



HAL
open science

MUSCL reconstruction and Haar wavelets

Laurent Gosse

► **To cite this version:**

Laurent Gosse. MUSCL reconstruction and Haar wavelets. Communications in Mathematical Sciences, 2015, 13 (6), pp.1501–1514. <10.4310/CMS.2015.v13.n6.a7>. <hal-01061889>

HAL Id: hal-01061889

<https://hal.science/hal-01061889v1>

Submitted on 8 Sep 2014

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.



HAL Authorization

MUSCL RECONSTRUCTION AND HAAR WAVELETS

LAURENT GOSSE *

Abstract. MUSCL extensions (*Monotone Upstream-centered Schemes for Conservation Laws*) of the Godunov numerical scheme for scalar conservation laws are shown to admit a rather simple reformulation when recast in the formalism of the Haar multi-resolution analysis of $L^2(\mathbb{R})$. By pursuing this wavelet reformulation, a seemingly new MUSCL-WB scheme is derived for advection-reaction equations which is stable for a Courant number up to 1 (instead of roughly $\frac{1}{2}$). However these high-order reconstructions aren't likely to improve the handling of delicate nonlinear wave interactions in the involved case of systems of Conservation/Balance laws.

Key words. Godunov scheme, Haar wavelets, Multi-resolution Analysis, MUSCL reconstruction, Second-order resolution (SOR), slope-limiter, wave interactions, Well-balanced (WB) scheme.

Subject classifications. 65M06, 65T60, 35Q35.

1. Introduction

The goal of this text is to recast the widely-used MUSCL high-order schemes for computing the entropy solution of a one-dimensional convex scalar conservation law,

$$\partial_t u + \partial_x f(u) = 0, \quad u(t=0, \cdot) = u_0 \in L^1 \cap BV(\mathbb{R}), \quad (t, x) \in \mathbb{R}_*^+ \times \mathbb{R}, \quad (1.1)$$

into the formalism of a multi-resolution analysis of $L^2(\mathbb{R})$ derived from the Haar wavelet. For convenience, we shall always work with a Cartesian uniform computational grid, determined by a space-step Δx and a time-step Δt satisfying the standard homogeneous CFL restriction. Let $J \in \mathbb{Z}$ be fixed, we select in a first stage:

$$\Delta x = 2^{-J}, \quad \max |f'(u)| \Delta t \leq \frac{\Delta x}{2} = 2^{-J-1}.$$

1.1. The standard Godunov scheme

By defining $C_k = (x_{k-\frac{1}{2}}, x_{k+\frac{1}{2}})$ as the generic computational cell of width Δx centered on $x_k = k\Delta x$, $k \in \mathbb{Z}$, one may apply the Divergence Theorem on any rectangle $C_k \times (t^n, t^{n+1})$ in order to derive a mass-preserving numerical scheme for (1.1):

$$\int_{C_k} u(t^{n+1}, x) dx = \int_{C_k} u(t^n, x) dx - \int_{t^n}^{t^{n+1}} f\left(u(\tau, x_{k+\frac{1}{2}})\right) - f\left(u(\tau, x_{k-\frac{1}{2}})\right) d\tau.$$

This is equivalent to writing down the weak formulation of (1.1) with test-functions being the indicator functions of C_k , denoted $\chi(C_k)$. Hereafter, we use the standard notation $u_k^n = \int_{C_k} u(t^n, x) \frac{dx}{\Delta x}$. Yet the observation leading to Godunov scheme is the following: in case $u(t^n, \cdot)$ is constant on each computational cell C_k , then the boundary flux terms can be explicitly computed by resolving all the discontinuities, that is to say, Riemann problems at both interfaces $x_{k\pm\frac{1}{2}}$. Moreover, since Riemann fans $\omega(\frac{x}{t}; u^L, u^R)$ display a self-similar structure, one has a nice simplification,

$$\int_{t^n}^{t^{n+1}} f\left(u(\tau, x_{k+\frac{1}{2}})\right) d\tau = \Delta t \cdot f\left(\omega(0; u_k^n, u_{k+1}^n)\right). \quad (1.2)$$

*IAC-CNR "Mauro Picone", Via dei Taurini 19, 00185 Roma, Italy, (l.gosse@ba.iac.cnr.it).

When seeking an explicit time-marching algorithm, one may want to get rid of the Riemann solution ω , and it can be shown that (1.2) defines a smooth and consistent numerical flux function denoted by $F: \mathbb{R}^2 \rightarrow \mathbb{R}$,

$$\forall u, v \in \mathbb{R}^2, \quad F(u, v) = f(\omega(0; u, v)) = \begin{cases} \min_{u \leq \xi \leq v} f(\xi) & \text{if } u \leq v \\ \max_{v \leq \xi \leq u} f(\xi) & \text{if } u > v \end{cases} \quad (1.3)$$

Now, let's consider another formulation of this numerical scheme: denote by \mathcal{P}_J

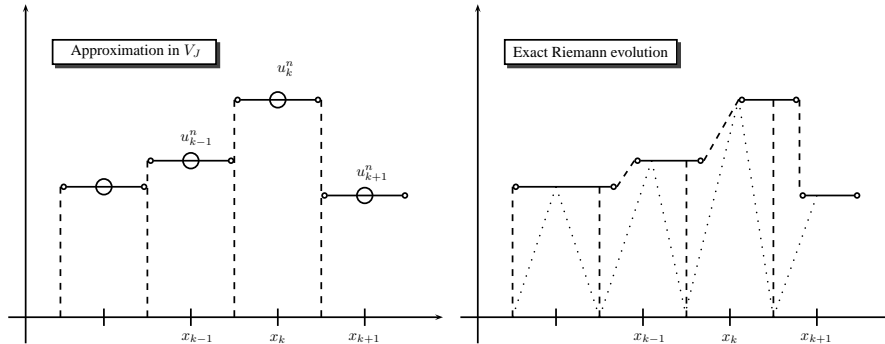


FIG. 1.1. Piecewise constant approximation (left) and exact Riemann fans (right) in (1.4).

the L^2 -projector onto the space of piecewise-constant functions on the computational grid, and $\mathcal{E}_J(t)$ the exact Riemann evolution operator¹ at each interface point $x_{k+\frac{1}{2}} = (k + \frac{1}{2})\Delta x = (k + \frac{1}{2})2^{-J}$ of the grid, the Godunov approximation reads:

$$\forall n \in \mathbb{N}, \quad u^{\Delta x}(t^n, \cdot) = [\mathcal{P}_J \circ \mathcal{E}_J(\Delta t)]^n \mathcal{P}_J(u_0). \quad (1.4)$$

Godunov wipes all the details at a finer scale than the grid by layered local averaging.

1.2. Scaling function and the Multi-Resolution formalism

We recall what is a *Multi-Resolution Analysis* (MRA, [25]) as $L^1 \cap BV(\mathbb{R}) \subset L^2(\mathbb{R})$.

DEFINITION 1.1. A sequence of nested (scale-limited) subspaces $V_j \subset L^2(\mathbb{R})$ is called a **Multi-Resolution Analysis** of $L^2(\mathbb{R})$ if $\{0\} \subset \dots \subset V_{-1} \subset V_0 \subset V_1 \subset \dots \subset L^2(\mathbb{R})$. Moreover, the following properties must hold:

- for all $f \in L^2(\mathbb{R})$, $\|\mathcal{P}_j f - f\|_{L^2} \rightarrow 0$ as $j \rightarrow +\infty$ also, $\mathcal{P}_j f \rightarrow 0$ as $j \rightarrow -\infty$.
- if $f(x) \in V_j$, then $f(\frac{x}{2}) \in V_{j-1}$ and for all $k \in \mathbb{Z}$, $f(x - 2^j k) \in V_j$.
- there exists a shift-invariant orthonormal base of V_0 given by the **scaling function** $\varphi_k(x) = \varphi(x - k)$ for $k \in \mathbb{Z}$.

Hence \mathcal{P}_j stands for the orthogonal projector onto the subspace V_j . Wavelet spaces W_j are defined as the **orthogonal complement** of V_j in V_{j+1} : $V_{j+1} = V_j \oplus W_j$. From φ_k , the base of V_0 , one deduces a base of V_j , $j \in \mathbb{Z}$, by simple dilatation,

$$\forall k \in \mathbb{Z}, \quad \varphi_{j,k}(x) = \sqrt{2^j} \varphi_k(2^j x) = \sqrt{2^j} \varphi(2^j x - k). \quad (1.5)$$

¹Actually the scale index J isn't indispensable and one may denote \mathcal{E} as the exact (entropy) solution operator: however, we shall keep on displaying J (or $J+1$) hereafter for ease of reading.

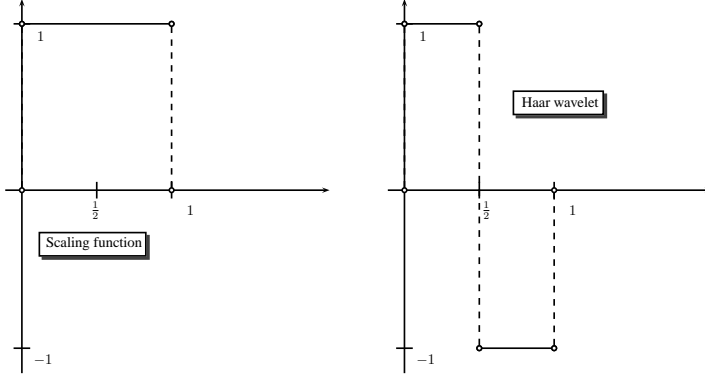


FIG. 1.2. Haar's scaling function (left) and wavelet (right).

Thus, the orthogonal projection of f onto the **scale-limited subspace** V_j reads:

$$\mathcal{P}_j f = \sum_{k \in \mathbb{Z}} \langle f, \varphi_{j,k} \rangle \varphi_{j,k}, \quad \langle f, \varphi_{j,k} \rangle = \int_{\mathbb{R}} f(x) \varphi_{j,k}(x) dx, \quad (1.6)$$

which is the best approximation of f in V_j in the least-squares sense. In the context of applications to the Godunov scheme (1.4), there exists a scaling function (hence a MRA) of particular interest, which is part of the ‘‘Haar system’’ (see Fig. 1.2),

$$\varphi(x + \frac{1}{2}) = \chi([0, 1]), \quad \psi(x + \frac{1}{2}) = \chi\left([0, \frac{1}{2})\right) - \chi\left([\frac{1}{2}, 1]\right). \quad (1.7)$$

The shift factor of $\frac{1}{2}$ is necessary in order to fit with the notation of (1.4), that is, to ensure that the Haar scaling functions match the indicator of each computational cell C_k . A simple observation is that, thanks to the definition (1.7), for a computational grid for which $\Delta x = 2^{-J}$, the Godunov approximation reads now:

$$\forall n \in \mathbb{N}, \quad u^{\Delta x}(t^n, \cdot) = \frac{1}{\sqrt{2^J}} \sum_{k \in \mathbb{Z}} u_k^n \varphi_{J,k}, \quad (1.8)$$

where the initial coefficients are obtained through:

$$u_k^{n=0} = \frac{1}{\Delta x} \int_{C_k} u_0(x) dx = \sqrt{2^J} \int_{\mathbb{R}} u_0(x) \sqrt{2^J} \varphi(2^J x - k) dx = 2^{\frac{J}{2}} \langle u_0, \varphi_{J,k} \rangle.$$

1.3. Main Theorem and plan of the Note

Having at hand the expressions (1.7) of both the Haar father and mother wavelets, we see that MUSCL reconstructions rewrite as a set of (mother) wavelet corrections:

THEOREM 1.2. *Let $\Delta x = 2^{-J}$ be the grid's parameter and $u_0 \in L^1 \cap BV(\mathbb{R})$ be Cauchy data for (1.1). For any TVD-admissible slope-limiter function $\phi: \mathbb{R} \rightarrow [0, 2]$, let's R_ϕ stand for the associated MUSCL reconstruction (2.2), then $R_\phi: V_J \rightarrow V_{J+1}$ and*

$$R_\phi \circ \mathcal{P}_J(u_0) = 2^{-\frac{J}{2}} \left(\sum_{k \in \mathbb{Z}} u_k^0 \varphi_{J,k} - \sum_{k \in \mathbb{Z}} \phi(r_k^0) \frac{u_{k+1}^0 - u_k^0}{2} \psi_{J,k} \right), \quad r_k^0 = \frac{u_k^0 - u_{k-1}^0}{u_{k+1}^0 - u_k^0}. \quad (1.9)$$

Accordingly, $Id - R_\phi \circ \mathcal{P}_J$ maps V_J into W_J and there is a “back-projection” property:

$$\mathcal{P}_J \circ R_\phi \circ \mathcal{P}_J = \mathcal{P}_J, \quad L^2(\mathbb{R}) \rightarrow V_J. \quad (1.10)$$

Both the equations (1.9) and (1.10) imply that for any limiter ϕ , MUSCL reconstructions induce only a “fluctuation component” in W_J , so it can’t recover the type of sub-grid details (in the elementary Riemann fans) which are discarded in the Godunov averaging step, like the ones displayed for instance on the right part of Fig. 1.1. It is possible to devise local projectors furnishing exact solutions at certain times; however, it doesn’t seem possible to recast them in this “Haar wavelet framework” because they result of an interpolation between \mathcal{P}_J and the random sampling of Glimm, see [13]. It is important to remember that, even if R_ϕ generates an approximation in V_{J+1} , usual MUSCL schemes still use the Riemann evolution operator \mathcal{E}_J , that is to say, the (new) discontinuities located in x_k aren’t resolved (see Fig. 2.2).

$$\forall n \in \mathbb{N}, \quad u^{\Delta x, \phi}(t^n, \cdot) = [\mathcal{P}_J \circ \mathcal{E}_J(\Delta t) \circ R_\phi]^n \mathcal{P}_J(u_0). \quad (1.11)$$

In Section 2, we prove the Main Theorem and in Section 3, advection-reaction equations and some issues raised by interaction of waves for systems [1, 12, 14] are studied. Finally, in Appendix A, some facts about “evolutionary errors” for discretizations based on the *Method of Lines* are recalled following mainly [9, 26, 29, 30, 39, 40].

REMARK 1.3. *Hereafter, ϕ stands for a slope-limiter. One may set up a flux-limiter instead, but showing an analogy with a wavelet formalism would result less easy. Generally, the term “flux-limiter” is used when it acts directly on fluxes, and “slope-limiter”, when it acts just on states. Both have the same mathematical form, and have the effect of limiting the solution’s gradient near shocks or local extrema.*

2. Proof of the Main Theorem

MUSCL-based numerical schemes extend the idea of using a linear piecewise approximation to each cell by using *slope limited* left and right extrapolated states². With the Godunov flux (1.3), they yield second-order resolution (SOR), Total-Variation Diminishing (TVD) time-marching processes after some approximations.

2.1. MUSCL reconstruction as an extrapolation process

Let’s recall how the “extrapolated states” are derived: for any indexes $k, n \in \mathbb{Z} \times \mathbb{N}$, the Godunov averaging furnished an approximate (formally first-order) value $u_k^n \simeq u(t^n, x_k)$, from which a piecewise-linear reconstruction is deduced in each cell,

$$v_k^n : C_k \rightarrow \mathbb{R}, \quad v_k^n(x) = u_k^n + (x - x_k)\sigma_k^n. \quad (2.1)$$

A first way to proceed is by analogy with Lax-Wendroff second-order scheme with $f(u) = u$: so, a convenient definition of the local slopes reads,

$$\sigma_k^n = \frac{u_{k+1}^n - u_k^n}{\Delta x} \phi(r_k^n), \quad r_k^n = \frac{u_k^n - u_{k-1}^n}{u_{k+1}^n - u_k^n}.$$

The slope-limiter ϕ must meet with several constraints in order to ensure both the TVD and second-order accuracy, see [35]. Recalling the construction of the Godunov scheme, we must face now the resolution of interfacial discontinuities separating linear

²The classical Donoho-Stark criterion suggests that, as its scaling function is discontinuous, performing scale-limited extrapolation in the Haar multi-resolution spaces may be unstable [10, 17]. Hereafter, extrapolation refers to *piecewise polynomial* extrapolation, and not to *scale-limited* one.

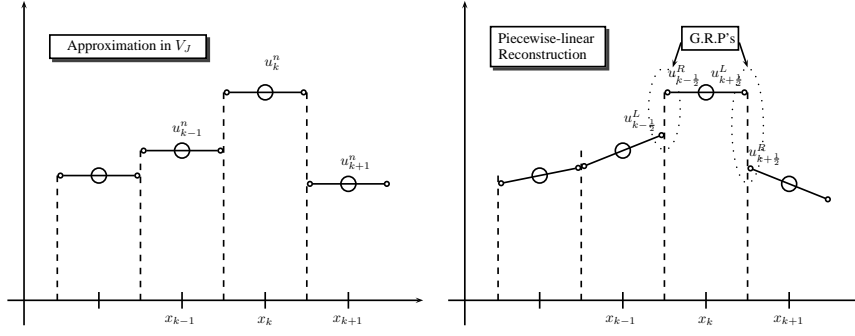


FIG. 2.1. Piecewise-linear reconstruction leading to Generalized Riemann Problems.

polynomials instead of constant states: this is usually called a *Generalized Riemann Problem* (see [36] pages 427–9 and Fig. 2.1). The issue is, quoting Osher “obtaining the exact solution to this nonlinear initial value problem with piecewise linear initial data is a nontrivial business” (see also [2, 3]): in particular, the simplification (1.2) is lost along with self-similarity property, except for the astute derivation presented in [16]. So in the vast majority of cases, the MUSCL algorithm limits itself to solving again usual (self-similar) *Riemann problems with extrapolated states* at each interface $x_{k\pm\frac{1}{2}}$ of the grid with $\Delta x = 2^{-J}$, hence $\mathcal{E}_J(\Delta t)$. For the advection equation, *i.e.* $f(u) = a \cdot u$, the global error generated by MINMOD reconstructions was analyzed in [28]: a convergence rate slightly greater than $\frac{1}{2}$ was obtained with weak solutions.

REMARK 2.1. Another way to motivate MUSCL piecewise-linear reconstructions is to work out the ODE system obtained by semi-discretization in space (the “Method of Lines”, evoked in [24]) in order to obtain a Local (space-) Truncation Error in Δx^2 for smooth exact solutions u : see our Appendix A and Verwer’s papers [29, 30, 39, 40].

2.2. The Haar wavelet fluctuation

As we explained in the former subsection, a more correct representation of the MUSCL algorithm is displayed on Fig. 2.2, where remain only a set of extrapolated states $u_{k-\frac{1}{2}}^{L/R}$ and the corresponding (usual, self-similar) Riemann problems:

$$\forall k \in \mathbb{Z}, \quad u_{k-\frac{1}{2}}^R = u_k^n - \phi(r_k^n) \frac{u_{k+1}^n - u_k^n}{2}, \quad u_{k+\frac{1}{2}}^L = u_k^n + \phi(r_k^n) \frac{u_{k+1}^n - u_k^n}{2} \quad (2.2)$$

Since the local reconstructions (2.1) are odd in the $x - x_k$ variable, it is now obvious that the states (2.2) rewrite by means of the Haar wavelet ψ . More precisely, given (1.8) as the Godunov approximation at time t^n in V_J , those states read:

$$u_{k-\frac{1}{2}}^R = u_k^n - \phi(r_k^n) \frac{u_{k+1}^n - u_k^n}{2} \psi\left(-\frac{1}{2}\right), \quad u_{k+\frac{1}{2}}^L = u_k^n + \phi(r_k^n) \frac{u_{k+1}^n - u_k^n}{2} \psi\left(\frac{1}{2}\right).$$

Since W_J is the orthogonal complement of V_J in V_{J+1} , the MUSCL extrapolated states furnish a piecewise-constant approximation in the finer scale-limited subspace,

$$2^{-\frac{J}{2}} \left(u_k^n \varphi_{J,k}(x) - \phi(r_k^n) \frac{u_{k+1}^n - u_k^n}{2} \psi_{J,k}(x) \right) \in V_{J+1} = V_J \oplus W_J, \quad (2.3)$$

because Haar wavelets satisfy the following relation:

$$2^{-\frac{J}{2}} \psi_{J,k}(x_{k\pm\frac{1}{2}}) = \psi \left(2^J (k \Delta x \pm \frac{\Delta x}{2}) - k \right) = \psi (2^J (k \cdot 2^{-J} \pm 2^{-J-1}) - k) = \psi(\pm \frac{1}{2}).$$

A consequence of the formulation (2.2) is that one has $\frac{1}{2}(u_{k-\frac{1}{2}}^R + u_{k+\frac{1}{2}}^L) = u_k^n$ in all the

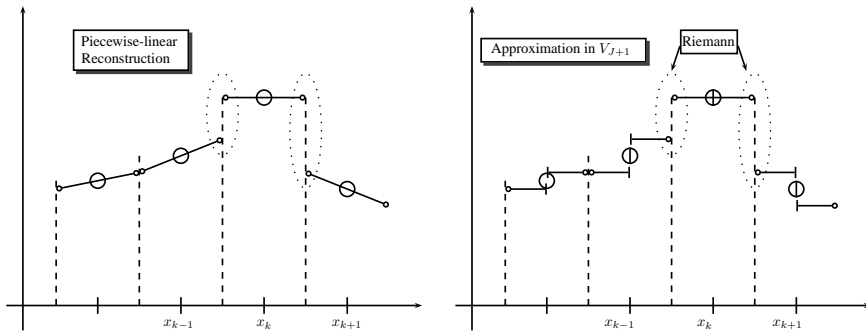


FIG. 2.2. MUSCL Piecewise-linear reconstruction and approximation in V_{J+1} .

cells C_k . However, such a symmetry should occur for a transonic rarefaction wave.

2.3. The back-projection (conservation) property

The property (1.10) is a direct consequence of the simple observation:

$$\forall k, k' \in \mathbb{Z}^2, \quad \int_{\mathbb{R}} \varphi_{J,k}(x) \psi_{J,k'}(x) dx = 0.$$

Indeed, thanks to the definition (1.6) and the expression (2.3), one sees that by linearity of the integral, the former orthogonality property implies

$$\forall k \in \mathbb{Z}, \quad \left\langle u_k^n \varphi_{J,k} - \phi(r_k^n) \frac{u_{k+1}^n - u_k^n}{2} \psi_{J,k}, \varphi_{J,k} \right\rangle = \langle u_k^n \varphi_{J,k}, \varphi_{J,k} \rangle = u_k^n.$$

This completes the proof of the Main Theorem 1.2. \square

3. Inhomogeneous equations, comments and outlook

Looking at Fig. 1.1, one sees that the reconstruction on Fig. 2.2 cannot yield an improvement in terms of elementary wave interactions. Indeed, in order to be compatible with the sub-grid details which are discarded in the averaging for the case of Fig. 1.1, the extrapolation process should address the right-half of the computational cells only. But such a reconstruction process wouldn't belong to W_J : instead, it would have a component in V_J and the "back-projection" (1.10) wouldn't hold (see [23]).

3.1. A new MUSCL-WB scheme for advection-reaction

A case where G.R.P.'s arising from piecewise-linear reconstructions like on Fig. 2.1 is when $f(u) = au$, with $a > 0$ taken for convenience: formula (1.2) modifies into,

$$\int_{t^n}^{t^{n+1}} a \cdot v_k^n \left(x_{k+\frac{1}{2}} - a(\tau - t^n) \right) d\tau = \underbrace{a \Delta t \left(u_k^n + \frac{\sigma_k^n \Delta x}{2} \right)}_{\text{usual MUSCL flux}} - \underbrace{a^2 \frac{\sigma_k^n \Delta t^2}{2}}_{\text{correction GRP}}, \quad (3.1)$$

so numerical fluxes do depend on time (see (5.3b) in [20]). Oppositely, when one simply substitutes G.R.P.'s with usual, self-similar Riemann problems, such a phenomenon doesn't show up, but the CFL number must be lowered to a value around $\frac{1}{2}$ (see [35]). Yet, according to Fig. 3.1, our reconstruction in V_{J+1} is different because

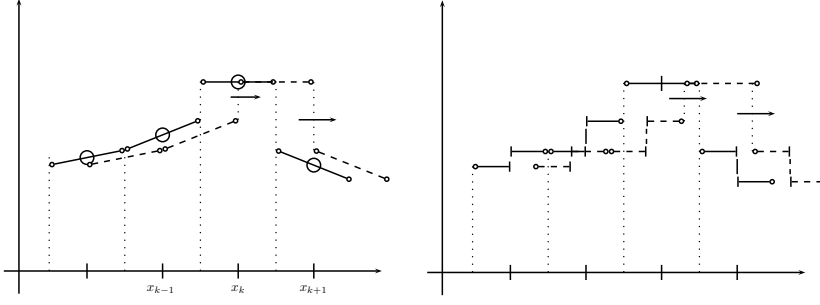


FIG. 3.1. Derivation of numerical fluxes: GRP's (left) and approximation in V_{J+1} (right).

it involves new discontinuities at the center of each cell, x_k . Hence it appears logical to replace the Riemann solver \mathcal{E}_J which handles only the discontinuities at interfaces $x_{k\pm\frac{1}{2}}$ by \mathcal{E}_{J+1} processing jumps at both interfaces and center. The resulting scheme,

$$\tilde{u}^{\Delta x, \phi}(t^n, \cdot) = [\mathcal{P}_J \circ \mathcal{E}_{J+1}(\Delta t) \circ R_\phi]^n \mathcal{P}_J(u_0), \quad \phi(r) = \frac{r + |r|}{1 + |r|}, \quad (3.2)$$

involves numerical fluxes still derived by modifying formula (1.2), see Fig. 3.2:

$$\begin{aligned} a \int_{t^n}^{t^{n+1}} u_k^n - 2^{\frac{j}{2}} \psi_{J,k}(x_{k+\frac{1}{2}} - a(\tau - t^n)) \frac{\sigma_k^n \Delta x}{2} d\tau &\stackrel{def}{=} a \Delta t \left(\frac{u_{k+1}^n + u_k^n}{2} - \tilde{Q}_{j+\frac{1}{2}}^n \frac{u_{k+1}^n - u_k^n}{2} \right) \\ &= a \Delta t \cdot u_k^n + \frac{\sigma_k^n \Delta x}{2} \left(\min(a \Delta t, \frac{\Delta x}{2}) - \max(0, a \Delta t - \frac{\Delta x}{2}) \right). \end{aligned}$$

A slope-limiter ϕ is indispensable because its numerical viscosity is only,

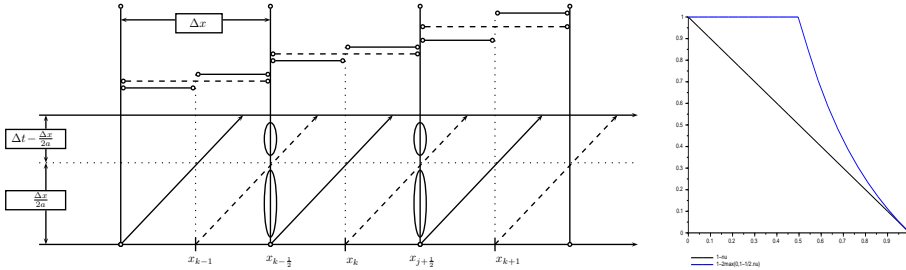


FIG. 3.2. Illustration of scheme (3.2) (left), and deviation w.r.t. exact GRP fluxes (right).

$$\begin{aligned} \tilde{Q}_{j+\frac{1}{2}}^n &= 1 - \phi(r_j^n) \left(\min(1, \frac{1}{2\nu}) - \max(0, 1 - \frac{1}{2\nu}) \right), \quad \nu = \frac{a \Delta t}{\Delta x}, \quad (3.3) \\ &= 1 + \phi(r_j^n) \left(2 \max(0, 1 - \frac{1}{2\nu}) - 1 \right), \quad \left(\text{as } 1 - \min(1, \frac{1}{2\nu}) = \max(0, 1 - \frac{1}{2\nu}) \right). \end{aligned}$$

One may compare it to $Q_{j+\frac{1}{2}}^n$, the one associated to the numerical flux (3.1):

$$Q_{j+\frac{1}{2}}^n = 1 - \phi(r_j^n)(1 - \nu), \quad \tilde{Q}_{j+\frac{1}{2}}^n = 1 - \phi(r_j^n) \left(1 - 2 \max(0, 1 - \frac{1}{2\nu}) \right),$$

see Fig. 3.2. All in all, this yields the following (and seemingly new) discretization,

$$u_k^{n+1} = u_k^n - \frac{a\Delta t}{\Delta x} \left[(u_k^n - u_{k-1}^n) + \left(2\min(1, \frac{1}{2\nu}) - 1 \right) \frac{\Delta x(\sigma_k^n - \sigma_{k-1}^n)}{2} \right], \quad (3.4)$$

which rewrites simply, $u_k^{n+1} = u_k^n - \frac{a\Delta t}{\Delta x} (\tilde{u}_{k+\frac{1}{2}}^n - \tilde{u}_{k-\frac{1}{2}}^n)$, after having defined,

$$\tilde{u}_{k+\frac{1}{2}}^n = (1 - \alpha_k^n)u_k^n + \alpha_k^n u_{k+1}^n, \quad \alpha_k^n = \left(\min(1, \frac{1}{2\nu}) - \frac{1}{2} \right) \phi(r_k^n) \in [0, 1].$$

Below we display numerical results for $a=1$, $u_0(x) = \sin^3(2\pi x)$, and 2^6 points in

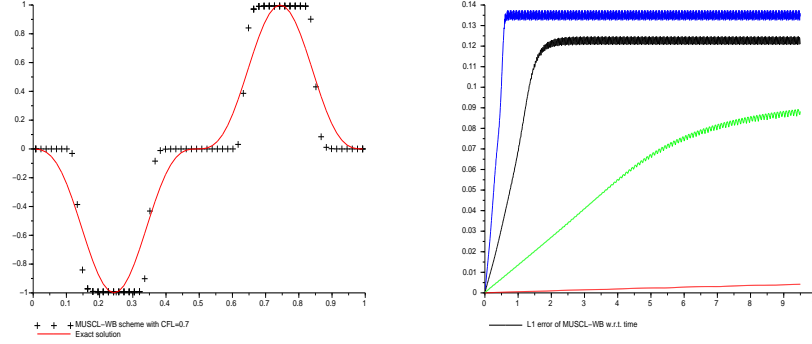


FIG. 3.3. Advection with $CFL=0.7$ (left) and L^1 -errors (right) for $CFL=0.55, 0.7, 0.85, 0.99$.

$x \in (0, 1)$ with periodic boundary conditions. To prevent the development of un-aesthetic staircases in the numerical solution, we sampled initial data on the grid $\{\frac{\Delta x}{2}, \frac{3\Delta x}{2}, \dots, 1 - \frac{\Delta x}{2}\}$: results are displayed on the left of Fig. 3.3 with $\Delta t = 0.7\Delta x$. On its right, one can see the dependence of the time-growth of the L^1 -error with respect to the Courant number; even if the error gets big when it goes lower than 0.8, at least it stops growing after a certain time. Since our MUSCL scheme isn't restricted by low CFL numbers, they are less vulnerable to numerical diffusion's bad effects.

REMARK 3.1. For $\nu=1$, the scheme (3.4) yields $u_k^{n+1} = u_{k-1}^n$, so it is exact. If $\nu \leq \frac{1}{2}$, it reduces to the usual second-order MUSCL scheme. Oppositely, for $\frac{1}{2} < \nu < 1$,

$$r_k^n = \frac{u_k^n - u_{k-1}^n}{u_{k+1}^n - u_k^n} = 1 - \frac{u_{k+1}^n - 2u_k^n + u_{k-1}^n}{u_{k+1}^n - u_k^n} = 1 + O(\Delta x), \quad \text{for } u(t^n, \cdot) \text{ smooth,}$$

so $\phi(r_k^n) = \phi(1) + O(\Delta x)$ because ϕ is Lipschitz-continuous and $\phi(1) = 1$ [35]. Yet,

$$\begin{aligned} \tilde{u}_{k+\frac{1}{2}}^n &= u_k^n + \frac{1}{2} \left(1 - \left(2 - \frac{1}{\nu} \right) \right) (u_{k+1}^n - u_k^n) (1 + O(\Delta x)) \\ &= \frac{1}{2} (u_{k+1}^n + u_k^n) + O(\Delta x^2) - \left(1 - \frac{1}{\nu} \right) (u_{k+1}^n - u_k^n) (1 + O(\Delta x)) \\ &= \frac{1}{2} (u_{k+1}^n + u_k^n) + \left(\frac{1}{\nu} - 1 \right) O(\Delta x) + O(\Delta x^2). \end{aligned}$$

According to (A.5), the L.T.E. is $O(\Delta x)$ for $\nu > \frac{1}{2}$ despite the weak viscosity (3.3). Next, when considering an inhomogeneous equation of the type $\partial_t u + \partial_x u = k(x)u$, the scheme (3.2) can match the WB framework [18] where one solves a ‘‘lifted equation’’,

$$\partial_t u + \partial_x u - u \partial_x a = 0, \quad \partial_t a = 0, \quad (\text{because } \partial_x a(x) = k(x)),$$

which induces a solver $\tilde{\mathcal{E}}$ now including a “standing wave” locally rendering the source,

$$\forall n \in \mathbb{N}, \quad \tilde{u}^{\Delta x, \phi, WB}(t^n, \cdot) = \left[\mathcal{P}_J \circ \tilde{\mathcal{E}}_{J+1}(\Delta t) \circ R_\phi \right]^n \mathcal{P}_J(u_0).$$

In the same manner as for (3.4), we get a simple expression of the resulting scheme,

$$u_k^{n+1} = u_k^n - \frac{\Delta t}{\Delta x} \left(\tilde{u}_{k+\frac{1}{2}}^n - \tilde{u}_{k-\frac{1}{2}}^n \cdot \exp(a(x_k) - a(x_{k-1})) \right).$$

The exact solution reads $u(t, x) = u_0(x-t) \exp(a(x-t) - a(x))$: on the left of Fig. 3.4,

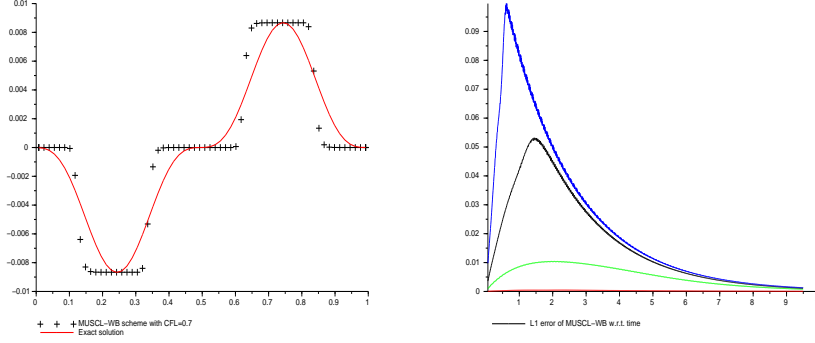


FIG. 3.4. Exponential decay, $CFL=0.7$ (left); L^1 -errors (right) for $CFL=0.55, 0.7, 0.85, 0.99$.

our WB scheme is set up for $k(x) \equiv -\frac{1}{2}$ and compared to it with 2^6 grid points and $CFL=0.7$ at $t=9.5$. The function $a(x)$ is discretized according to the grid correspond-

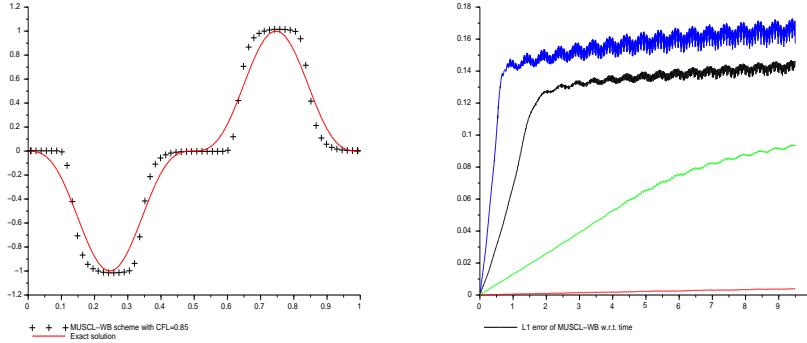


FIG. 3.5. Oscillating $k(x)$ with $CFL=0.85$ (left); L^1 -errors (right) for $CFL=0.55, 0.7, 0.85, 0.99$.

ing to V_J , so it jumps only at interfaces $x_{k \pm \frac{1}{2}}$: this means that the discontinuities in x_k are resolved with the homogeneous Riemann solver. A more accurate scheme would be produced if $a(x)$ is sampled on the fine scale of V_{J+1} , too. In this case, even the discontinuities in x_k 's are resolved with the WB Riemann solver. Such a scheme would be well suited for source terms containing an oscillating coefficient $k(x)$: for instance, the rather delicate case where $k(x) = \frac{\cos(4\pi x)}{2}$ is presented on Fig. 3.5. Such a MUSCL-WB scheme may be useful for kinetic models involving slow particles [18].

3.2. Wave-interactions and Engquist-Sjogreen counter-example

Now we switch to the more involved case of nonlinear systems of conservation laws: left apart the Temple class, shock curves aren't straight lines in the Hugoniot space. A first problem materializes because intermediate points resulting from the numerical "viscous smearing" of jumps generally don't belong to these curves. This creates spurious (small) waves of other characteristic families in the numerical solution: see Fig. 3.6 and [1, 22, 31]. Clearly MUSCL reconstructions can't improve noticeably this situation which occurs mainly for large discontinuities, though. About

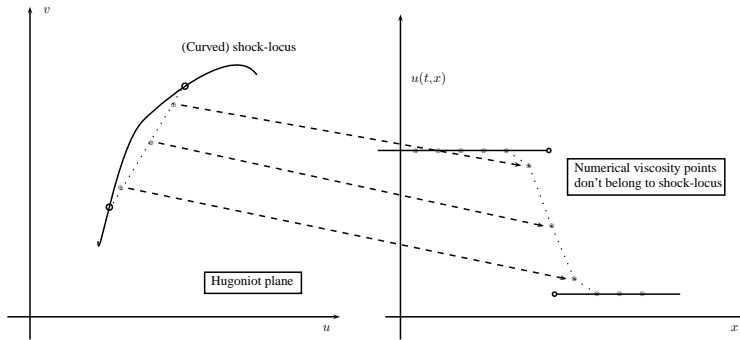


FIG. 3.6. Numerical viscosity and its effects for systems with curved shock loci

numerical wave interactions, P. L. Roe writes: *It is natural still to feel some anxiety about replacing the discrete jump conditions (Rankine-Hugoniot relations), that hold across an infinitesimally thin shock, with a 'smeared-out' statement of conservation. It seems likely that such a strategy will lead to some sort of unavoidable error. Currently, rather delicate computations of the interaction between a strong shock and weak acoustic waves are not successful unless the shock is either represented as an explicit discontinuity or else the grid spacing is greatly reduced in its vicinity ([21], page 15).* SOR can fail when Glimm's interaction potential is positive like in a p -system, an

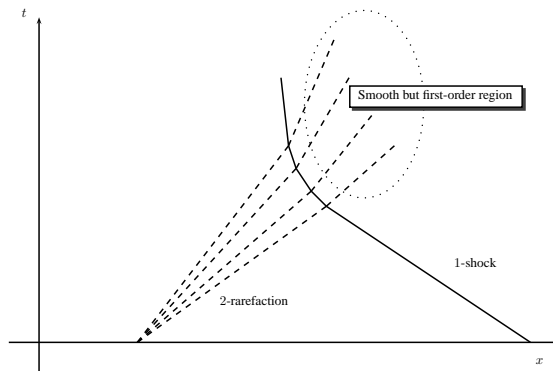


FIG. 3.7. Nonlinear 2×2 interaction and loss of SOR in a smooth region

interaction between a 1-shock and a 2-rarefaction. Along with numerical smearing of the profile, the local truncation error can reduce from second-order to first-order in

the vicinity of the shock wave. When interaction occurs (see Fig. 3.7), all the variables are coupled, and the formal accuracy of the outgoing rarefaction wave may not be second-order: this is the essence of what appears in [12], and later in [11, 32, 33].

3.3. Multi-dimensional issues

Most of existing 2D MUSCL schemes don't completely fit in our Haar wavelet formalism as the scaling function and mother wavelets read $\varphi^{2D}(x, y) = \varphi(x)\varphi(y)$ and

$$\psi^{2D,vert}(x, y) = \varphi(x)\psi(y), \quad \psi^{2D,hor}(x, y) = \psi(x)\varphi(y), \quad \psi^{2D,diag}(x, y) = \psi(x)\psi(y).$$

Hence it perceives a diagonal direction. The issues of numerical wave interactions still exist in 2D, see [7, 32, 38]. Multi-resolution schemes were studied in [8] (also [19, 4]).

Appendix A. Formal analysis of the evolutionary error.

Hereafter we follow the canvas of Cullen and Morton [9] in order to shed some light onto the various mechanisms of error creation/propagation (see also [5, 15, 40]).

A.1. Semi-discretization in space (Method of Lines)

Let a Cauchy problem for a given partial differential operator \mathcal{L} be,

$$\partial_t u = \mathcal{L}u, \quad u(t=0, \cdot) = u_0. \quad (\text{A.1})$$

For $\Delta x = 2^{-J}$ fixed and the corresponding gridding of the real line, a finite-differences approximation of \mathcal{L} acting on $\Delta x \cdot \mathbb{Z}$ is denoted by \mathcal{L}_J , so (A.1) reduces to an (infinite) differential system (*Method of Lines*, Ch. 17 in [24]), with $\tilde{u}(t, \cdot) \in \ell^\infty(\mathbb{Z})$, say:

$$\frac{d}{dt} \tilde{u} = \mathcal{L}_J \tilde{u}, \quad \tilde{u}(t=0, \cdot) = \mathcal{P}_J u_0, \quad (\text{A.2})$$

for which one can legitimately wonder about the global error $u - \tilde{u}$ at each time $t > 0$.

- one “triangulates” $u(t, \cdot) - \tilde{u}(t, \cdot)$ by inserting $\mathcal{P}_J u(t, \cdot)$,

$$u - \tilde{u} = (Id - \mathcal{P}_J)u + (\mathcal{P}_J u - \tilde{u}) := a_J + e_J,$$

where a_J is purely an approximation error, which belongs to the wavelet subspace $\cup_{j \geq J} W_j$. On the contrary, e_J stands for an evolutionary error, which may accumulate in time, and satisfies a differential equation,

$$\frac{d}{dt} e_J = \frac{d}{dt} \mathcal{P}_J u - \frac{d}{dt} \tilde{u} = \mathcal{P}_J \mathcal{L}u - \mathcal{L}_J \tilde{u}. \quad (\text{A.3})$$

- Triangulating again, one gets $\frac{de_J}{dt} = (\mathcal{P}_J \mathcal{L}u - \mathcal{L}_J \mathcal{P}_J u) + (\mathcal{L}_J \mathcal{P}_J u - \mathcal{L}_J \tilde{u})$, so

$$\frac{d}{dt} e_J + (\mathcal{L}_J \tilde{u} - \mathcal{L}_J \mathcal{P}_J u) = (\mathcal{P}_J \mathcal{L}u - \mathcal{L}_J \mathcal{P}_J u) := L.T.E.,$$

and by substituting \tilde{u} by $\mathcal{P}_J u - e_J$, we get finally:

$$\frac{d}{dt} e_J + [\mathcal{L}_J(\mathcal{P}_J u - e_J) - \mathcal{L}_J \mathcal{P}_J u] = L.T.E., \quad (\text{Local Truncation Error}). \quad (\text{A.4})$$

Hence, the L.T.E. is just a source term inside the differential equation (A.4) governing the scheme's evolutionary error; this was noted in [26, 29, 40], too.

In case both (A.1) and its (consistent) discrete approximation \mathcal{L}_J , are dissipative (“contractive” [39, 30], “strongly stable” in a terminology of [26]) in some norm, this source term is responsible for most of the error e_J ; if, on the contrary, (A.1) happens to be accretive, for instance if $\|u(t) - v(t)\| \leq K\|u_0 - v_0\|$ with $K > 1$ like in Bressan-Glimm’s theory of strictly hyperbolic systems of conservation laws [6], then both \mathcal{L}_J and the L.T.E. can contribute to the increase of the evolutionary error, see again [40].

REMARK A.1. *If the approximation \mathcal{L}_J is linear, then (A.4) simplifies into,*

$$\forall t > 0, \quad \frac{d}{dt} e_J(t) = \mathcal{L}_J e_J(t) + \tau_u(t),$$

where $\tau_u(t)$ stands for the L.T.E. related to (x -derivatives of) the exact solution $u(t, \cdot)$ to (A.1) at time t . Duhamel’s principle yields an expression of the evolutionary error,

$$e_J(t) = \exp(t \cdot \mathcal{L}_J) \left(e_J(t=0) + \int_0^t \exp(-s \cdot \mathcal{L}_J) \tau_u(s) ds \right).$$

Quantities like $\exp(t \cdot \mathcal{L}_J)$ are usually estimated by “logarithmic norms”, see e.g. [30].

A.2. Local Truncation Error (LTE) and second-order accuracy

Second-order accuracy in space for 1D scalar conservation laws (or linear advection equations) was studied in [27] (see also [21, 36]). These equations are dissipative in L^1 , so the former analysis yielding (A.4) indicates that the local truncation error is probably the main source of evolutionary error. For $\mathcal{L}u = -\partial_x f(u)$, it reads:

$$\forall k \in \mathbb{Z}, \quad \mathcal{P}_J \mathcal{L}u(t, x_k) = -\frac{1}{\Delta x} \int_{x_{k-\frac{1}{2}}}^{x_{k+\frac{1}{2}}} \partial_x f(u) dx = -\frac{f(u(t, x_{k+\frac{1}{2}})) - f(u(t, x_{k-\frac{1}{2}}))}{\Delta x},$$

by exact integration of the conservation law (1.1). Now, since high-order accuracy is only concerned with smooth exact solutions u , one approximates this expression with a second-order mid-point rule by taking advantage of $x_{k+\frac{1}{2}} = \frac{x_{k+1} + x_k}{2}$,

$$\mathcal{P}_J \mathcal{L}u(t, x_k) = \frac{f\left(\frac{u(t, x_{k+1}) + u(t, x_k)}{2}\right) - f\left(\frac{u(t, x_k) + u(t, x_{k-1})}{2}\right) + O(\Delta x^2)}{\Delta x},$$

and so, the L.T.E. is the difference between this approximation and the numerical scheme \mathcal{L}_J applied to the piecewise constant projection of the exact solution, $\mathcal{P}_J u$. Since \mathcal{L}_J needs to be conservative and consistent with \mathcal{L} , we assume it is given by a (smooth) numerical flux which reads, in standard notation,

$$\tilde{F}_{k+\frac{1}{2}} = F(u_{k+\frac{1}{2}}^L, u_{k+\frac{1}{2}}^R), \quad \mathcal{L}_J \mathcal{P}_J u(t, x_k) = \frac{\tilde{F}_{k+\frac{1}{2}}(t) - \tilde{F}_{k-\frac{1}{2}}(t)}{\Delta x},$$

where $u_{k+\frac{1}{2}}^{L/R}$ are obtained from the set of cell-centered values $\mathcal{P}_J u$ by means of a reconstruction like (2.2) and F is, for instance, the exact Godunov flux (1.3). Hence,

$$L.T.E. = \frac{[f\left(\frac{u(t, x_{k+1}) + u(t, x_k)}{2}\right) - \tilde{F}_{k+\frac{1}{2}}] - [f\left(\frac{u(t, x_k) + u(t, x_{k-1})}{2}\right) - \tilde{F}_{k-\frac{1}{2}}]}{\Delta x}.$$

As the CFL condition imposes $\Delta t = O(\Delta x)$, second-order accuracy asks for,

$$\left| f\left(\frac{u(t, x_{k+1}) + u(t, x_k)}{2}\right) - \tilde{F}_{k+\frac{1}{2}}(t) \right| = O(\Delta x^2),$$

which, by the smoothness of the flux functions, reduces simply to,

$$\forall t, k \in \mathbb{R}^+ \times \mathbb{Z}, \quad \left| u_{k+\frac{1}{2}}^{L/R}(t) - \frac{u(t, x_{k+1}) + u(t, x_k)}{2} \right| = O(\Delta x^2). \quad (\text{A.5})$$

And this meets with the definition used by Osher (see Lemma 2.1, page 953 in [27]) and Sjogreen (see Theorem 3.9 in [34], page 47). A slightly different derivation of a second-order scheme for smooth solutions is given in [5] (page 53), essentially by keeping the term $\frac{d}{dt} \mathcal{P}_J u$ in (A.3) inside the expression of the L.T.E as follows:

$$\begin{aligned} \frac{d}{dt} \mathcal{P}_J u(t, \cdot) &= \lim_{\Delta t \rightarrow 0} \left(\frac{\mathcal{P}_J u(t + \Delta t, \cdot) - \mathcal{P}_J u(t, \cdot)}{\Delta t} \right) \\ &= - \frac{F(u(t, \cdot + \Delta x), u(t, \cdot)) - F(u(t, \cdot), u(t, \cdot - \Delta x))}{\Delta x}, \end{aligned}$$

where F is the exact flux defined in (1.3). The L.T.E. is now defined like,

$$\forall k \in \mathbb{Z}, \quad \frac{d}{dt} \mathcal{P}_J u(t, x_k) - \mathcal{L}_J \mathcal{P}_J u(t, x_k) = - \frac{\mathcal{F}_{k+\frac{1}{2}}(t) - \mathcal{F}_{k-\frac{1}{2}}(t)}{\Delta x},$$

where $\mathcal{F}_{k+\frac{1}{2}}(t) = F(u(t, x_k + \Delta x), u(t, x_k)) - F_{k+\frac{1}{2}}(t)$. The scheme induced by the numerical flux $F_{k+\frac{1}{2}}$ is called second-order in space as soon as, for any smooth exact solution $u(t, \cdot)$, $\mathcal{F}_{k+\frac{1}{2}}$ is a quadratic quantity (possibly depending on $|\partial_{xx} u(t, \cdot)|$),

$$\forall t \geq 0, \quad |\mathcal{F}_{k+\frac{1}{2}}(t)| = O(\Delta x^2). \quad (\text{A.6})$$

Clearly, both criteria pick up variants of the (unstable) ‘‘centered scheme’’ which is second-order, but unstable because it lets the total variation increase strongly: despite its L.T.E. is quite small, its evolutionary error quickly grows with the left-hand side of (A.3). MUSCL reconstructions, involving a slope limiter, allow to keep both L.T.E. and other terms in the O.D.E. (A.3) governing e_J rather small (in smooth regions).

REFERENCES

- [1] M. Arora, P.L. Roe, *On postshock oscillations due to capturing schemes in unsteady flows*, J. Comput. Phys. **130** (1997) 25–40.
- [2] M. Ben-Artzi, J. Falcovitz, **Generalized Riemann problems in computational fluid dynamics**, Cambridge monographs on applied and computational mathematics **11** (2003).
- [3] C. Berthon, C. Sarazin, R. Turpault, *Space-time Generalized Riemann Problem Solvers of Order k for Linear Advection with Unrestricted Time Step*, J. Sci. Comput. **55** (2013) 268–308
- [4] B.L. Bihari, A. Harten, *Multiresolution schemes for the numerical solution of 2D conservation laws*, SIAM J. Scient. Comput. **18** (1997) 315–354.
- [5] F. Bouchut, **Nonlinear stability of finite volume methods for hyperbolic conservation laws, and well-balanced schemes for sources**, Frontiers in Mathematics series, Birkhäuser, 2004, ISBN 3-7643-6665-6.
- [6] A. Bressan. **Hyperbolic Systems of Conservation Laws – The one-dimensional Cauchy problem**, Oxford Lecture Series in Mathematics and its Applications **20** (2000).
- [7] Mark H. Carpenter, Jay H. Casper, *Accuracy of Shock Capturing in Two Spatial Dimensions*, AIAA Journal **37** (1999) 1072–1079.
- [8] Cohen, Albert; Kaber, Sidi Mahmoud; Müller, Siegfried; Postel, Marie. *Fully adaptive multiresolution finite volume schemes for conservation laws*, Math. Comp. **72** (2003) 183–225
- [9] M.J.P. Cullen, K.W. Morton, *Analysis of Evolutionary Error in Finite Element and Other Methods*, J. Comput. Phys. **34** (1980) 245–267.
- [10] D.L. Donoho, P.B. Stark, *Uncertainty principles and signal recovery*, SIAM J. Appl. Math. **49** (1989) 906–931.

- [11] G. Efrainsson, G. Kreiss, *A remark on numerical errors downstream of slightly viscous shocks*, SIAM J. Numer. Anal. **36** (1999) 853–863.
- [12] B. Engquist, B. Sjögreen, *The Convergence Rate of Finite Difference Schemes in the Presence of Shocks*, SIAM J. Numer. Anal. **35** (1998) 2464–2485.
- [13] H. Gilquin, *Une famille de schémas numériques T.V.D. pour les lois de conservation hyperboliques*. RAIRO - Model. Math. & Anal. Num., **20** (1986), 429–460
- [14] J. Glimm, *The interaction of nonlinear hyperbolic waves*, Comm. Pure Appl. Math. **41** (1988) 569–590.
- [15] Sonia M. Gomes, *Unified overview of wavelet-based methods for differential equations*, Proc. SPIE 3078, Wavelet Applications IV, **730** (April 3, 1997); doi:10.1117/12.271758
- [16] J.B. Goodman, R.J. LeVeque, *A geometric approach to high resolution TVD schemes*, SIAM J. Numer. Anal. **25** (1988) 268–284.
- [17] L. Gosse, *A Donoho–Stark criterion for stable signal recovery in discrete wavelet subspaces*, J. Comput. Appl. Math. **235** (2011) 5024–5039.
- [18] L. Gosse, **Computing Qualitatively Correct Approximations of Balance Laws**, Springer (2013) ISBN 978-88-470-2891-3
- [19] A. Harten, *Discrete multi-resolution analysis and generalized wavelets*, Appl. Numer. Math. **12** (1993) 153–192
- [20] A. Harten, S. Osher, *Uniformly high-order accurate non-oscillatory schemes*, SIAM J. Numer. Anal. **24** (1987), no. 2, 279–309 (paper available in [21])
- [21] Hussaini, M.Yousuff (ed.); van Leer, Bram (ed.); Van Rosendale, John (ed.), **Upwind and high-resolution schemes**. (English) Berlin: Springer (1997).
- [22] Shi Jin, Jian-Guo Liu, *The effects of numerical viscosities. I. Slowly moving shocks*, J. Comput. Phys. **126** (1996) 373–389.
- [23] F. Lagoutiere, *Non-dissipative entropy satisfying discontinuous reconstruction schemes for hyperbolic conservation laws*, Preprint www.math.u-psud.fr/~lagoutie/publi.html
- [24] R.J. LeVeque, **Numerical methods for conservation laws**, ETH Zurich, Birkhauser 1992.
- [25] S. Mallat, **A wavelet tour of signal processing**, Academic Press, 1998.
- [26] Keith W. Morton, *On the analysis of finite volume methods for evolutionary problems*, SIAM J. Numer. Anal. **35** (1998), 2195–2222.
- [27] S. Osher, *Convergence of generalized MUSCL schemes*, SIAM J. Numer. Anal. **22** (1985), 947–961. (paper available in [21])
- [28] B. Popov, O. Trifonov, *Order of convergence of second order schemes based on the MINMOD limiter*, Math. of Comp. **75** (2006) 1735–1753.
- [29] J.M. Sanz-Serna, J.G. Verwer, *Convergence Analysis of One-Step Schemes in the Method of Lines*, Applied Math. Comput. **31** (1989) 183–196.
- [30] J.M. Sanz-Serna, J.G. Verwer, *Stability and convergence at the PDE/stiff ODE interface*, Applied Numer. Math. **5** (1989) 117–132.
- [31] M. Siklosi, B. Batzorig, G. Kreiss, *An investigation of the internal structure of shock profiles for shock capturing schemes*, J. Comput. Appl. Math. **201** (2007) 8–29.
- [32] M. Siklosi, G. Efrainsson, *Analysis of first order errors in shock calculations in two space dimensions*, SIAM J. Numer. Anal. **43** (2005) 672–685.
- [33] M. Siklosi, G. Kreiss, *Elimination of first order errors in time dependent shock calculations*, SIAM J. Numer. Anal. **41** (2003) 2131–2148.
- [34] B. Sjögreen, *Lecture notes*, at www.math.fsu.edu/~sussman/Bjorn.Sjogreen.Notes.pdf
- [35] P.K. Sweby, *High Resolution Schemes Using Flux Limiters for Hyperbolic Conservation Laws*, SIAM J. Numer. Anal. **21** (1984) 995–1011.
- [36] E.F. Toro, **Riemann Solvers and Numerical Methods for Fluid Dynamics: A Practical Introduction**, Third Edition, Springer (2009).
- [37] B. Van Leer, *Towards the Ultimate Conservative Difference Scheme, V. A Second Order Sequel to Godunov’s Method*, J. Comp. Phys. **32** (1979) 101–136.
- [38] B. Van Leer, *Upwind and High-Resolution Methods for Compressible Flow: From Donor Cell to Residual-Distribution Schemes*, Commun. Comput. Phys. **1** (2006) 192–206.
- [39] J.G. Verwer, *Contractivity in locally one-dimensional splitting methods*, Numer. Math. **44** (1984) 247–259.
- [40] J.G. Verwer, J.M. Sanz-Serna, *Convergence of Method of Lines Approximations to Partial Differential Equations*, Computing **33** (1984) 297–313.