 \times \mathbf{1}_{x < 0}, \\ p^0(x) = \mathbf{1}_{x < 0} + 0.1 \times \mathbf{1}_{x \geq 0}, \end{cases} \quad (25)$$



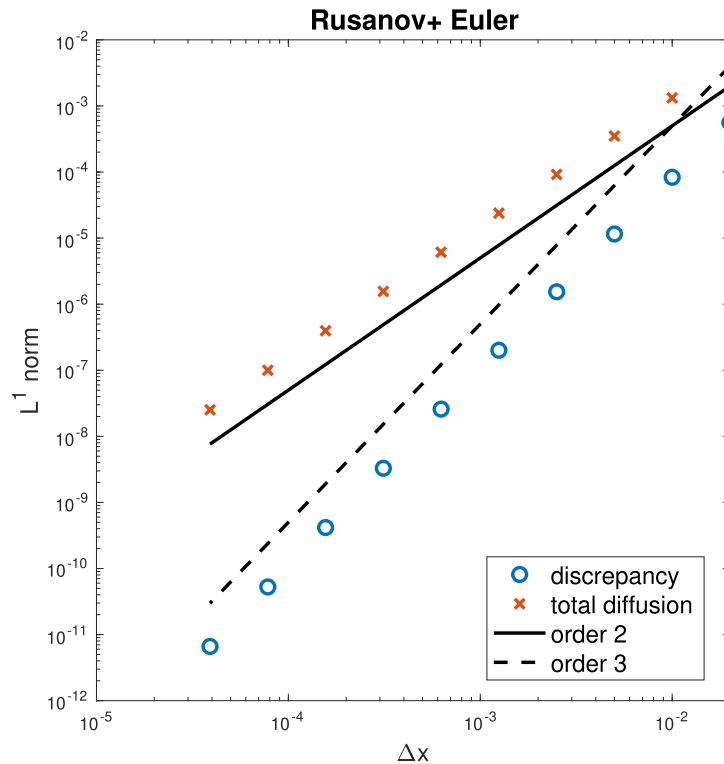


FIGURE 4. Total numerical diffusion and total discrepancy between two numerical diffusions on the Testcase (22).

on the time interval  $[0, 0.2]$  and the space interval  $[-1, 1]$ . The CFL number is  $\alpha = 1/6$  and the number of cells is  $M = 400$ . The discontinuity at  $x = 0$  creates a 1-rarefaction wave, a 2-contact discontinuity and a 3-shock. We stick to periodic boundary conditions, so there is another discontinuity at  $x = -1$ . It creates a 1-shock, a 2-contact discontinuity and a 3-rarefaction wave.

On Figure 5 we compare the first order HLLC scheme Section 10.4.2 of [24] and the Roe scheme without entropy fix Section 11.2 of [24] with a forward Euler time stepping. Then we consider the second order MUSCL scheme with a RK2 time stepping. The slopes are limited on the primitive variable  $(\rho, u, p)$  and the underlying first order scheme is the HLLC scheme. We compare the results obtained with a minmod limiter and a superbee limiter.

The a posteriori quantification of the numerical diffusion gives once again positive values near the stationary nonphysical shock created by the Roe scheme. It also detects the overcompressive behavior of the superbee limiter, with a spike of positive numerical diffusion located on the oscillations in density around the central contact discontinuity. The superbee limiter is often too strong and may prevent the scheme from converging, see [2]. For second order schemes, the numerical diffusion is located around one or two spikes around each discontinuity. This depends on the initialization of the optimization algorithm, see Figure 8 below.

Then we consider a testcase where the solution does not contain a shock, but only a contact discontinuity. The velocity and the pressure are initially constant, equal to 0.1 and 1. The initial density is

$$\rho^0(x) = 1 + 0.2x + 0.05 \sin(6\pi x) + 0.4 \times \mathbf{1}_{x < 0}. \quad (26)$$

The final time is 2 seconds, the CFL number is  $\alpha = 0.75$ . On Figure (6), we compare the HLL scheme and the HLLC scheme and find unsurprisingly that the latter is much less diffusive. As usual, we observe this fact

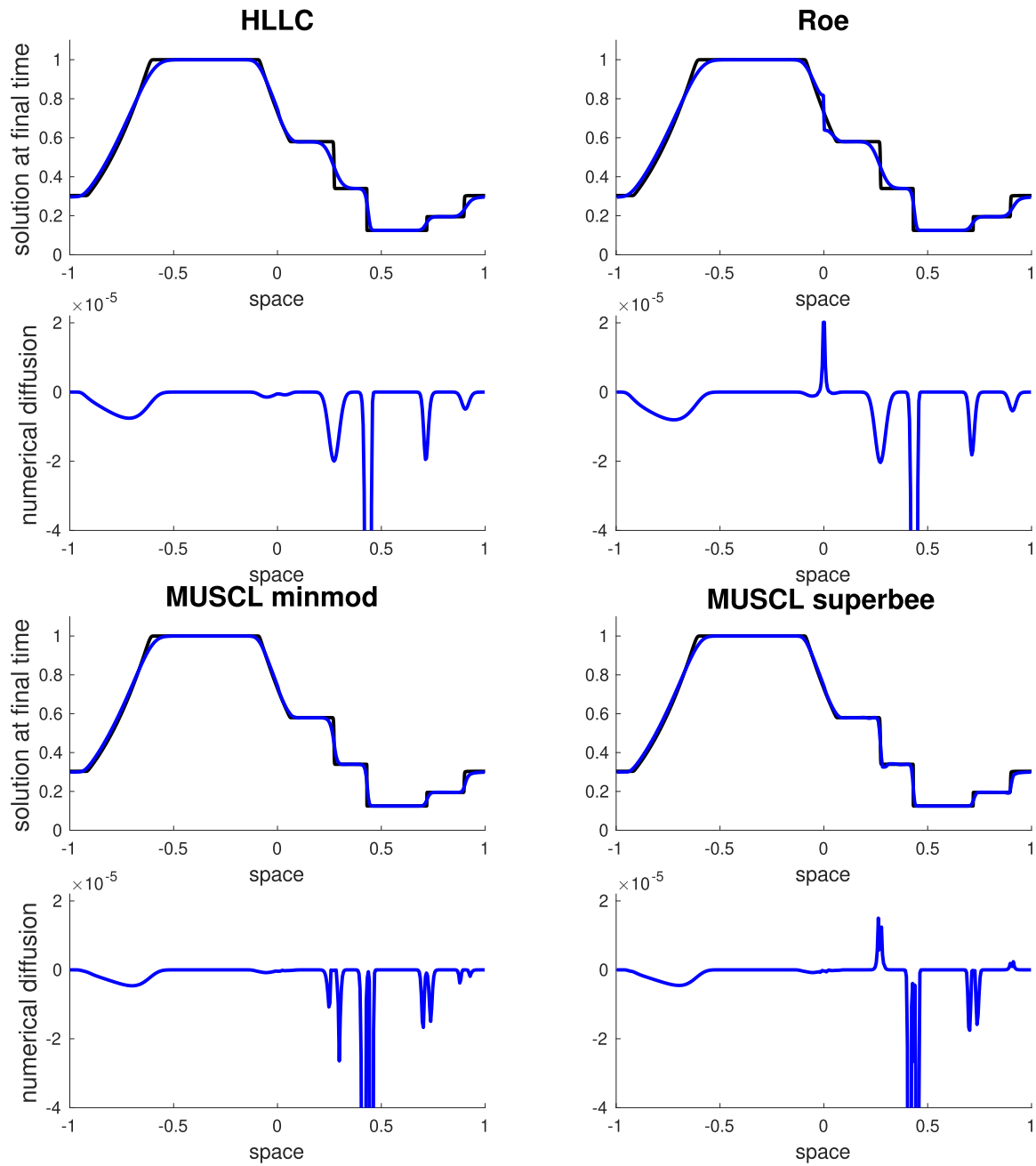


FIGURE 5. Densities (lines 1 and 3) and a posteriori numerical diffusion (lines 2 and 4) for several numerical schemes on Testcase (25).

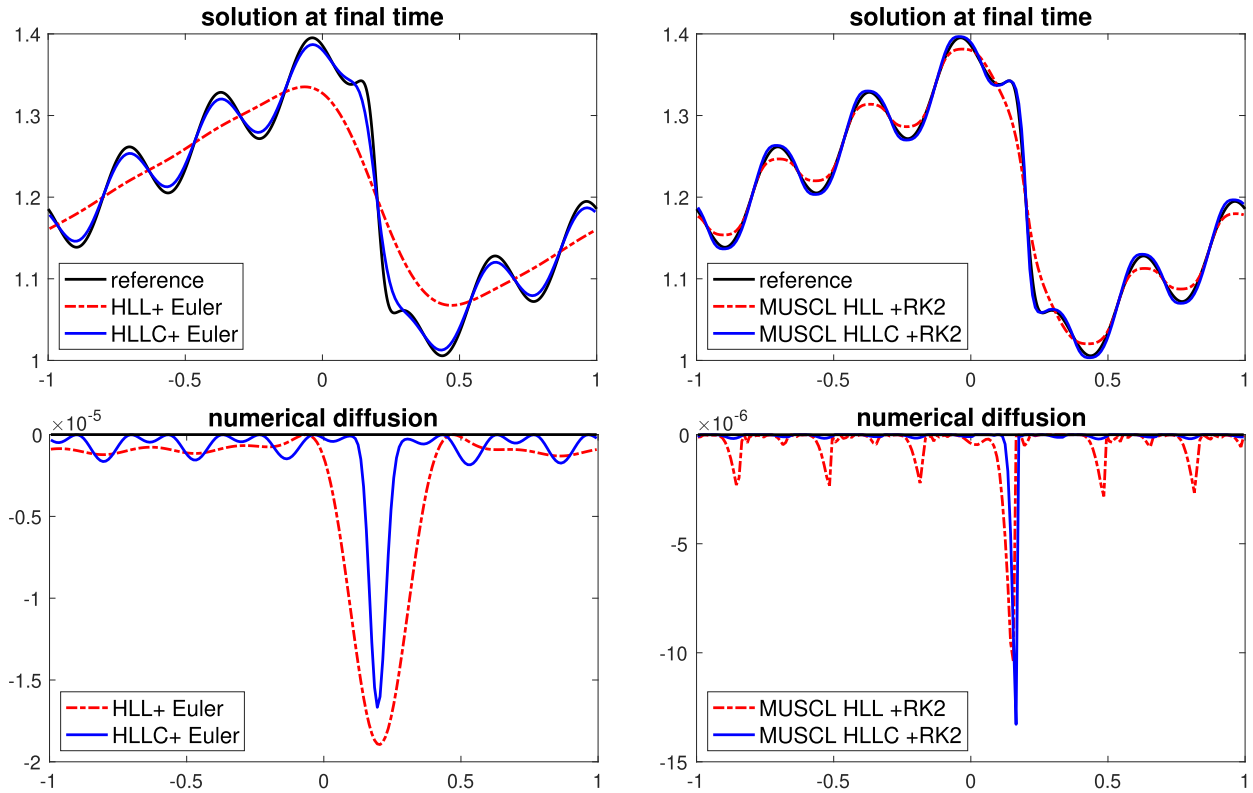


FIGURE 6. Densities (first line) and numerical diffusion (second line) for Testcase (26). The discontinuity is a slowly moving contact.

through the higher level of detail captured by the HLLC scheme at the final time. The quantification of numerical diffusion confirms it quantitatively, with lower and more localized  $D^{opt,n}$  values. The same holds for their second order extensions with a MUSCL procedure, using either the HLL or the HLLC flux as the underlying first order scheme. The slope limitation is on the conservative variables  $(\rho, \rho u, E)$  and we use a minmod limiter.

### 4.3. A naive *a priori* quantification of the numerical diffusion

In this section we turn to the practical consideration of the computational cost of the minimization presented above. The functional  $\mathcal{J}^D$  couples all the cells of the mesh because each numerical entropy flux  $\gamma_{j+1/2}$  appears in two consecutive terms of the sum. The cost grows quadratically with the meshsize, and becomes larger than the minute for meshes with more than 4000 cells, see Figure 7, left.

This observation disqualifies the use of this procedure in realistic codes where the quantification of  $D_j^n$  would be useful. In ocean global circulation models (OGCMs) for example, the meshes contain billion of cells and the numerical diffusion under scrutiny here is strongly linked with the question of spurious mixing identified as a key issue in such codes [13]. This explains why several works deal with its quantification. Some of them ([5, 18] and others) are directly related to the mathematical theory of discrete entropy inequality in the sense that specific choices of numerical entropy fluxes are imposed in (6) to define the numerical diffusion  $D_j^n$ . Here we propose a complementary approach: while these works focus on  $G_{j+1/2}^n$  at the risk of making a suboptimal choice and quantifying  $D_j^n$  incorrectly, we present a more robust choice of  $D_j^n$  at the cost of forgetting the numerical entropy fluxes  $G_{j+1/2}^n$ . As a consequence the theoretical results of the previous section are almost completely lost; however the quantitative results are very satisfactory and the computational cost is drastically

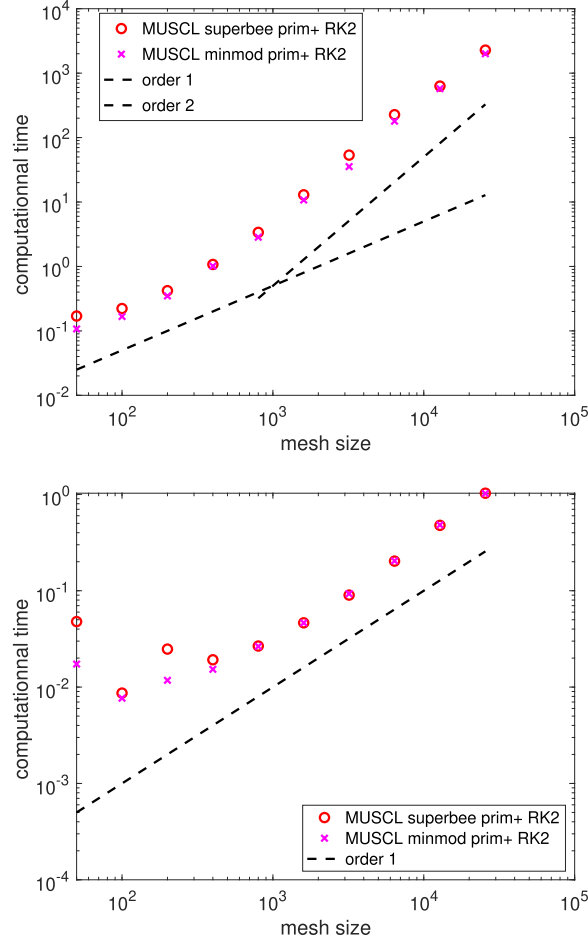


FIGURE 7. Computational cost of the minimization procedure (left) and of the *a priori* guess (28) at the first iteration of Testcase (26).

diminished, see Figure 7, right. We believe this fastest quantification could be used for some applications in the future, including the quantification of spurious mixing in OGCM or as a metric for mesh refinement.

In the derivation of the consistency part of the functional  $\mathcal{J}^C$ , we proposed bounds on the numerical entropy flux  $G_{j+1/2}^n \in [m_{j+1/2}^n, M_{j+1/2}^n]$ . Under the hypothesis that the bounds are correctly ordered, it follows that  $\underline{D}_j^n \leq D_j^{opt,n} \leq \overline{D}_j^n$ , with

$$\underline{D}_j^n = \eta(u_j^{n+1}) - \eta(u_j^n) + \lambda(m_{j+1/2}^n - M_{j-1/2}^n) \quad (27)$$

and

$$\overline{D}_j^n = \eta(u_j^{n+1}) - \eta(u_j^n) + \lambda(M_{j+1/2}^n - m_{j-1/2}^n).$$

An important point is that  $\underline{D}_j^n$  and  $\overline{D}_j^n$  are computationally affordable, see right of Figure 7. Indeed,  $m_{j+1/2}^n$  and  $M_{j+1/2}^n$  are computed with the scheme (4) on a small initial data centered around the interface  $j + 1/2$  with  $s_L + s_R$  cells on its left and on its right, see Lemma 1.

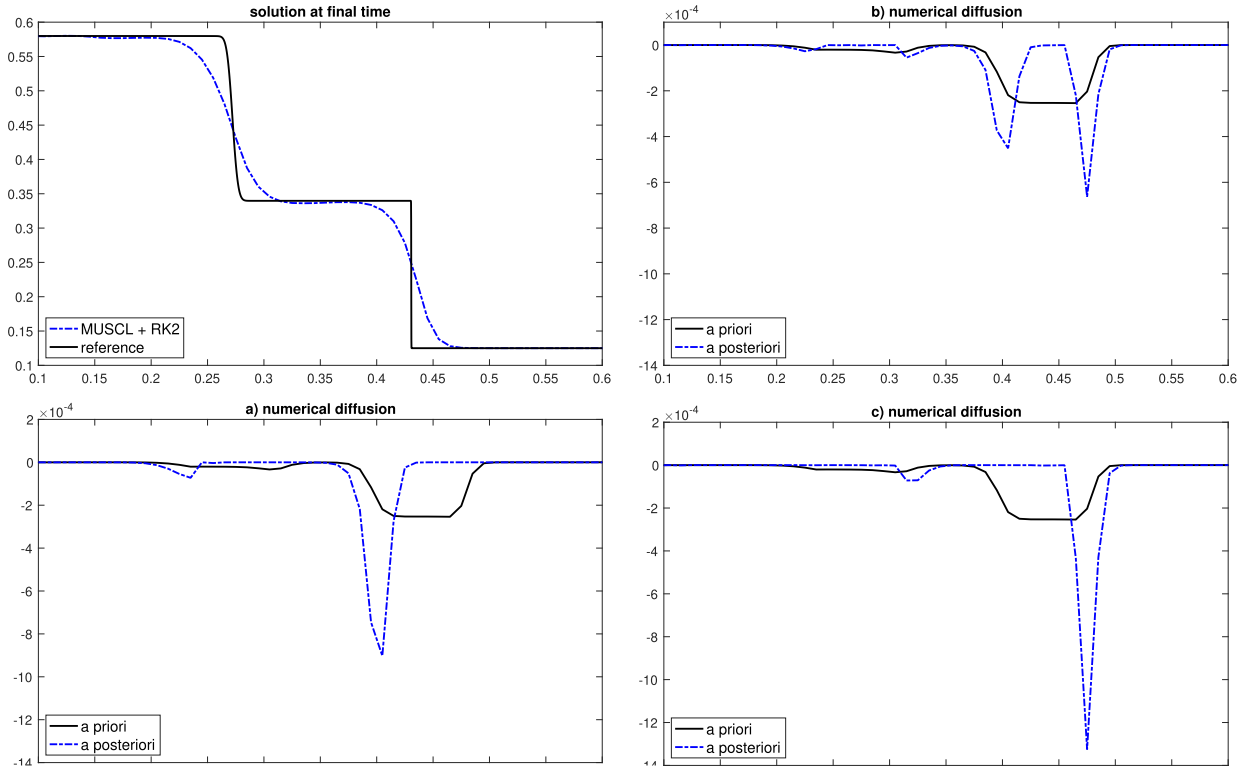


FIGURE 8. Comparison of  $D_j^{opt,n}$  and  $D_j^{a\,priori,n}$  on Testcase (25). For the computation of  $D_j^{opt,n}$  the optimization is initialized with (a)  $\gamma_{j+1/2}^n = m_{j+1/2}^n$  or (b)  $\gamma_{j+1/2}^n = \frac{1}{2} (m_{j+1/2}^n + M_{j+1/2}^n)$  or (c)  $\gamma_{j+1/2}^n = M_{j+1/2}^n$ .

We observe that the lower bound (27) is particularly interesting for two reasons. First, if the left hand side is positive it indicates that (5) cannot hold on this particular cell. We detail and capitalize on that idea in Section 5. Second, even though it is much lower than  $D_j^{opt,n}$ , the variations of  $\underline{D}_j^n$  are similar.

As a final stage, we perform a naive renormalization by a constant coefficient  $\alpha$  and define  $D_j^{a\,priori,n} = \alpha \underline{D}_j^n$  in such a way that the total amount of numerical diffusion is correct:

$$\sum_j D_j^{a\,priori,n} = \alpha \sum_j \underline{D}_j^n = \sum_j \eta(u_j^{n+1}) - \eta(u_j^n).$$

It yields to the *a priori* quantification of the numerical diffusion

$$D_j^{a\,priori,n} = \frac{\sum_k \eta(u_k^{n+1}) - \eta(u_k^n)}{\sum_k \underline{D}_k^n} \underline{D}_j^n. \quad (28)$$

On Figure 8 we compare  $D_j^{a\,priori,n}$  and  $D_j^{opt,n}$ . The *a posteriori* quantification  $D_j^{opt,n}$  depends on the initialization of the minimization procedure and produces narrow spikes of numerical diffusion. The *a priori* quantification  $D_j^{a\,priori,n}$  gives smoother results, with a diffusion spread out the whole length of the discrete

discontinuity. The computation of this quantity only requires the bounds  $m_{j+1/2}^n$  and  $M_{j+1/2}^n$  and does not require any optimization. We insist on the fact  $D_j^{a\text{ priori},n}$  is pertinent if one is only interested in the numerical diffusion and not in the numerical entropy fluxes, which are not provided by this approach.

**Remark 3.** One could think that the upper bound  $\bar{D}_j^n$  is also interesting. In particular if this quantity was nonpositive in all cells it would be a strong indication that the scheme is entropy satisfying. However, we observed that  $\bar{D}_j^n$  is in most cases strictly positive, even for first order entropy satisfying scheme. It makes this quantity far less interesting both qualitatively and theoretically.

### 5. AN ENTROPY STRESS TEST FOR NUMERICAL SCHEMES

In this section we go back to the classical question of the existence of discrete entropy inequality (5). We explore further the implication of the *a priori* consistency bounds  $m_{j+1/2}^n \leq G_{j+1/2}^n \leq M_{j+1/2}^n$ . It leads to a novel criterion and an algorithm which allows us to detect that entropy inequality (5) is not verified.

#### 5.1. Worst initial data in terms of entropy

**Proposition 4.** Consider a numerical flux  $\mathcal{F}$  with a stencil of  $s_L$  points on the left and  $s_R$  on the right. If there exists an initial data

$$(\dots, u_{-s_L}^0, u_{-s_L}^0, u_{-s_L+1}^0, \dots, u_0^0, \dots, u_{s_R-1}^0, u_{s_R}^0, u_{s_R}^0, \dots) \tag{29}$$

such that

$$\min \left( \eta(u_0^0) - \frac{\Delta t}{\Delta x} (m_{1/2}^0 - M_{-1/2}^0) - \eta(u_0^1), M_{1/2}^0 - m_{1/2}^0 \right) < 0, \tag{30}$$

where  $M_{1/2}^0$ ,  $m_{1/2}^0$  and  $M_{-1/2}^0$  are given by (9) with  $n = 0$ , then the finite volume scheme does not have a discrete entropy inequality.

*Proof.* By contraposition, let us assume that the scheme has a discrete entropy inequality. Therefore, for all initial data  $(u_j^n)_j$ , there exists consistent entropy numerical fluxes  $(G_{j+1/2})_j$  such that

$$\eta(u_j^1) \leq \eta(u_j^0) - \frac{\Delta t}{\Delta x} (G_{j+1/2} - G_{j-1/2}).$$

According to Lemma 1, we have  $m_{j+1/2}^0 \leq G_{j+1/2} \leq M_{j+1/2}^0$ . It follows especially  $m_{1/2}^0 \leq M_{-1/2}^0$  and

$$\eta(u_0^1) \leq \eta(u_0^0) - \frac{\Delta t}{\Delta x} (m_{1/2}^0 - M_{-1/2}^0).$$

Therefore the opposite of (30) is satisfied. This is true in particular for the initial data of the form (29). □

As a consequence, to determine if a scheme is entropy satisfying or not, we introduce the new functional

$$\mathcal{E}(u_{-s_L}^0, \dots, u_{s_R}^0) = \min \left( \eta(u_0^0) - \frac{\Delta t}{\Delta x} (m_{1/2}^0 - M_{-1/2}^0) - \eta(u_0^1), m_{1/2}^0 - M_{1/2}^0 \right).$$

For a hyperbolic system of  $d$  equations, it has  $d(s_L + s_R + 1)$  unknowns in  $\Omega$ . Following Proposition 4, if the minimization of  $\mathcal{E}$  returns a strictly negative result then the scheme is not entropy satisfying. In other words this new optimization constructs the worst initial data in terms of entropy.

**Remark 4.** The first part of (30) is similar to the definition (27) of  $\underline{D}_0^0$ . We already knew that if this quantity is positive the scheme is not entropy satisfying. This result goes further in that direction by restricting the size of the data to  $s_L + s_R + 1$  values and by shifting the optimization procedure on the initial data.

TABLE 1. Summary of the obtained results for different variants of MUSCL scheme.

	1st order scheme	reconstructed variables	entropy satisfying	most negative value of $\mathcal{E}$
(a)	HLL	conservative $(\rho, \rho u, E)$	no	-2.228
(b)	HLL	primitive $(\rho, u, p)$	no	$-4.352 \times 10^{-5}$
(c)	HLL	entropy $(\rho, u, s)$	yes?	$\geq 0$
(d)	HLLC	conservative $(\rho, \rho u, E)$	no	$-5.257 \times 10^{-1}$
(e)	HLLC	primitive $(\rho, u, p)$	no	$-3.518 \times 10^{-4}$
(f)	HLLC	entropy $(\rho, u, s)$	no	$-2.73 \times 10^{-4}$

## 5.2. Exploration of the RK2+MUSCL scheme

We apply this procedure to the gas dynamic (23) approximated by the common MUSCL scheme in space and a two step Runge-Kutta march in time. Some details are given in Section 4.2 and in the appendix. The six versions have been tested; counterexamples are constructed for five of them. The results are summarized in Table 1.

More precisely, we consider the limitations in primitive variables  $(\rho, v, p)$ , entropy variables  $(\rho, v, s)$  and conservative variables  $(\rho, q, E)$  described in [1] where the positiveness of the pressure and density are guaranteed. The underlying first order scheme is either the Rusanov scheme, the HLL or the HLLC scheme. We implemented two versions of the latter: one based on a pressure estimate and the other one on the estimation of extremal wave speeds. The CFL number is fixed to  $\alpha = 0.1$ . Details and references are given in the appendix. For the overall scheme we have  $s_L = s_R = 4$  and the minimization of Proposition 4 has a total of 18 parameters with the positiveness of density and pressure as constraints.

We focus on counterexamples with small total variation. The unknowns are constrained to

$$\forall j \in \{-4, \dots, 4\} \quad (\rho_j^0, u_j^0, p_j^0) \in (\rho_0, u_0, p_0) + [-0.1, 0.1]^3, \quad (31)$$

where  $(\rho_0, u_0, p_0)$  is selected randomly in a larger domain. The main reason for considering almost constant initial data is that if

$$u_{-s_L}^0 = u_{-s_L+1}^0 = \dots u_0^0 = \dots = u_{s_R-1}^0 = u_{s_R}^0,$$

then  $M_{\pm 1/2}^0 = m_{\pm 1/2}^0$ ,  $\eta(u_0^1) = \eta(u_0^0)$  and the three components of Proposition 4 vanish. Thus counterexamples may be found in the vicinity of constant states if the scheme is not entropy satisfying. This region somehow corresponds to regular solutions.

We chose randomly 30 000 initializations in the region

$$(\rho_0, u_0, p_0) \in [0.11, 5] \times [-0.2, 10] \times [0.11, 10] \quad (32)$$

and another 30 000 initializations in the smaller region

$$(\rho_0, u_0, p_0) \in [0.8, 1.2] \times [-0.2, 1] \times [0.8, 1.2]. \quad (33)$$

On local extremum the MUSCL scheme is identical to the chosen first order scheme. Thus it seems more likely to find counterexamples on monotonic data. Half of the initialization are modified so that each variable is increasing or decreasing. The rearrangement is made either in the primitive, entropy or conservative variables, in equal proportion.

The minimization procedure constructs numerous initial data that violate (5) in the first time step. They are gathered on Figure 9 when  $(\rho_0, u_0, p_0)$  belongs to the large region (32) and on Figure 10 when it belongs to the smaller region (33). Each point on these Figures represents the mean value

$$\left( \frac{1}{9} \sum_{k=-4}^4 \rho_j^0, \frac{1}{9} \sum_{k=-4}^4 v_j^0, \frac{1}{9} \sum_{k=-4}^4 p_j^0 \right)$$

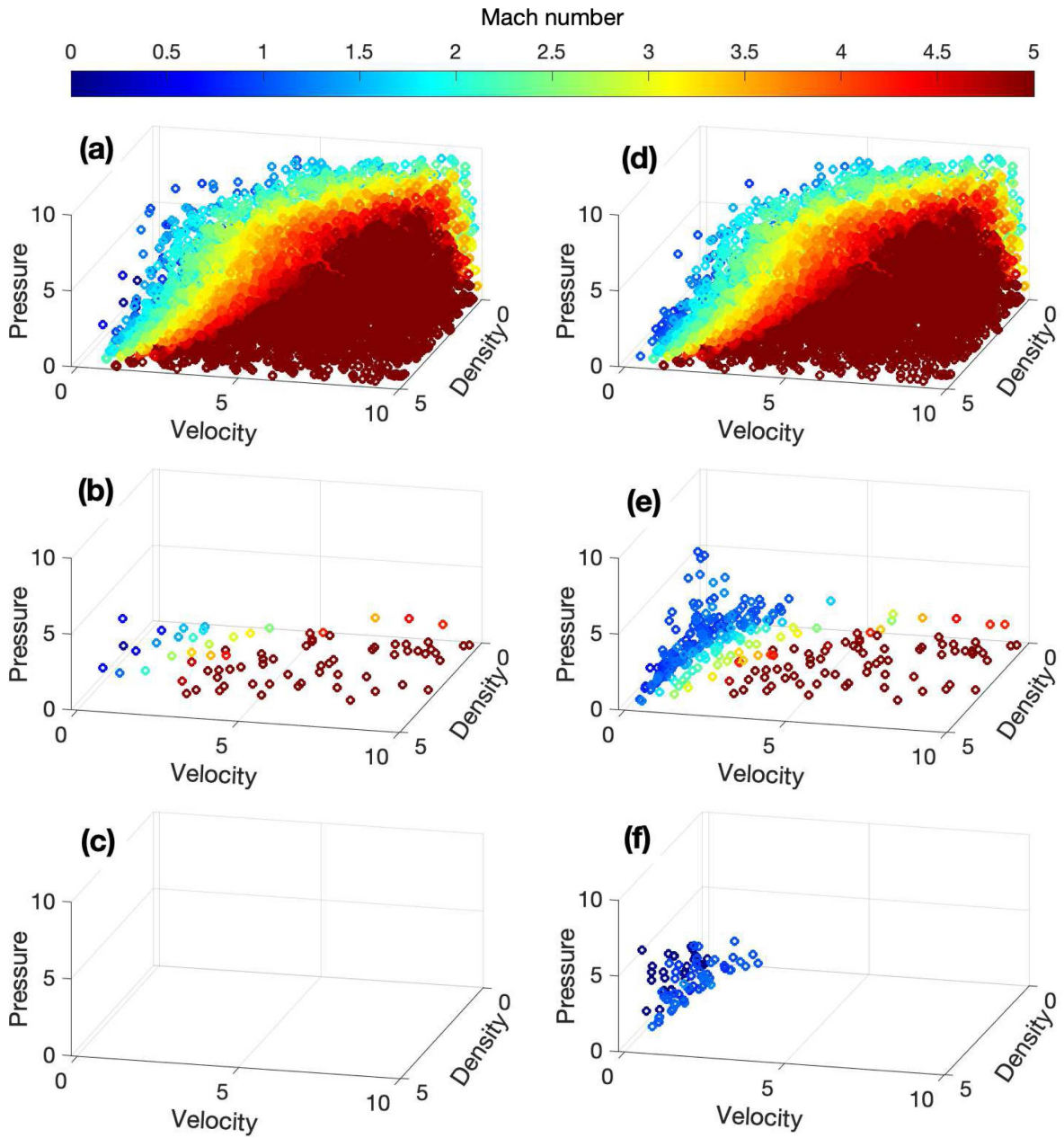


FIGURE 9. Counterexamples on the region  $[0.11, 5] \times [-0.2, 10] \times [0.11, 10]$ . The limited variables are  $(\rho, q, E)$  in the first line,  $(\rho, u, p)$  in the second line,  $(\rho, u, s)$  in the third line. The first order underlying scheme is HLL on the left and HLLC on the right.



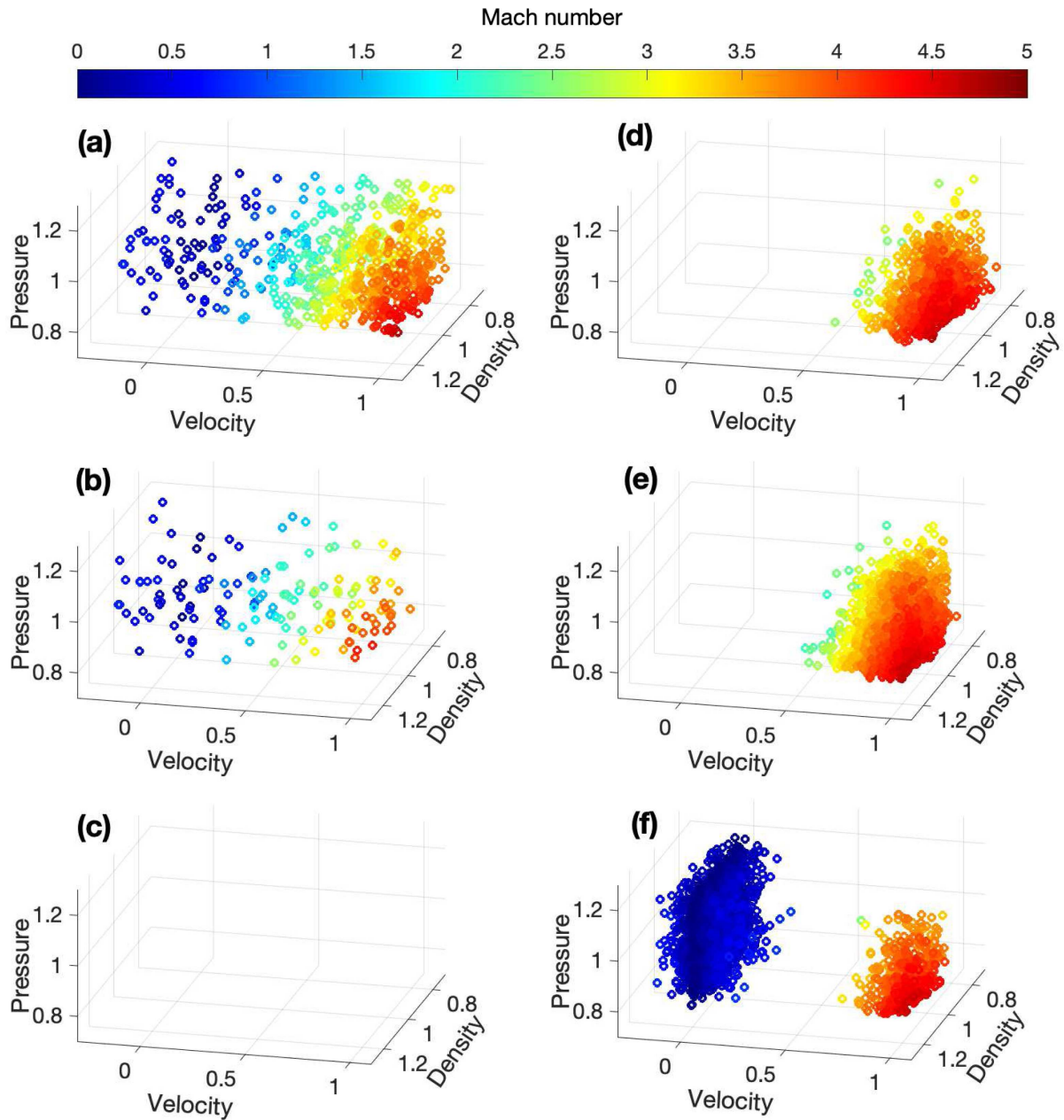


FIGURE 10. Zoom on a the region  $[0.8, 1.2] \times [-0.2, 1] \times [0.8, 1.2]$ . The limited variables are  $(\rho, q, E)$  in the first line,  $(\rho, u, p)$  in the second line,  $(\rho, u, s)$  in the third line. The first order underlying scheme is HLL on the left and HLLC on the right.

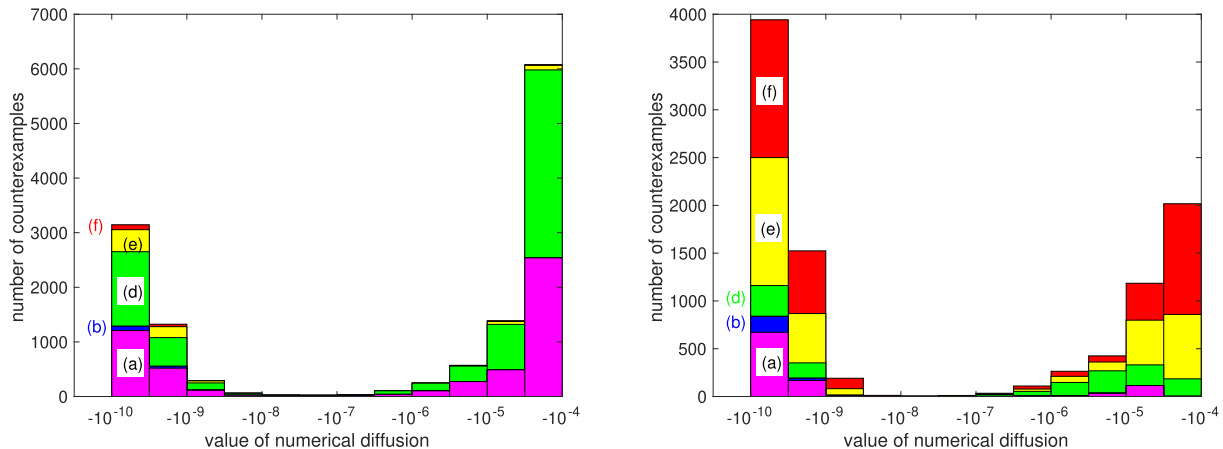


FIGURE 11. Histogram of the values of  $\mathcal{E}(u_{-4}^0, \dots, u_4^0)$  for the counterexamples found in the large region (32) (left) and in the small region (33) (right). A small number of counterexamples for scheme (b) and (f) (left) or (a) and (b) (right) are found in the region  $\mathcal{E} \approx -10^{-4}$ . The histogram are stopped artificially at  $-10^{-4}$ .

of a counterexample. The results for the two versions of the HLLC schemes are very similar and we include only the version based on the wave-speed estimates. Similarly the Rusanov and HLL schemes behave similarly and we only include the HLL scheme. We used the `fmincon` of MATLAB [23] to minimize the functional  $\mathcal{E}$  under the constraints (31). The histograms of the values corresponding to (30) for those counterexamples are given on Figure 11. We only kept counterexamples where this value is greater than  $-10^{-10}$  which is several order of magnitude larger than the machine precision estimated at  $2.2 \times 10^{-16}$ .

The influence of the choice of the set variables which are limited is striking on Figure 9, with more counterexamples for the conservative variables  $(\rho, q, E)$  than for the primitive variables  $(\rho, v, p)$ . There are even fewer counterexamples for the entropy variables  $(\rho, v, s)$ . This hierarchy is less clear on the zoom of Figure 10. The distribution on counterexamples around  $|v| = 0$  depends on the numerical choices.

The counterexample isolated on Figure 12 has a very small total variation and a null velocity. We checked that it remains a counterexample for smaller and smaller timestep (or equivalently for smaller and smaller CFL number). This may indicate that the limit  $\Delta t \rightarrow 0$  in (5) does not hold for the MUSCL+RK2 scheme based on a HLL Riemann solver with a limitation on the conservative variables.

On the other hand we see many counterexamples with large Mach number  $|v|/\sqrt{\frac{\gamma p}{\rho}}$  on Figure 9. The chosen first order scheme is not to blame for the lack of discrete entropy inequality. Indeed in that case the exact flux at interface  $j + 1/2$  is either  $f(u_j^0)$  or  $f(u_{j+1}^0)$  depending on the sign of the velocity. Most numerical schemes reproduce that.

Eventually Functional (30) remains nonnegative for all the random initializations when the limitation is in the entropy variable  $(\rho, v, s)$  and the first order scheme is HLL or Rusanov. To explore further if this scheme is entropy satisfying, we relaxed the constraint on the total variation and searched for counterexamples in the much larger domain

$$\forall j \in \{-4, \dots, 4\} \quad (\rho_j^0, u_j^0, p_j^0) \in [0.001, 10] \times [-50, 50] \times [0.001, 20].$$

This search was unsuccessful with 57 000 random initializations. Interestingly it holds for the overlimited version of the MUSCL scheme of [1] but also for the simpler original scheme. It supports the following conjecture.

**Conjecture 1.** *The RK2+MUSCL scheme for the gas dynamic equations (23) with*

$\rho$	0.9985	0.9991	0.9998	0.9998	1.0005	1.0005	1.0009	1.0018
$v$	0	0	0	0	0	0	0	0
$p$	1.2519	1.2531	1.2549	1.2566	1.2584	1.2601	1.2618	1.2635

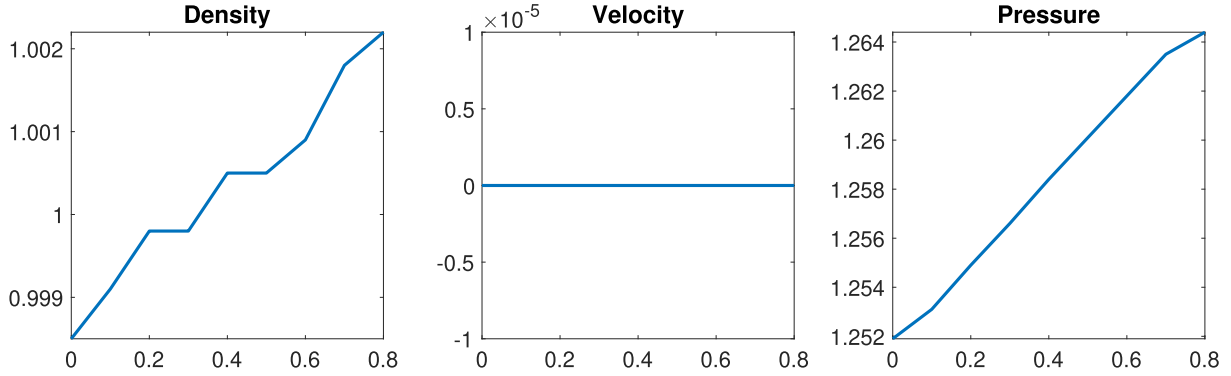


FIGURE 12. MUSCL+RK2 scheme with this initial data is non entropy satisfying, with a HLL first order scheme, a minmod limiter and a reconstruction in  $(\rho, u, p)$ . In this case  $\mathcal{E}(u_{-4}^0, \dots, u_4^0) \approx -4.10^{-9}$ .

- (1) a MUSCL scheme in space based on a HLL first order scheme with a minmod limiter on the entropy variables  $(\rho, u, s)$
- (2) the RK2 march in time described in the Appendix

verifies a discrete entropy inequality for a CFL number  $\alpha = 0.1$ .

The same type of numerical experiment has been performed for the scalar Burgers equation (21). Three entropy satisfying two points schemes (Rusanov, Osher and Godunov) and three non entropy satisfying schemes (Roe, Lax-Wendroff and Mac Cormack) have been tested. We performed 30 000 optimization of the functional  $\mathcal{E}$  and found no negative values for the entropy satisfying schemes, and counterexamples for the non entropy satisfying schemes, which supports the following conjecture.

**Conjecture 2.** Consider the RK2+MUSCL scheme for the Burgers equation (21) with a minmod limiter on  $u$  and the RK2 march in time described in the appendix. For a Courant number up to 1, this scheme is entropy satisfying when the chosen underlying first order scheme is entropy satisfying, and non entropy satisfying when it is not.

## APPENDIX A. DESCRIPTION OF THE NUMERICAL SCHEMES

The Rusanov scheme (first order)

This scheme is one of the simplest approximate Riemann Solver, see [14]. The numerical flux is

$$\mathcal{F}(u_L, u_R) = \frac{f(u_L) + f(u_R)}{2} - \frac{A(u_L, u_R)}{2} (u_R - u_L) \quad (\text{A.1})$$

and the numerical entropy flux is

$$\mathcal{G}(u_L, u_R) = \frac{G(u_L) + G(u_R)}{2} - \frac{A(u_L, u_R)}{2} (\eta(u_R) - \eta(u_L)). \quad (\text{A.2})$$

The scalar quantity  $A(u_L, u_R)$  should be large enough. For scalar equation, we set  $A(u_L, u_R) = \max(|f'(u_L)|, |f'(u_R)|)$  and for the equations of gas dynamic (23) we take

$$A(u_L, u_R) = \max\left(|v_L| + \sqrt{\frac{\gamma p_L}{\rho_L}}, |v_R| + \sqrt{\frac{\gamma p_R}{\rho_R}}\right).$$

*The HLL scheme (first order)*

This scheme is also a simple two waves approximate Riemann solver. We present the scheme for gas dynamic. Following Section 10.3 of [24] we introduce the sound speeds  $a_L = \sqrt{\gamma p_L / \rho_L}$  and  $a_R = \sqrt{\gamma p_R / \rho_R}$  and the Roe averages  $\bar{a} = \frac{\sqrt{\rho_L} a_L + \sqrt{\rho_R} a_R}{\sqrt{\rho_L} + \sqrt{\rho_R}}$  and  $\bar{v} = \frac{\sqrt{\rho_L} v_L + \sqrt{\rho_R} v_R}{\sqrt{\rho_L} + \sqrt{\rho_R}}$ . The wave speed are estimates as  $S_L = \min(v_L - a_L, v_R - a_R, \bar{v} - \bar{a})$  and  $S_R = \min(v_L + a_L, v_R + a_R, \bar{v} + \bar{a})$  and the flux is given by equation (10.21) of [24]. With  $-S_L = S_R = A(u_L, u_R)$  the HLL scheme is the Rusanov scheme.

*The HLLC scheme for gas dynamics (first order)*

This scheme is a commonly used approximate Riemann solver for (23). This solver is exact on some important particular solutions of (23), the stationary contact discontinuities. We implemented two versions of the schemes and obtained similar results. The first one corresponds to Paragraph 10.4.2 of [24]. The extremal wave speed are estimated with (10.49) of [24] and the pressure and velocities are constant in the “star region”. The second one is based on a pressure estimate and corresponds to Paragraph 10.6, variant 1 of [24].

*The Roe scheme (first order)*

In the case of scalar conservation laws, the numerical scheme is

$$\mathcal{F}(u_L, u_R) = f(u_L) \mathbf{1}_{\sigma \geq 0} + f(u_R) \mathbf{1}_{\sigma < 0}, \quad \sigma = \frac{f(u_R) - f(u_L)}{u_R - u_L}.$$

We refer the reader to [22] and Chapter 11 of [24] for the presentation of the Roe scheme for the Euler equations (23). This scheme does not have a discrete entropy inequality since it preserves entropy violating shock ( $f(u_L) = f(u_R)$ ,  $G(u_R) > G(u_L)$ ). In can also produce negative pressure [12], in which case the simulation fails entirely.

**The MUSCL scheme (first order)**

The MUSCL procedure is a commonly used procedure to obtain a second ordre scheme in space, and can be easily combined with a second order time scheme such as a 2 step Runge-Kutta method. It is a 4-points flux with  $s_L = s_R = 2$ .

For scalar conservation laws, the procedure is the following. At the beginning of each time step, the constant value  $u_j^n$  in cell  $\mathcal{C}_j = [x_j - \frac{\Delta x}{2}, x_j + \frac{\Delta x}{2}]$  is replaced by the affine function

$$x \mapsto u_j^n + \sigma_j^n (x - x_j),$$

where  $\sigma_j \in \mathbf{R}$  is a slope, determined in such a way that no new extrema are created and the scheme remains total variation diminishing. A common choice is

$$\begin{aligned} \sigma_j^n &= \text{minmod}(u_j^n - u_{j-1}^n, u_{j+1}^n - u_j^n), \\ &= \max(0, \min(u_j^n - u_{j-1}^n, u_{j+1}^n - u_j^n)) + \min(0, \max(u_j^n - u_{j-1}^n, u_{j+1}^n - u_j^n)). \end{aligned}$$

Other limiters are possible, see Section 13.7.3 of [24]. The MUSCL flux is given by

$$F_{j+1/2} = \mathcal{F}(u_{j,+}^n, u_{j+1,-}^n), \quad u_{j,\pm}^n = u_j^n \pm \frac{\Delta x}{2} \sigma_j^n, \tag{A.3}$$

where  $\mathcal{F}$  is a first order two points scheme.

In the case of hyperbolic system with  $p > 1$ , this strategy is mimicked componentwise. For the Euler equation (23) we followed the strategy of [1]. The reconstruction with a minmod limiter can be applied on the three conservative variables  $(\rho, \rho v, E)$ . It is also common to reconstruct the primitive variables  $(\rho, v, p)$  or in the entropy variables  $\left(\rho, v, s = \frac{p}{\rho^\gamma}\right)$  and to deduce the values  $(\rho v)_{j,\pm}^n = \rho_{j,\pm}^n v_{j,\pm}^n$  with

$$E_{j,\pm}^n = \frac{p_{j,\pm}^n}{\gamma - 1} + \frac{\rho_{j,\pm}^n (v_{j,\pm}^n)^2}{2} \quad \text{or} \quad E_{j,\pm}^n = \frac{(\rho_{j,\pm}^n)^\gamma s_{j,\pm}^n}{\gamma - 1} + \frac{\rho_{j,\pm}^n (v_{j,\pm}^n)^2}{2}.$$

A variation of the discrete entropy inequality (5) is obtained in [1] but is not sufficient to obtain a Lax-Wendroff theorem.

### A second order Runge-Kutta method in time

To achieve second order in time one can use a two step Runge-Kutta method in time. Consider a numerical flux  $\mathcal{F}$  with a stencil of  $s_L$  points to the left and  $s_R$  points to the right. Two substeps are computed

$$\bar{u}_j^n = u_j^n - \frac{\Delta t}{\Delta x} (\mathcal{F}(u_{j-s_L+1}^n, \dots, u_{j+s_R}^n) - \mathcal{F}(u_{j-s_L}^n, \dots, u_{j+s_R-1}^n)),$$

$$\bar{\bar{u}}_j^n = \bar{u}_j^n - \frac{\Delta t}{\Delta x} (\mathcal{F}(\bar{u}_{j-s_L+1}^n, \dots, \bar{u}_{j+s_R}^n) - \mathcal{F}(\bar{u}_{j-s_L}^n, \dots, \bar{u}_{j+s_R-1}^n)).$$

The final update is  $u_j^{n+1} = \frac{u_j^n + \bar{\bar{u}}_j^n}{2}$ .

The whole procedure rewrites in the compact form (4) with

$$\bar{\mathcal{F}}(u_{j-2s_L+1}^n, \dots, u_{j+2s_R}^n) = \frac{\mathcal{F}(u_{j-s_L+1}^n, \dots, u_{j+s_R}^n) + \mathcal{F}(\bar{u}_{j-s_L+1}^n, \dots, \bar{u}_{j+s_R}^n)}{2}.$$

If the numerical flux  $\mathcal{F}$  is entropy satisfying with a numerical entropy flux  $\mathcal{G}$ , the same holds for  $\bar{\mathcal{F}}$  with the numerical entropy flux

$$\bar{\mathcal{G}}(u_{j-2s_L+1}^n, \dots, u_{j+2s_R}^n) = \frac{\mathcal{G}(u_{j-s_L+1}^n, \dots, u_{j+s_R}^n) + \mathcal{G}(\bar{u}_{j-s_L+1}^n, \dots, \bar{u}_{j+s_R}^n)}{2}.$$

For scalar conservation laws  $d = 1$ , the maximum of the wave speeds decreases with time, thus the CFL restriction for the computation of  $\bar{u}_j$  is more restrictive than the one for  $\bar{\bar{u}}$ . This is not true when  $p \geq 2$ , and we adopt the time evolution of equation (2.12) of [1], with the slight modification that  $\Delta t_2$  cannot exceed  $\Delta t_1$ .

### REFERENCES

- [1] C. Berthon, Stability of the MUSCL schemes for the Euler equations. *Commun. Math. Sci.* **3** (2005) 133–157.
- [2] C. Berthon and V. Desveaux, An entropy preserving MOOD scheme for the Euler equations. *Int. J. Finite.* **11** (2014) 39.
- [3] F. Bouchut, Nonlinear Stability of Finite Volume Methods for Hyperbolic Conservation Laws and Well-balanced Schemes for Sources. *Frontiers in Mathematics*. Birkhäuser Verlag, Basel (2004).
- [4] F. Bouchut, C. Bourdarias and B. Perthame, A MUSCL method satisfying all the numerical entropy inequalities. *Math. Comput.* **65** (1996) 1439–1461.
- [5] H. Burchard and H. Rennau, Comparative quantification of physically and numerically induced mixing in ocean models. *Ocean Model.* **20** (2008) 293–311.
- [6] C. Chalons and P.G. LeFloch, A fully discrete scheme for diffusive-dispersive conservation laws. *Numer. Math.* **89** (2001) 493–509.
- [7] S. Clain, S. Diot and R. Loubere, A high-order finite volume method for systems of conservation laws-multi-dimensional optimal order detection (mood). *J. Comput. Phys.* **230** (2011) 4028–4050.

- [8] P. Colella and P.R. Woodward, The piecewise parabolic method (ppm) for gas-dynamical simulations. *J. Comput. Phys.* **54** (1984) 174–201.
- [9] F. Coquel and P.G. LeFloch, An entropy satisfying MUSCL scheme for systems of conservation laws. *Numer. Math.* **74** (1996) 1–33.
- [10] F. Couderc, A. Duran and J.P. Vila, An explicit asymptotic preserving low froude scheme for the multilayer shallow water model with density stratification. *J. Comput. Phys.* **343** (2017) 235–270.
- [11] S. Diot, S. Clain and R. Loubere, Improved detection criteria for the multi-dimensional optimal order detection (mood) on unstructured meshes with very high-order polynomials. *Comput. Fluids* **64** (2012) 43–63.
- [12] B. Einfeldt, C. Munz, P.L. Roe and B. Sjogreen, On godunov-type methods near low-densities. *J. Comput. Phys.* **92** (1991) 273–295.
- [13] B. Fox-Kemper, A. Adcroft, C.W. Böning, E.P. Chassignet, E. Curchitser, G. Danabasoglu, C. Eden, M.H. England, R. Gerdes, R.J. Greatbatch and S.M. Griffies, Challenges and prospects in ocean circulation models. *Front. Mar. Sci.* **6** (2019) 65.
- [14] E. Godlewski and P.-A. Raviartm, Numerical approximation of hyperbolic systems of conservation laws, 2nd edition. In Vol. 118 of *Applied Mathematical Sciences*. Springer-Verlag, New York (2021).
- [15] A. Harten and S. Osher, Uniformly high-order accurate nonoscillatory schemes .1. *SIAM J. Numer. Anal.* **24** (1987) 279–309.
- [16] A. Hildebrand and S. Mishra, Entropy stable shock capturing space-time discontinuous Galerkin schemes for systems of conservation laws. *Numer. Math.* **126** (2014) 103–151.
- [17] R.M. Holmes, J.D. Zika, S.M. Griffies, A.M.C.C. Hogg, A.E. Kiss and M.H. England, The geography of numerical mixing in a suite of global ocean models. *J. Adv. Model. Earth Syst.* **13** (2021) e2020MS002333.
- [18] K. Klingbeil, M. Mohammadi-Aragh, U. Graewe and H. Burchard, Quantification of spurious dissipation and mixing - discrete variance decay in a finite-volume framework. *Ocean Model.* **81** (2014) 49–64.
- [19] P. Lax and B. Wendroff, Systems of conservation laws. *Commun. Pure Appl. Math.* **13** (1960) 217–237.
- [20] L. Martaud, M. Badsì, C. Berthon, A. Duran and K. Saleh, Global entropy stability for class of unlimited high-order schemes for hyperbolic systems of conservation laws. Preprint: [arXiv:hal-03206727](https://arxiv.org/abs/hal-03206727) (2021).
- [21] B. Perthame and C.W. Shu, On positivity preserving finite volume schemes for Euler equations. *Numer. Math.* **73** (1996) 119–130.
- [22] P.L. Roe, Approximate riemann solvers, parameter vectors, and difference-schemes. *J. Comput. Phys.* **43** (1981) 357–372.
- [23] The Mathworks, Inc., Natick, Massachusetts. MATLAB version 9.11.0.1769968 (R2021b) (2021).
- [24] E.F. Toro, Riemann Solvers and Numerical Methods for Fluid Dynamics, 3rd edition. A practical introduction, Springer-Verlag, Berlin (2009).
- [25] B. Van Leer, Towards the ultimate conservative difference scheme. 5. 2nd-order sequel to Godunovs method. *J. Comput. Phys.* **32** (1979) 101–136.
- [26] X. Zhang and C.-W. Shu, Positivity-preserving high order finite difference weno schemes for compressible Euler equations. *J. Comput. Phys.* **231** (2012) 2245–2258.



**Please help to maintain this journal in open access!**

This journal is currently published in open access under the Subscribe to Open model (S2O). We are thankful to our subscribers and supporters for making it possible to publish this journal in open access in the current year, free of charge for authors and readers.

Check with your library that it subscribes to the journal, or consider making a personal donation to the S2O programme by contacting [subscribers@edpsciences.org](mailto:subscribers@edpsciences.org).

More information, including a list of supporters and financial transparency reports, is available at <https://edpsciences.org/en/subscribe-to-open-s2o>.