

 Open access • Journal Article • DOI:10.1016/J.ANIHPC.2015.01.001

## **Stringent error estimates for one-dimensional, space-dependent $2 \times 2$ relaxation systems** — [Source link](#)

Debora Amadori, Laurent Gosse

**Institutions:** University of L'Aquila, IAC

**Published on:** 01 May 2016 - Annales De L Institut Henri Poincare-analyse Non Lineaire (Elsevier Masson)

**Topics:** Conservation law

Related papers:

- [Computing Qualitatively Correct Approximations of Balance Laws](#)
- [On the zero relaxation limit for a system modeling the motions of a viscoelastic solid](#)
- [Asymptotic behavior of hyperbolic boundary value problems with relaxation term](#)
- [Global Solutions and Zero Relaxation Limit for a Traffic Flow Model](#)
- [Bv solutions and relaxation limit for a model in viscoelasticity](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/stringent-error-estimates-for-one-dimensional-space-ifvpgko9u3>



**HAL**  
open science

## Stringent error estimates for one-dimensional, space-dependent $2 \times 2$ relaxation systems

Debora Amadori, Laurent Gosse

► **To cite this version:**

Debora Amadori, Laurent Gosse. Stringent error estimates for one-dimensional, space-dependent  $2 \times 2$  relaxation systems. *Annales de l'Institut Henri Poincaré (C) Non Linear Analysis*, Elsevier, 2015, pp.23. 10.1016/j.anihpc.2015.01.001 . hal-01057737

**HAL Id: hal-01057737**

**<https://hal.archives-ouvertes.fr/hal-01057737>**

Submitted on 25 Aug 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.

# Stringent error estimates for one-dimensional, space-dependent $2 \times 2$ relaxation systems

Debora Amadori

*DISIM, Università degli Studi dell'Aquila, L'Aquila, Italy*

Laurent Gosse\*

*IAC-CNR "Mauro Picone" (sezione di Roma)*

*Via dei Taurini, 19 - 00185 Rome, Italy*

---

## Abstract

Sharp and local  $L^1$  *a-posteriori* error estimates are established for so-called "well-balanced"  $BV$  (hence possibly discontinuous) numerical approximations of  $2 \times 2$  space-dependent Jin-Xin relaxation systems under sub-characteristic condition. According to the strength of the relaxation process, one can distinguish between two complementary regimes: 1/ a weak relaxation, where local  $L^1$  errors are shown to be of first order in  $\Delta x$  and uniform in time, 2/ a strong one, where numerical solutions are kept close to entropy solutions of the reduced scalar conservation law, and for which Kuznetsov's theory indicates a behavior of the  $L^1$  error in  $t \cdot \sqrt{\Delta x}$ . The uniformly first-order accuracy in weak relaxation regime is obtained by carefully studying interaction patterns and building up a seemingly original variant of Bressan-Liu-Yang's functional, able to handle  $BV$  solutions of arbitrary size for these particular inhomogeneous systems. The complementary estimate in strong relaxation regime is proven by means of a suitable extension of methods based on entropy dissipation for space-dependent problems.

*Key words:* Bressan-Liu-Yang functional; Entropy dissipation; Kuznetsov's method;  $L^1$  error estimate; space-dependent relaxation model; Well-balanced scheme.

*1991 MSC:* 35L60, 65M06

---

## Contents

1	Introduction	2
1.1	Space-dependent $2 \times 2$ Jin-Xin relaxation model	2

---

\* Corresponding Author

*Email addresses:* amadori@univaq.it (Debora Amadori), l.gosse@ba.iac.cnr.it (Laurent Gosse).

*Preprint submitted to Elsevier Science*

1.2	Main result and plan of the paper	3
2	Construction of the Well-Balanced approximation	5
2.1	First considerations on the $3 \times 3$ Riemann problem	5
2.2	Shape of positively invariant domains for $3 \times 3$ Riemann problems	7
2.3	Total variation estimate of the WB approximation	9
3	An $L^1$ error estimate through a Lyapunov functional	10
3.1	A lemma based on sub-characteristic condition	11
3.2	Accurate interaction estimates for WB approximations	14
3.3	Lyapunov functional and error estimate for weak relaxation	19
4	Complementary $L^1$ error estimate through entropy dissipation	25
4.1	Quasi-monotonicity and entropy inequalities	25
4.2	Derivation of the complementary $L^1$ error estimate	27
A	An elementary example	29
	References	30

## 1 Introduction

### 1.1 Space-dependent $2 \times 2$ Jin-Xin relaxation model

We consider the simplest 1D kinetic model involving a space-dependent Knudsen number, which can be written with the following system:

$$\begin{cases} \partial_t \rho + \partial_x J = 0 \\ \partial_t J + \partial_x \rho = 2k(x)g(\rho, J) \end{cases} \quad (1)$$

under the assumption that for some  $c > 0$ ,

$$k \in L^1 \cap BV(\mathbb{R}), \quad k(x) \geq 0 \quad (2)$$

$$\partial_J g \leq -c < 0, \quad |\partial_\rho g| < |\partial_J g|. \quad (3)$$

Moreover we assume that there exists a  $C^1$  bounded map  $A(\rho)$  such that

$$g(\rho, A(\rho)) = 0 \quad \text{for all } \rho. \quad (4)$$

The curve  $J = A(\rho)$  will be called *equilibrium curve*. By taking the derivative of (4) and using (3), it follows that the so-called *sub-characteristic condition* holds:

$$|A'(\rho)| < 1. \quad (5)$$

A typical choice for  $g$  is given by the *relaxation term*:

$$g(\rho, J) = A(\rho) - J. \quad (6)$$

**Remark 1** *The system (1) with the assumptions (2)-(3)-(4) perfectly matches the two-scale relaxation framework studied in [12]. A variant of the expression (6) would be for instance:*

$$g(x, \rho, J) = A(x, \rho) - J, \quad x \mapsto A(x, \cdot) \in C^1(\mathbb{R}), \quad \sup_x |\partial_\rho A(x, \rho)| < 1.$$

*Such a model, which appears for instance in [10,31], wouldn't strongly modify our interaction estimates and consequently our error estimates. Another field of application would be an elementary semi-conductor*

model, for which the convective part would correspond to a lattice temperature  $\theta_0$ ,

$$\partial_t J + \partial_x(\theta_0 \rho) = \frac{1}{\tau(x)} (\tau(x) E(x) \rho - J), \quad \tau(x) |E(x)| < \sqrt{\theta_0}, \quad (7)$$

with  $E(x)$  a small static electric field, and  $\tau(x)$  standing for a space-dependent relaxation time depending on the local doping concentration.

In terms of “microscopic diagonal” variables  $f^\pm$ , defined by

$$\rho = f^+ + f^-, \quad J = f^+ - f^-$$

the system (1) rewrites as a discrete-velocity kinetic model:

$$\begin{cases} \partial_t(f^-) - \partial_x(f^-) = -k(x) G(f^-, f^+) \\ \partial_t(f^+) + \partial_x(f^+) = k(x) G(f^-, f^+) \end{cases} \quad (8)$$

where  $G(f^-, f^+) := g(f^+ + f^-, f^+ - f^-)$ . Initial data for (8) are chosen such that

$$f^\pm(t=0, \cdot) = f_0^\pm \in L^1 \cap BV(\mathbb{R}). \quad (9)$$

We close this section by indicating that our semi-linear, space-dependent model (1) belongs to the class of relaxation systems [21], which was intensively studied both analytically and numerically more a decade ago, mostly for constant coefficients  $\partial_x k \equiv 0$ , though: see [3,20,22,23,27,29,33,35], also [4,5,6] and the survey by Natalini [30]. Rigorous error estimates for inhomogeneous hyperbolic problems follow from papers dealing with homogeneous ones, like [7,11,18,28,34]; however, a new strategy, partly inspired by Laforest [26], consists in taking advantage of the Bressan-Liu-Yang  $L^1$  stability theory [8,9] in order to derive sharp error estimates for space-dependent source terms problems: see [1,2]. Here we address a model which is motivated by recent applications like for instance the ones presented in [12] or [31].

## 1.2 Main result and plan of the paper

The main theorem of this paper is the following.

**Theorem 1** *Under the assumptions (2), (6) and the sub-characteristic condition (5), the WB algorithm defined in §2.1 (see Fig. 3) satisfies the following local error estimates:*

$$\begin{aligned} \int_{x_1+t}^{x_2-t} |f_{\Delta x}^\pm(t, x) - f^\pm(t, x)| dx &\leq \int_{x_1}^{x_2} |f_{\Delta x}^\pm(0, x) - f_0^\pm(x)| dx \\ &+ \min \left\{ (K-1) \int_{x_1}^{x_2} |f_{\Delta x}^\pm(0, x) - f_0^\pm(x)| dx + \Delta x \cdot \mathcal{E}_1; 2t\sqrt{\Delta x} \cdot \mathcal{E}_2 \right\} \end{aligned} \quad (10)$$

for any  $x_1, x_2, t \in \mathbb{R}^2 \times \mathbb{R}^+$ ,  $\Delta x$  is sufficiently small (see (47)), with the definitions,

$$\begin{aligned} \mathcal{E}_1(t, x_1, x_2) &= (2KC_0 + K_0) \|k\|_{L^1(x_1, x_2)} + (2C_0 - 1) \|k\|_{L^1(x_1+t, x_2-t)}, \\ \mathcal{E}_2(t, x_1, x_2) &= \sqrt{C_0 \|k\|_{L^1(x_1, x_2)} A(t) + \sqrt{\Delta x} C_0 \|k\|_{L^1(x_1, x_2)} \|k\|_{L^\infty(x_1, x_2)}} \end{aligned}$$

where  $C_0$  stands for the Maxwellian gap (see (23)),  $C_1 = \frac{4}{3 \log(\frac{3}{2})}$  and

$$\begin{aligned} K &= \frac{1}{\max\{0, 1 - 4C_1 \|k\|_{L^1(x_1, x_2)}\}} \geq 1, \quad \mathcal{A}_0 = \text{TV} \{f_0^\pm; (x_1, x_2)\} + 2C_0 \|k\|_{L^1(x_1, x_2)}, \\ K_0 &= 1 + \frac{16K^2 C_1 \mathcal{A}_0}{3K + 1}, \quad A(t) = \frac{32}{C_0 t} \text{TV} \{f_0^\pm; (x_1, x_2)\} + \text{TV} \{k; (x_1, x_2)\}. \end{aligned}$$

Some comments on the main estimate (10) are now in order. Whenever  $4C_1\|k\|_{L^1(x_1,x_2)} \geq 1$  it results  $K = +\infty$  and then only the estimate with  $\mathcal{E}_2$  is meaningful. Clearly, thanks to (9), the initial error is bounded by  $\Delta x \cdot \text{TV}\{f_0^\pm; (x_1, x_2)\}$  as soon as the algorithm is initialized with a convenient sampling of  $f_0^\pm$ , see (68). The first estimate in (10) is **uniform in time  $t$  and first-order in  $\Delta x$** , but is meaningful only for "weak relaxation regime", where  $K$  remains finite. The complementary estimate is **linear in time  $t$  and half-order in  $\Delta x$** , which was to be expected as, in strong relaxation regime, the system (1) behaves like the reduced scalar conservation law for which optimal convergence order is studied in [32]. Let's see at once on an elementary example how they behave:

- Assume first that  $\text{TV}\{f_0^\pm; (x_1, x_2)\} = 0$ : the error estimate (10) boils down to

$$\int_{x_1+t}^{x_2-t} |f_{\Delta x}^\pm(t, x) - f^\pm(t, x)| dx \leq \min \left\{ \Delta x \cdot \mathcal{E}_1; 2t\sqrt{\Delta x} \cdot \mathcal{E}_2 \right\}.$$

Accordingly, there is no error at time  $t = 0$ . For semi-conductor models like (7), initial data usually are  $J(t = 0, \cdot) \equiv 0$  and  $\rho(t = 0, \cdot) = d$ , a piecewise-constant doping profile so the initial error vanishes, too.

- The initial Maxwellian gap is fixed to  $C_0 = \frac{1}{2}$  so that the last term of the time-uniform estimate  $\mathcal{E}_1$  cancels.
- On the coefficient  $k$ , we assume that  $\|k\|_{L^\infty(x_1,x_2)} = 1$  and  $\|k\|_{L^1(x_1,x_2)} = 1/8C_1$ , so that  $K = 2$ ; however, we don't restrict its total variation.

Based on all these assumptions, we can estimate the terms  $\mathcal{E}_1, \mathcal{E}_2$ . It is found that

$$\mathcal{A}_0 = \|k\|_{L^1(x_1,x_2)} = \frac{1}{8C_1}, \quad K = 2, \quad K_0 = \frac{15}{7}, \quad A(t) \equiv \text{TV}\{k; (x_1, x_2)\}$$

and then

$$\mathcal{E}_1 = \frac{2 + K_0}{8C_1}, \quad \mathcal{E}_2 = \frac{1}{16C_1} \left( \sqrt{16C_1 \text{TV}\{k; (x_1, x_2)\}} + \sqrt{\Delta x} \right).$$

Therefore, the time-dependent estimate dominates the other one as soon as

$$t \geq \sqrt{\Delta x} \frac{\mathcal{E}_1}{2\mathcal{E}_2} = \frac{29}{7 + 28\sqrt{C_1 \text{TV}\{k; (x_1, x_2)\}}/\Delta x}.$$

Hence, the time-uniform estimate  $\mathcal{E}_1$  is sharper in case the relaxation term is multiplied by a small, but oscillating (or at least, displaying areas of strong variation) coefficient. This meets with early implementations of the so-called "generalized Glimm scheme" by Weinan E [13] in a context of homogenization of scalar balance laws.

Such a time-uniform estimate is a consequence of applying the Bressan-Liu-Yang  $L^1$ -stability theory to a modified, homogeneous but non-conservative, version of system (1), see (11). A Godunov scheme can be set up, relying on a Riemann solver where the effects of the localized relaxation term are handled by means of a supplementary, static, jump relation, sometimes called "standing wave", or "zero-wave". Our estimate (10) shows that *in a context where  $\text{TV}(k)$  can be (locally) big, more accurate approximations can be obtained (perhaps beyond a certain time) by means of numerical schemes relying on this type of Riemann solvers, where the source term is handled like a "local scattering center" inducing a stationary discontinuity*, as suggested by Glimm and Sharp in [14].

The remaining part of the paper is entirely devoted to the proof of Theorem 1:  $\mathcal{E}_1$  is established in §2 and §3 whereas  $\mathcal{E}_2$  is derived by means of Kuznetsov's method [25,28] in §4.

More precisely, within the assumptions (2)–(4), the Riemann problem for the non-conservative reformulation of (1) is studied in §2.1, its positively-invariant domains are described in §2.2 (see Fig. 2) and a time-uniform total-variation estimate is shown in §2.3.

In §3 we set up the Lyapunov functional inspired by the Bressan-Liu-Yang  $L^1$  stability theory for general homogeneous  $n \times n$  hyperbolic systems: a technical Gronwall lemma allows to derive wave scattering estimates in §3.1, then accurate interaction estimates are proved in §3.2 (for simplicity, the analysis is specialized to the assumptions (2), (6)). This part culminates in §3.3 where the decay of our Lyapunov functional is established for possibly big  $BV$  data in weak relaxation regime (that is, for  $\|k\|_{L^1}$  suitably bounded).

A complementary estimate is needed for strong relaxation: §4.1 contains the derivation of entropy inequalities corresponding to the Godunov scheme for (11) and then, §4.2 indicates how they lead to another type of  $L^1$  error estimates, fully compatible with Kuznetsov's half-order accuracy for numerical approximations of scalar conservation laws. This last part is carried out in the more general framework of (2)–(4).

**Remark 2** (*algorithmic implications*) *Our new estimate  $\mathcal{E}_1$  in (10) appears like being specific to so-called "well-balanced" methods where source terms are concentrated onto interfacial discontinuities: the fact that  $\mathcal{E}_1$  doesn't grow in time and is independent of  $\text{TV}(k)$ , at least when  $\|k\|_{L^1}$  is small enough, suggests that this type of algorithms should outperform more conventional (time-splitting, see e.g. [15,28]) ones when  $k(x)$  display strong variations. This difference was already seen in [2] on a simpler model of damped wave equation. Moreover, it can give hints on why WB algorithms deliver high accuracy results on shallow water equations in presence of a steep topography: our results don't strictly apply to such a quasi-linear model, though.*

## 2 Construction of the Well-Balanced approximation

In this context, the WB approach consists in dealing with the inhomogeneous system (8) by means of a non-conservative homogeneous  $3 \times 3$  system, which turns out to be equivalent for smooth  $a(x)$ ,

$$\begin{cases} \partial_t \rho + \partial_x J & = 0, \\ \partial_t J + \partial_x \rho - 2g(\rho, J) \partial_x a & = 0, \\ \partial_t a & = 0, \end{cases} \quad a = a(x) \doteq \int_{-\infty}^x k(y) dy, \quad (11)$$

or equivalently, since  $G(f^+, f^-) = g((f^+ + f^-), (f^+ - f^-))$ ,

$$\partial_t f^\mp \mp \partial_x f^\mp \pm G(f^+, f^-) \partial_x a = 0, \quad \partial_t a = 0. \quad (12)$$

From assumption (2) one has that

$$a(x) \in BV(\mathbb{R}) \cap C(\mathbb{R}), \quad a_x \geq 0. \quad (13)$$

The characteristic speeds of system (12) are  $\lambda = \{-1, 0, 1\}$  with corresponding eigenvectors

$$\vec{r}_- = (0, 1, 0)^t, \quad \vec{r}_0 = (G, G, 1)^t, \quad \vec{r}_+ = (1, 0, 0)^t,$$

where we denote  $G(f^+, f^-) := g((f^+ + f^-), (f^+ - f^-))$ . The  $0$ -wave curves are those characteristic curves corresponding to  $\lambda = 0$ . One can easily check that the characteristic curves for  $\lambda = \pm 1$  are straight lines, while for  $\lambda = 0$  they are straight lines whenever  $A \equiv 0$  (see [2]).

**Remark 3** *The above mentioned procedure appears to trace back to Glimm and Sharp [14]. It consists in localizing a source term of bounded extent into a countable collection of "local scattering centers" rendered by Dirac masses, in order to integrate it inside a Riemann solver by means of an elementary (obviously very linearly degenerate) wave. It is extensively used for weakly nonlinear kinetic equations in Part II of [16].*

### 2.1 First considerations on the $3 \times 3$ Riemann problem

As usual, let

$$U_\ell = (f_\ell^-, f_\ell^+, a_\ell), \quad U_r = (f_r^-, f_r^+, a_r)$$

be a given a Riemann data for (12). The Riemann problem for system (12) is solved in terms of the three characteristic families, resulting in three waves: the two  $\pm 1$ -waves, with corresponding speed  $\pm 1$ , where only  $f^\pm$  can change its value; and the  $0$ -wave, corresponding to the stationary field of (12), evolving along the stationary equations

$$\partial_x f^\pm = k(x)G(f^-, f^+), \quad (14)$$

or equivalently, from (1):

$$\partial_x J = 0, \quad \partial_x \rho = 2k(x)g(\rho, J). \quad (15)$$

Notice that  $J$  is constant along stationary solutions. In terms of the diagonal variables  $f^\pm$ , the equilibrium curve  $J = A(\rho)$ , i.e. the level curve  $G = 0$ , is clearly expressed by

$$f^+ - f^- = A(f^+ + f^-).$$

By (5),  $|A'| < 1$ , and thanks to the implicit function theorem, the corresponding curve realizes a graph in the  $f^\pm$  plane. Indeed for each  $f^-$  the map  $\mathbb{R} \rightarrow \mathbb{R}$ ,

$$x \mapsto x - f^- - A(x + f^-)$$

is strictly increasing and tends to  $\pm\infty$  as  $x \rightarrow \pm\infty$ . Hence there exists a function  $f^+ = E(f^-)$ , globally defined on  $\mathbb{R}$ , along which the source vanishes. This map  $E$  is smooth and its derivative equals

$$E'(f^-) = \frac{1 + A'(y)}{1 - A'(y)}, \quad y = E(f^-) + f^-.$$

Thanks to (5),  $E' > 0$  thus  $E$  is strictly increasing, and the range of  $E'$  is  $(0, +\infty)$ . Now we observe the following interesting feature: from (3) it follows that

$$\frac{\partial G}{\partial f^-} = \partial_\rho g - \partial_J g > 0, \quad \frac{\partial G}{\partial f^+} = \partial_\rho g + \partial_J g < 0 \quad (16)$$

and therefore the gradient of  $G$  “points” in the bottom-right direction of the  $(f^-, f^+)$ -plane. Consequently  $G > 0$  below the graph and  $G < 0$  above the graph (see the arrows on Fig. 2).

The intermediate states in the Riemann fan are (see Fig. 1)

$$U_1 = (f_*^-, f_\ell^+, a_\ell), \quad U_2 = (f_r^-, f_*^+, a_r),$$

while the waves appearing in the solution are as follows:  $U_\ell$  and  $U_1$  are connected by a  $(-1)$ -wave of size  $\sigma_{-1}$ ,  $U_1$  and  $U_2$  are connected by a  $0$ -wave of size  $\sigma_0$ , and  $U_2$  and  $U_r$  are connected by a  $1$ -wave of size  $\sigma_1$  where

$$\begin{cases} \sigma_{-1} = f_*^- - f_\ell^- = (f_*^- - f_\ell^+) - (f_\ell^- - f_\ell^+) = J_\ell - J_* = \rho_{*,\ell} - \rho_\ell \\ \sigma_0 = a_r - a_\ell \\ \sigma_1 = f_r^+ - f_*^+ = (f_r^+ - f_r^-) - (f_*^+ - f_r^-) = J_r - J_* = \rho_r - \rho_{*,r}. \end{cases} \quad (17)$$

Here the “\*” signals that the corresponding value is related to the  $0$ -wave: more precisely,  $(\rho_{*,\ell}, J_*)$  and  $(\rho_{*,r}, J_*)$  denote the left and right states separated by the  $0$ -wave, respectively, in terms of the macroscopic variables  $(\rho, J)$ .

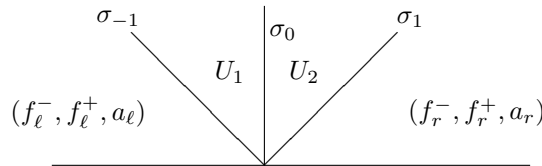


Figure 1. Schematic view of a Riemann problem for (12).

**Remark 4** *There exists a practical way to construct the Riemann problem for small  $\delta$ . If  $\delta = 0$ , there is no zero-wave thus  $U_1 = U_2$  is given by the state  $P = (f_r^-, f_\ell^+)$ , that corresponds to the intersection of the  $(-1)$ -wave issued from  $(f_\ell^-, f_\ell^+)$  and the  $(+1)$ -wave issued from  $(f_r^-, f_r^+)$ . Clearly here  $J_* = f_\ell^+ - f_r^-$ .*

*For  $\delta > 0$  small, the value of  $J_*$  can be obtained by perturbation as follows:*

- *In the very special case where  $G(P) = 0$ , that is, the intersection point  $P$  lies on the equilibrium curve, then the intermediate states  $U_1, U_2$  again coincide whatever is the value of  $\delta > 0$ .*



- Assume now that  $G(P) > 0$ , the other case being similar. Then, for a convenient  $J_*$ , one has to solve the equation

$$\partial_a \rho = 2g(\rho, J_*), \quad \rho(a_l) = \rho_{*,\ell}.$$

- (1) Define  $B(\rho, J)$  by integrating up to a constant  $\frac{1}{2g}$  with respect to  $\rho$ ,

$$B(\rho, J) = \int^{\rho} \frac{d\rho'}{2g(\rho', J)}. \quad (18)$$

For each value of the parameter  $J$ , the above function is well defined and monotone in a neighborhood of a point  $(\bar{\rho}, J)$  such that  $g(\bar{\rho}, J) \neq 0$ . Then, the left and right states of the 0-wave satisfy the relation

$$B(\rho_{*,r}, J_*) - B(\rho_{*,\ell}, J_*) = \int_{\rho_{*,\ell}}^{\rho_{*,r}} \frac{d\rho'}{2g(\rho', J)} = a_r - a_\ell.$$

Notice that, even if  $B$  is defined by (18) up to a function depending on  $J$ , the above difference does not depend on the choice of the particular function.

- (2) Now, by the very definition of  $f^\pm$ , we have

$$\rho_{*,\ell} + J_* = 2f_\ell^+, \quad \rho_{*,r} - J_* = 2f_r^- \quad (19)$$

then, by taking advantage of the fact that  $f^+$  (resp.  $f^-$ ) is constant across a  $-1$ -wave (resp. across a  $1$ -wave), we can write an implicit equation for  $J_*$ , in terms of the parameters  $f_\ell^+$  and  $f_r^-$ :

$$B(2f_r^- + J_*, J_*) - B(2f_\ell^+ - J_*, J_*) = a_r - a_\ell. \quad (20)$$

The equation (20) already appeared in the context of diffusive numerical approximations, in a slightly different form: see the book [16], page 150.

## 2.2 Shape of positively invariant domains for $3 \times 3$ Riemann problems

We first need to establish control on the amplitude of the WB approximations. A standard roadmap for doing so is to seek a positively invariant domain for the Riemann problem: see Fig. 2.

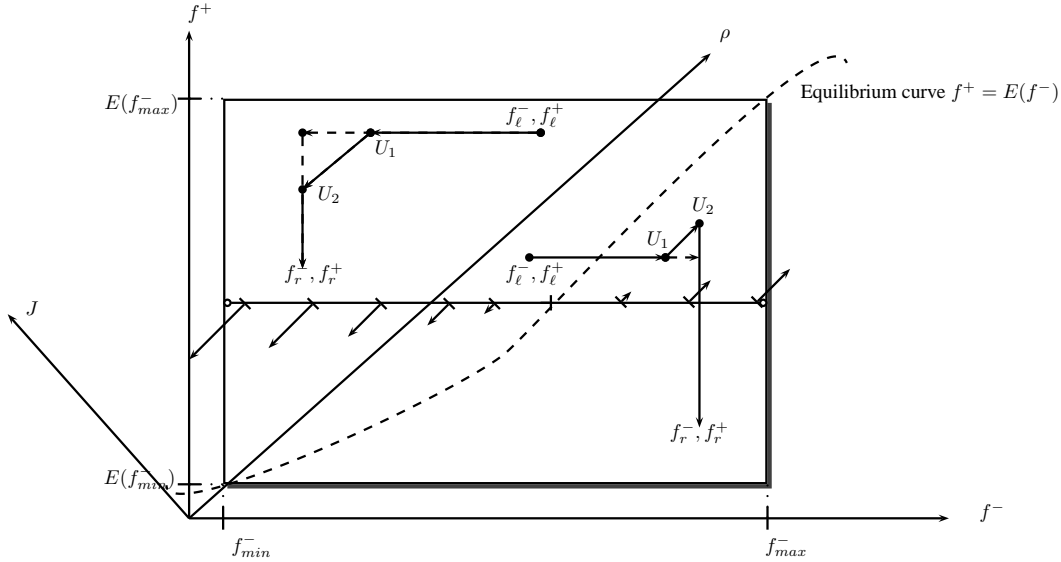


Figure 2. Invariant domain for  $3 \times 3$  system (12) and 2 sets of initial/final states. Diagonal arrows stand for the projection of the vector  $\vec{r}_0$  on the  $f^\pm$  plane.

**Proposition 1** Let  $f_{min}^- < f_{max}^-$ ,  $P_1 = (f_{min}^-, E(f_{min}^-))$ ,  $P_2 = (f_{max}^-, E(f_{max}^-))$ , and  $\delta := a_r - a_\ell > 0$ . Then the rectangle

$$D = [f_{min}^-, f_{max}^-] \times [E(f_{min}^-), E(f_{max}^-)] \quad (21)$$

is positively invariant for the unique solution of the Riemann problem. More precisely, for any pair of states  $(f_\ell^-, f_\ell^+)$  and  $(f_r^-, f_r^+) \in D$  and for  $\delta = a_r - a_\ell > 0$ , there exists a single choice of the intermediate states  $U_1, U_2$  for which one has

$$||f_r^- - f_\ell^-| - |\sigma_{-1}|| \leq C_0\delta, \quad ||f_r^+ - f_\ell^+| - |\sigma_1|| \leq C_0\delta \quad (22)$$

where  $C_0$  measures the ‘‘Maxwellian gap’’ in  $L^\infty$ :

$$C_0 = \max\{|G(f^-, f^+)|; (f^-, f^+) \in D\}. \quad (23)$$

*Proof.* Thanks to (19), all the intermediate states in the Riemann problem can be deduced from the knowledge of  $J_*$ , and the values  $f_*^+$  and  $f_*^-$  are defined by the identity

$$f_*^+ - f_r^- = f_\ell^+ - f_*^- = J_*.$$

Clearly, if  $G(f_r^-, f_\ell^+) = 0$ , then  $J_* = f_\ell^+ - f_r^-$  for every  $\delta > 0$ . On the other hand, when  $G(f_r^-, f_\ell^+) \neq 0$ , we aim at showing that the value of  $\tilde{J}$  is actually implicitly defined by the equation (20). Indeed let  $(f^-, f^+) \in D$  be such that  $G(f^-, f^+) > 0$  (that is,  $(f^-, f^+)$  **below** the equilibrium curve), the opposite case being similar, and define the function  $F$  as follows:

$$\begin{aligned} (J, \delta) \mapsto F(J, \delta; f^\pm) &= B(2f^- + J, J) - B(2f^+ - J, J) - \delta \\ &= \int_{2f^+ - J}^{2f^- + J} \frac{d\rho'}{2g(\rho', J)} - \delta \end{aligned} \quad (24)$$

(subscripts in  $f^\pm$  were dropped). One easily finds a solution of the particular equation for  $\delta = 0$ ,

$$F(J_0, 0; f^\pm) = 0 \Leftrightarrow J_0 = f^+ - f^-.$$

This solution corresponds to the case where there is no zero-wave because  $\delta = a_r - a_\ell$  vanishes. Let us verify that the following property holds:

$$0 \neq \frac{\partial F}{\partial J}(J_0, 0; f^\pm) = (\partial_J B + \partial_\rho B)(2f^- + J_0, J_0) - (\partial_J B - \partial_\rho B)(2f^+ - J_0, J_0),$$

but since  $2f^+ - J_0 = f^+ + f^- = 2f^- + J_0$ , this expression reduces to

$$\frac{\partial F}{\partial J}(J_0, 0; f^\pm) = 2\partial_\rho B(f^+ + f^-, f^+ - f^-) = \frac{1}{G(f^+, f^-)} \neq 0.$$

The implicit functions theorem ensures existence and uniqueness of a smooth function,

$$\tilde{J} : (0, \varepsilon) \times D \rightarrow \mathbb{R}, \quad \delta, f^\pm \mapsto \tilde{J}(\delta; f^\pm), \quad (25)$$

such that, for  $0 < \varepsilon \ll 1$  and  $G(f^+, f^-) \neq 0$ , one has  $J(0; f^\pm) = J_0 = f^+ - f^-$  and

$$F(\tilde{J}, \delta; f^\pm) = 0 \Leftrightarrow \int_{2f^+ - \tilde{J}}^{2f^- + \tilde{J}} \frac{d\rho'}{2g(\rho', \tilde{J})} = \delta. \quad (26)$$

Moreover, since  $\partial F / \partial \delta = -1$ , we make explicit the derivative of  $\tilde{J}$ :

$$\frac{\partial \tilde{J}}{\partial \delta} = \frac{1}{\frac{\partial F}{\partial J}(\tilde{J}, \delta; f^\pm)}, \quad \frac{\partial \tilde{J}}{\partial \delta}(\delta = 0; f^\pm) = G(f^+, f^-) > 0.$$

Under those smallness restrictions and  $(f^-, f^+)$  being fixed below the equilibrium curve, the restriction  $\delta \mapsto \tilde{J}$  is increasing because  $G(f^+, f^-) > 0$ . Moreover, as soon as  $\tilde{J}(\delta, f^\pm)$  is defined, the segment in the state space  $\rho, J$  along which the integral in (24) is computed, parametrized by  $\rho$  as follows,

$$[2f^+ - \tilde{J}, 2f^- + \tilde{J}] \ni \rho \mapsto (\rho, \tilde{J}) \quad (27)$$

does not intersect the equilibrium curve: otherwise, the integral in (26) would blow up, instead of being equal to  $\delta$ . Now we intend to verify that  $\frac{\partial F}{\partial J}(\tilde{J}, \delta; f^\pm) > 0$ , in order to establish that  $\tilde{J}$  is indeed increasing as soon as it is defined. Using the definition of  $B$ , (see 18), we have

$$\begin{aligned} 2\frac{\partial F}{\partial J}(J, \delta; f^\pm) &= 2(\partial_J B + \partial_\rho B)(2f^- + J, J) - 2(\partial_J B - \partial_\rho B)(2f^+ - J, J) \\ &= \int_{2f^+ - J}^{2f^- + J} \frac{1}{g^2(\rho', J)} |\partial_J g(\rho', J)| d\rho' + \frac{1}{g(2f^- + J, J)} + \frac{1}{g(2f^+ - J, J)}. \end{aligned} \quad (28)$$

Assuming again that  $G(f^+, f^-) > 0$ , if  $J$  is set to  $J = \tilde{J}(\delta; f^\pm)$  then the extrema of the integral above satisfy  $2f^+ - \tilde{J} < 2f^- + \tilde{J}$  (see Fig. 2). Hence we can take advantage of the last condition in (3),

$$\begin{aligned} \int_{2f^+ - J}^{2f^- + J} \frac{1}{g^2(\rho', J)} |\partial_J g(\rho', J)| d\rho' &\geq \int_{2f^+ - J}^{2f^- + J} \frac{1}{g^2(\rho', J)} |\partial_\rho g(\rho', J)| d\rho' \\ &\geq \int_{2f^+ - J}^{2f^- + J} \frac{1}{g^2(\rho', J)} \partial_\rho g(\rho', J) d\rho' \\ &= -\frac{1}{g(2f^- + J, J)} + \frac{1}{g(2f^+ - J, J)}. \end{aligned}$$

We can therefore estimate from below the integral in (28) and get

$$\frac{\partial F}{\partial J}(J, \delta; f^\pm) \geq \frac{1}{g(2f^+ - J, J)} > 0.$$

We now deduce that the function  $\delta \mapsto \tilde{J}(\delta; f^\pm)$  is actually defined on  $\mathbb{R}^+$ :

- there exists  $J_{\max}$  such that the interval in (27) has no intersection with the equilibrium curve for  $J_0 \leq J < J_{\max}$ .
- For  $J = J_{\max}$  there exists a point of the interval, say  $(\bar{\rho}, J_{\max})$  such that  $g(\bar{\rho}, J_{\max}) = 0$ . Since  $g$  is  $C^1$ , then  $g(\rho, J_{\max}) = O(1)(\rho - \bar{\rho})$  and therefore the corresponding integral is not finite.

Hence  $\tilde{J}(\delta)$  is defined for every  $\delta > 0$ , and  $\tilde{J}(\delta) \rightarrow J_{\max}$  as  $\delta \rightarrow \infty$ .

Analogously  $\delta \mapsto \tilde{J}(\delta; f^\pm)$  is decreasing when  $(f^-, f^+)$  is **above** the equilibrium curve, and that it is defined for all  $\delta > 0$  finite. The monotonicity of  $\tilde{J}$  implies that the domain  $D$  is positively invariant for the Riemann problem, see Fig. 2. Finally, concerning (22), we use (14) and (17) to estimate the jump in the  $f^\pm$  coordinate across the 0-wave, that is:

$$|f_*^+ - f_\ell^+| = |f_r^- - f_*^-| \leq \sup_D |G| \cdot \delta$$

and this yields,

$$||f_r^+ - f_\ell^+| - |\sigma_1|| = ||f_r^+ - f_\ell^+| - |f_r^+ - f_*^+|| \leq |f_*^+ - f_\ell^+| \leq C_0 \delta$$

with  $C_0$  as (23). An analogous estimate holds for  $\sigma_{-1}$  thus we end up with (22).  $\square$

### 2.3 Total variation estimate of the WB approximation

Let  $D$  be a rectangle as in (21) that contains the values of initial data  $(f_0^-, f_0^+)$ . By means of Prop. 1, since  $a_x \geq 0$ , up to a suitable choice of the initial data, the approximate solution remains confined inside the region  $D$

$$\forall t > 0, \quad (f^-, f^+)(t, \cdot) \in D. \quad (29)$$

Previous results on positively invariant domains for the  $3 \times 3$  Riemann problem for (12) allow to easily derive uniform bounds on the total variation of the corresponding WB approximation thanks to its peculiar structure. The method hereafter is taken from [19], p. 643, and we recall it now for completeness:

- Differentiate in time each equation on  $f^\pm$  in (12), multiply by  $(\text{sgn}(\partial_t f^-), \text{sgn}(\partial_t f^+))^t$  and then integrate on  $x \in \mathbb{R}$ . It comes

$$\forall t > 0, \quad \partial_t \int_{\mathbb{R}} (|\partial_t f^+(t, x)| + |\partial_t f^-(t, x)|) \cdot dx \leq 0,$$

thanks to the quasi-monotonicity of  $G$ .

- An easy observation is that, by taking moduli, it comes:

$$|\partial_x f^\pm| - |G(f^-, f^+) \partial_x a| \leq |\partial_t f^\pm| \leq |\partial_x f^\pm| + |G(f^-, f^+) \partial_x a|.$$

- It remains to integrate in space in order to obtain the estimate:

$$\int_{\mathbb{R}} (|\partial_x f^+(t, x)| + |\partial_x f^-(t, x)|) \cdot dx \leq \int_{\mathbb{R}} (|\partial_x f^+(0, x)| + |\partial_x f^-(0, x)|) \cdot dx + 4C_0 \text{TV } a$$

where  $C_0$  is given as in (23).

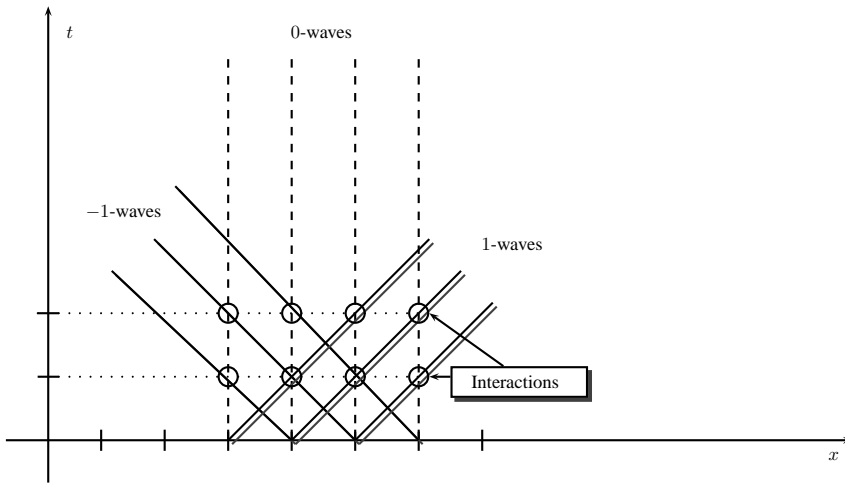


Figure 3. Schematic view of a WB approximate solution: circles indicate Riemann problems studied in Prop. 1. Since the Courant number is 1, constant states always lie in between them.

Recall also that  $\text{TV } a = \|k\|_{L^1}$ . The last point to clarify addresses the fact that we seek a BV-bound on an *approximation* depicted on Fig. 3, and not the exact solutions  $f^\pm$  of (12). However, since we chose to work with the Courant number equal to 1, the only difference between the WB approximation and the exact solutions lies in the sampling of initial data. Hence one gets the following estimate that does **not depend on time** for the WB approximation:

$$\boxed{\text{TV } f^+(t, \cdot) + \text{TV } f^-(t, \cdot) \leq \text{TV } f^+(0, \cdot) + \text{TV } f^-(0, \cdot) + 4C_0 \|k\|_{L^1}.} \quad (30)$$

This BV-bound is identical to the one obtained in [2] by means of rather different methods, though.

### 3 An $L^1$ error estimate through a Lyapunov functional

In this section we assume for simplicity that (6) and (5) hold, that is  $g(\rho, J) = A(\rho) - J$  with  $|A'| < 1$ . The extension to more general relaxation terms (2)–(4) follows without substantial difficulties, but at the price of tedious computations. We first study interactions between various patterns of waves for the system (12) in order to, in a second step, estimate the time-variation of the Lyapunov  $L^1$ -functional.

### 3.1 A lemma based on sub-characteristic condition

We start with a Gronwall-type lemma that exploits the sub-characteristic condition (5). Since  $|A'| < 1$  and  $\rho$  ranges over a compact set, there exists a positive constant  $\alpha < 1$  such that  $|A'| \leq \alpha$ .

**Lemma 1** *Let  $|A'| \leq \alpha$  and for  $a \in [a_\ell, a_r]$ ,  $\rho(a, J)$ , satisfy the parameter-dependent differential equation*

$$\forall J \in \text{Range}(\tilde{J}), \quad \frac{\partial \rho(a, J)}{\partial a} = 2(A(\rho) - J).$$

For some  $J_2 > J_1$ ,  $a \in [a_\ell, a_r]$ , define  $\phi(a) = \rho(a, J_2) - \rho(a, J_1)$ , and assume that

$$\phi(a_r) = \rho(a_r, J_2) - \rho(a_r, J_1) = J_2 - J_1 > 0. \quad (31)$$

Then the following inequalities hold:

$$\forall a \in [a_\ell, a_r], \quad \phi(a_r) - \phi(a) \leq -\tilde{c}(J_2 - J_1)(a_r - a) \quad (32)$$

with

$$\tilde{c} = (1 - \alpha) \frac{e^{2\alpha(a_\ell - a_r)} - 1}{\alpha(a_\ell - a_r)} > 0, \quad (33)$$

and

$$\forall a \in [a_\ell, a_r], \quad 0 < \phi(a) - \phi(a_r) \leq \tilde{C}(J_2 - J_1)(a_r - a) \quad (34)$$

with

$$\tilde{C} = (1 + \alpha) \frac{e^{2\alpha(a_r - a_\ell)} - 1}{\alpha(a_r - a_\ell)}. \quad (35)$$

*Proof.* As soon as  $\phi(a) > 0$  (which is true by continuity for  $a$  close to  $a_r$ ),  $\phi$  satisfies

$$\phi'(a) = 2(A(\rho(a, J_2)) - A(\rho(a, J_1))) - 2(J_2 - J_1) \leq 2\alpha\phi(a) - 2(J_2 - J_1). \quad (36)$$

The Gronwall lemma yields<sup>1</sup>

$$e^{2\alpha(a - a_r)}\phi(a_r) - \phi(a) \leq \frac{J_2 - J_1}{\alpha} (e^{2\alpha(a - a_r)} - 1). \quad (37)$$

By summing, subtracting and then using (37) and (31), we infer that

<sup>1</sup> Proof of (37). We have

$$\left( e^{-2\alpha(a - a_r)}\phi(a) \right)' = e^{-2\alpha(a - a_r)} (\phi' - 2\alpha\phi) \leq -2(J_2 - J_1)e^{-2\alpha(a - a_r)}.$$

By integrating in the interval  $[a, a_r]$  we find that

$$\phi(a_r) - e^{-2\alpha(a - a_r)}\phi(a) \leq \frac{J_2 - J_1}{\alpha} (1 - e^{-2\alpha(a - a_r)}).$$

It remains to multiply by  $e^{2\alpha(a - a_r)}$  on both sides and get (37).

$$\begin{aligned}
\phi(a_r) - \phi(a) &= \left[ e^{2\alpha(a-a_r)} \phi(a_r) - \phi(a) \right] - \left( e^{2\alpha(a-a_r)} - 1 \right) \phi(a_r) \\
&\leq \frac{J_2 - J_1}{\alpha} \left( e^{2\alpha(a-a_r)} - 1 \right) - \left( e^{2\alpha(a-a_r)} - 1 \right) \phi(a_r) \\
&= (J_2 - J_1) \frac{e^{2\alpha(a-a_r)} - 1}{\alpha} (1 - \alpha) \\
&= -(a_r - a)(J_2 - J_1) \left\{ \frac{e^{2\alpha(a-a_r)} - 1}{\alpha(a - a_r)} \right\} (1 - \alpha) \leq -\tilde{c}(a_r - a)(J_2 - J_1)
\end{aligned}$$

with  $\tilde{c}$  as in (33). This proves (32). Such inequality, rewritten as

$$\phi(a_r) + (a_r - a)(J_2 - J_1)\tilde{c}(1 - \alpha) \leq \phi(a), \quad a \in [a_\ell, a_r],$$

shows also that  $\phi$  remains positive, hence the above argument is valid as soon  $\phi(a)$  is defined. To prove (34), we start again from computing  $\phi'$  and find the opposite inequality to (36):

$$\phi'(a) \geq -2\alpha\phi(a) - 2(J_2 - J_1),$$

where we used also that  $\phi(a) > 0$ . The Gronwall lemma yields

$$e^{-2\alpha(a-a_r)}\phi(a_r) - \phi(a) \geq -\frac{J_2 - J_1}{\alpha} \left( e^{-2\alpha(a-a_r)} - 1 \right).$$

By proceeding as in the first part of the proof, we obtain

$$\begin{aligned}
\phi(a_r) - \phi(a) &= \left[ e^{2\alpha(a_r-a)} \phi(a_r) - \phi(a) \right] - \left( e^{2\alpha(a_r-a)} - 1 \right) \phi(a_r) \\
&\geq -\frac{J_2 - J_1}{\alpha} \left( e^{2\alpha(a_r-a)} - 1 \right) - \left( e^{2\alpha(a_r-a)} - 1 \right) \phi(a_r) \\
&= -(J_2 - J_1) \left\{ \frac{e^{2\alpha(a_r-a)} - 1}{\alpha(a_r - a)} \right\} (1 + \alpha) (a_r - a) \geq -\tilde{C}(J_2 - J_1)(a_r - a)
\end{aligned}$$

with  $\tilde{C}$  as in (35). It remains to change sign in the inequality above, and then get (34).  $\square$

**Remark 5** The function  $\phi(a)$  quantifies the dependence of  $\rho$  with respect to the parameter  $J$ . Accordingly, one can notice that, formally,

$$\phi(a) \simeq \frac{\partial \rho}{\partial J}(J_2 - J_1), \quad \phi(a_r) - \phi(a) \simeq \frac{\partial^2 \rho}{\partial J \partial a}(J_2 - J_1)(a_r - a).$$

Hence, assuming all the necessary smoothness, the estimates (32) and (34) state in a rigorous manner the informal statement that the mixed derivative is strictly negative,

$$-\tilde{C} \leq \frac{\partial^2 \rho}{\partial a \partial J} \leq -\tilde{c} < 0.$$

**Lemma 2** Consider the elementary interaction pattern displayed on Fig. 4, with  $\delta = a_r - a_\ell > 0$ : the conservation law holds,

$$\boxed{|\tilde{\sigma}_1| + |\tilde{\sigma}_{-1}| = |\sigma_1|}. \quad (38)$$

In particular, the reflected wave has always opposite sign with respect to the transmitted one:

$$\text{sgn}(\tilde{\sigma}_1) = -\text{sgn}(\tilde{\sigma}_{-1}).$$

Moreover, the amplitude of the reflected wave is estimated by

$$|\tilde{\sigma}_{-1}| \leq \tilde{C}_1 \delta |\sigma_1| \quad \tilde{C}_1 = \frac{\tilde{C}}{2} = (1 + \alpha) \frac{e^{2\alpha\delta} - 1}{2\alpha\delta} \quad (39)$$

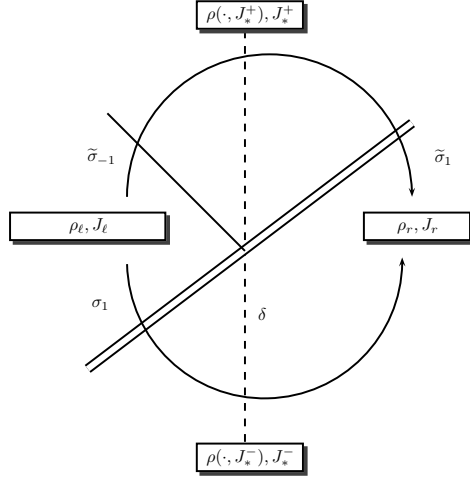


Figure 4. Interaction pattern corresponding to the scattering of a linear wave  $\sigma_1$  by a source-term discontinuity of size  $\delta = a_r - a_\ell > 0$ .

with  $\tilde{C}$  as in (35). The symmetric case of a  $(-1)$ -wave interacting with the 0-wave is completely analogous. Such a lemma expresses a strong conservation law for the scattering process where an incoming wave  $\sigma_1$  is scattered by a zero-wave  $\delta$  giving birth to reflected/transmitted waves  $\tilde{\sigma}_{\pm 1}$ .

*Proof.* It splits into several steps.

- Let  $a \mapsto \rho(a, J)$ ,  $a \in [a_\ell, a_r]$ , stand for solutions of the ODE problem along the 0-wave associated with a flux value  $J$ . More precisely, we have, before and after interaction, respectively:

$$\begin{aligned} \frac{\partial}{\partial a} \rho(a, J_*^-) &= 2(A(\rho^-) - J_*^-), & \rho^-(a_r, J_*^-) &= \rho_r, \\ \frac{\partial}{\partial a} \rho(a, J_*^+) &= 2(A(\rho^+) - J_*^+), & \rho(a_r, J_*^+) &= \rho_r - \tilde{\sigma}_1. \end{aligned}$$

Notice also that  $J_*^- = J_r$ , and so  $J_*^- - J_*^+ = \tilde{\sigma}_1$ .

- Assume now that  $\tilde{\sigma}_1 > 0$ : we can apply Lemma 1 with

$$\phi(a) = \rho(a, J_*^-) - \rho(a, J_*^+).$$

Since  $\phi(a_r) = J_*^- - J_*^+ = \tilde{\sigma}_1 > 0$ , the estimate (32) for  $a = a_\ell$  leads to

$$\phi(a_r) - \phi(a_\ell) \leq -\tilde{c} \delta \tilde{\sigma}_1, \quad (40)$$

while the estimate (34) for  $a = a_\ell$  lead to

$$\phi(a_r) - \phi(a_\ell) \geq -\tilde{C} \delta \tilde{\sigma}_1. \quad (41)$$

- Oppositely, if  $\tilde{\sigma}_1 < 0$ , Lemma 1 can still be applied with  $\tilde{\phi}(a) = \rho(a, J_*^+) - \rho(a, J_*^-)$ , and

$$-\tilde{C} \delta |\tilde{\sigma}_1| \leq \tilde{\phi}(a_r) - \tilde{\phi}(a_\ell) \leq -\delta |\tilde{\sigma}_1| \tilde{c},$$

leading to

$$-\tilde{C} \delta \tilde{\sigma}_1 \geq \phi(a_r) - \phi(a_\ell) \geq -\delta \tilde{\sigma}_1 \tilde{c}. \quad (42)$$

- By equating  $\rho_r - \rho_\ell$  before and after the interaction, and by using the definition of the size of waves, (17) in terms of jumps of  $\rho$ , we get

$$\begin{aligned} \rho_r - \rho_\ell &= \sigma_1 + (\rho(a_r, J_*^-) - \rho(a_\ell, J_*^-)), & (\text{lower curved arrow on Fig. 4}), \\ &= \tilde{\sigma}_1 + \tilde{\sigma}_{-1} + (\rho(a_r, J_*^+) - \rho(a_\ell, J_*^+)), & (\text{upper curved arrow on Fig. 4}). \end{aligned}$$

Henceforth, one deduces:

$$\tilde{\sigma}_1 + \tilde{\sigma}_{-1} + (\rho(a_r, J_*^+) - \rho(a_\ell, J_*^+)) = \sigma_1 + (\rho(a_r, J_*^-) - \rho(a_\ell, J_*^-)). \quad (43)$$

Moreover, by equating  $J_r - J_\ell$  before and after the interaction, we find that

$$\tilde{\sigma}_1 - \tilde{\sigma}_{-1} = \sigma_1. \quad (44)$$

Subtracting (44) from (43), it comes

$$2\tilde{\sigma}_{-1} = (\rho(a_r, J_*^-) - \rho(a_\ell, J_*^-)) - (\rho(a_r, J_*^+) - \rho(a_\ell, J_*^+)). \quad (45)$$

• If  $\tilde{\sigma}_1 > 0$ , one uses (40) and gets

$$2\tilde{\sigma}_{-1} \leq -\tilde{c}\tilde{\sigma}_1\delta < 0,$$

and therefore, from (41):

$$2|\tilde{\sigma}_{-1}| \leq \tilde{C}\delta|\tilde{\sigma}_1|. \quad (46)$$

• while, for  $\tilde{\sigma}_1 < 0$ , the second inequality in (42) leads to

$$2\tilde{\sigma}_{-1} = (\rho(a_r, J_*^-) - \rho(a_\ell, J_*^-)) - (\rho(a_r, J_*^+) - \rho(a_\ell, J_*^+)) \geq \tilde{c}|\tilde{\sigma}_1|\delta > 0,$$

and therefore, using the first inequality in (42), we get again (46).

From the above study of the sign of  $\tilde{\sigma}_{\pm 1}$  we conclude that

$$\text{sgn}(\tilde{\sigma}_1) = -\text{sgn}(\tilde{\sigma}_{-1}).$$

By using again (44) one finds that  $\text{sgn}(\tilde{\sigma}_1) = \text{sgn}(\sigma_1)$  and hence we get (38):

$$|\tilde{\sigma}_1| + |\tilde{\sigma}_{-1}| = |\sigma_1|.$$

• Finally, to complete the estimate (39) on the amplitude of the reflected wave, it is enough to recall (46) and use (38) to get

$$|\tilde{\sigma}_{-1}| \leq \tilde{C}_1\delta|\tilde{\sigma}_1| \leq \tilde{C}_1\delta|\sigma_1|.$$

□

### 3.2 Accurate interaction estimates for WB approximations

Lemma 2 allows to consider more intricate interaction patterns, as we shall see hereafter.

**Proposition 2** *Let  $U_\ell$  and  $U_m$  be connected by a complete Riemann pattern of size  $q_{\pm 1}^-$  and  $q_0$ . Let  $U_m$  and  $U_r$  be connected by a single wave as described in the cases below. Finally let  $q_{\pm 1}^+$  be the sizes of the  $\pm 1$ -waves solving the Riemann problem for  $U_\ell, U_r$  (see Figures 5 and 6). Under the hypotheses of Proposition 1 and for*

$$2\|k\|_\infty \Delta x \leq \log\left(\frac{3}{2}\right), \quad C_1 = \frac{4}{3\log(3/2)} \simeq 3.29 \quad (47)$$

then the following properties hold:

(a) *If  $U_m$  and  $U_r$  be connected by a  $-1$ -wave of size  $\sigma_{-1}$ , then*

$$|q_{-1}^+ - q_{-1}^- - \sigma_{-1}| = |q_1^+ - q_1^-| \leq C_1 q_0 |\sigma_{-1}|. \quad (48)$$

(b) *If  $U_m$  and  $U_r$  be connected by a  $0$ -wave of size  $\sigma_0$ , then*

$$|q_{-1}^+ - q_{-1}^-| = |q_1^+ - q_1^-| \leq C_1 |q_1^-| \sigma_0. \quad (49)$$



(c) If  $U_m$  and  $U_r$  be connected by a 1-wave of size  $\sigma_1$ , then

$$q_{-1}^+ = q_{-1}^-, \quad q_1^+ = q_1^- + \sigma_1.$$

Deriving an explicit constant  $C_1$ , given in (47), was the main reason for setting up Lemma 1.

*Proof.* Denote by  $J_*^-$ ,  $J_*^+$  the intermediate values of  $J$  in the Riemann problem for  $(U_\ell, U_m)$  and  $(U_\ell, U_r)$  respectively. Then the following identities are valid for the sizes of waves:

$$\begin{cases} q_{-1}^+ - q_{-1}^- &= J_*^- - J_*^+, \\ q_1^+ - q_1^- &= (J_*^- - J_*^+) + (J_r - J_m). \end{cases} \quad (50)$$

Indeed, it is sufficient to recall the definitions (17); for instance we get  $q_{-1}^+ - q_{-1}^- = (J_\ell - J_*^+) - (J_\ell - J_*^-)$  and hence the first identity. Similar for the second one. We proceed in order of increasing difficulty.

(1) **Case (c).** One has  $J_r - J_m = \sigma_1$  and  $J_*^- = J_*^+$ . Hence the claim simply follows from (50), being  $q_1^+ - q_1^- - \sigma_1 = 0 = q_{-1}^+ - q_{-1}^-$ .

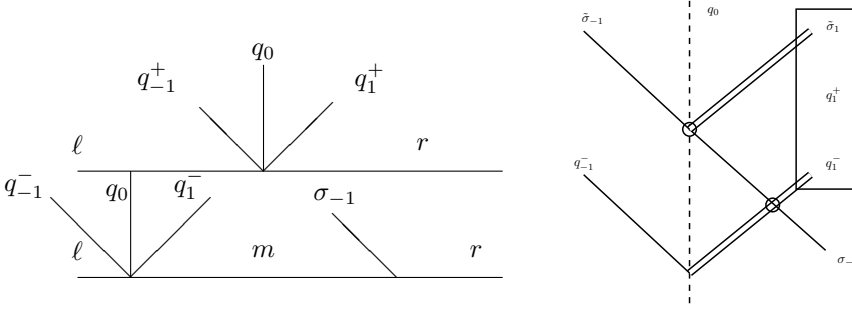


Figure 5. Illustration of Case (a).

(2) **Case (a).** Recalling the definition of sizes (17) one has that  $\sigma_{-1} = J_m - J_r$ , so identities (50) lead to

$$q_1^+ - q_1^- = q_{-1}^+ - q_{-1}^- - \sigma_{-1}. \quad (51)$$

- Let us proceed by letting both the linear waves  $\sigma_{-1}$  and  $q_1^-$  cross each other (without changing their size). Later,  $\sigma_{-1}$  interacts with  $q_0$ : denote by  $\tilde{\sigma}_{\pm 1}$  the resulting waves so that the final sizes  $q_{\pm 1}^+$  satisfy

$$q_{\pm 1}^+ = \tilde{\sigma}_{\pm 1} + q_{\pm 1}^- \quad \Rightarrow \quad q_1^+ - q_1^- = \tilde{\sigma}_1 = q_{-1}^+ - q_{-1}^- - \sigma_{-1} = \tilde{\sigma}_{-1} - \sigma_{-1}$$

Accordingly, equality (51) rewrites  $q_1^+ - q_1^- = \tilde{\sigma}_1 = \tilde{\sigma}_{-1} - \sigma_{-1}$ . Applying (39) in Lemma 2 we get

$$|\tilde{\sigma}_1| \leq \tilde{C}_1 q_0 |\sigma_{-1}|,$$

so (48) holds with

$$C_1 \geq \tilde{C}_1(q_0) \quad \forall q_0.$$

The choice of the constant  $C_1$  will be finalized in the next Case (b).

(3) **Case (b).** Here, the scattering processes related to 2 distinct zero-waves, of sizes  $q_0$  and  $\sigma_0$  respectively, are "glued" altogether into a unique one. In a linear context, this can be processed by means of the "Redheffer products" already set up in [17].

- In this last case we have  $J_r = J_m$  and hence (50) reduces to

$$q_1^+ - q_1^- = q_{-1}^+ - q_{-1}^- = J_*^- - J_*^+, \quad (52)$$

which already yields the left part of (49).

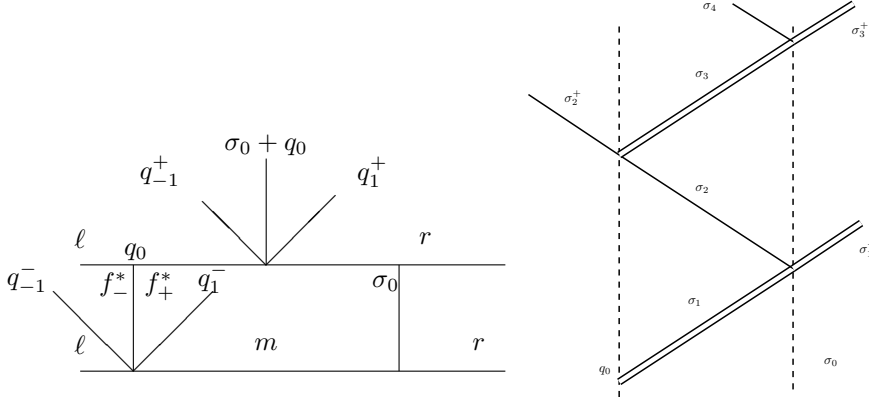


Figure 6. Illustration of Case (b).

- Without loss of generality, one can safely assume that  $q_{-1}^- = 0$ : indeed, let us show that the seemingly more intricate case  $q_{-1}^- \neq 0$  simply reduces to it. Let  $\tilde{q}_{\pm 1}^+$  be the result of the reduced interaction involving only  $q_0, q_1^-, \sigma_0$ ; estimates involve only quantities  $\tilde{q}_1^+ - q_1^-$  and  $\tilde{q}_{-1}^+$ . Now, if  $q_{-1}^- \neq 0$ , then by linearity, resulting waves  $q_{\pm 1}^+$  as in Figure 6 satisfy

$$q_1^+ = \tilde{q}_1^+, \quad q_{-1}^+ = \tilde{q}_{-1}^+ + q_{-1}^-.$$

- Accordingly, we assume the situation depicted in Fig. 6 where  $\sigma_1 = q_1^-$ : Lemma 2 gives that,

$$\text{sgn}(\sigma_1^+) = \text{sgn}(\sigma_1) = -\text{sgn}(\sigma_2).$$

By induction, this property propagates at each scattering event, so for all  $n \in \mathbb{N}$ ,

$$\text{sgn}(\sigma_n^+) = \text{sgn}(\sigma_n) = -\text{sgn}(\sigma_{n+1}), \quad \text{sgn}(\sigma_{2n+1}^+) = \text{sgn}(\sigma_1) = -\text{sgn}(\sigma_{2n}^+).$$

Next, let's consider quadratic interaction estimates on the right zero-wave (with size  $\sigma_0$ ):

$$|\sigma_{2n+2}| \leq \tilde{C}_1(\sigma_0)\sigma_0 |\sigma_{2n+1}| \leq \tilde{C}_1(\sigma_0)\tilde{C}_1(q_0)\sigma_0 q_0 |\sigma_{2n}|$$

where  $\tilde{C}_1$  is as in (39), and one has

$$\tilde{C}_1(x) = (1 + \alpha) \frac{\exp(2\alpha x) - 1}{2\alpha x} \simeq (1 + \alpha) \quad \text{as } x \rightarrow 0.$$

Also, we can estimate both  $\sigma_0, q_0$  as follows:  $\sigma_0, q_0 \leq \Delta x \|k\|_\infty$ . Hence

$$|\sigma_{2n+2}| \leq \gamma |\sigma_{2n}|, \quad \gamma \doteq (\tilde{C}_1(\bar{x})\bar{x})^2, \quad \bar{x} \doteq \Delta x \cdot \|k\|_\infty$$

and notice that  $\gamma \rightarrow 0$  as  $\Delta x \rightarrow 0$  in weak relaxation regime. This immediately implies that

$$|\sigma_{2n+2}| \leq \gamma^n |\sigma_2| \leq \gamma^n \tilde{C}_1(\bar{x}) |q_1^-| \cdot |\sigma_0|.$$

- It now remains to sum all the even terms:

$$\begin{aligned} \left| \sum_{n=1}^{\infty} \sigma_{2n}^+ \right| &= \sum_{n=1}^{\infty} |\sigma_{2n}^+| \leq \sum_{n=1}^{\infty} |\sigma_{2n}| \\ &\leq \left( \sum_{n=1}^{\infty} \gamma^{n-1} \right) \tilde{C}_1(\bar{x}) |q_1^-| \cdot |\sigma_0| = \underbrace{\left( \frac{\tilde{C}_1(\bar{x})}{1 - \gamma} \right)}_{=C_1} |q_1^-| \cdot |\sigma_0|. \end{aligned}$$

To estimate the above defined constant  $C_1$ , we assume that  $\bar{x}$  satisfies

$$\tilde{C}_1(\bar{x})\bar{x} = (1 + \alpha) \frac{\exp(2\alpha\bar{x}) - 1}{2\alpha} \leq \frac{1}{2}.$$

Since the above function of  $\alpha$  is increasing, and  $\alpha < 1$ , we let  $\alpha \rightarrow 1$  in the previous equation and define our quantities to be uniform in  $\alpha$  as follows:

$$\tilde{C}_1(\bar{x})\bar{x} = \exp(2\bar{x}) - 1 = \frac{1}{2},$$

that gives

$$\Delta x \cdot \|k\|_\infty = \bar{x} = \frac{1}{2} \log\left(\frac{3}{2}\right), \quad \tilde{C}_1(\bar{x}) = \frac{1}{\log\left(\frac{3}{2}\right)}.$$

Recalling the above definition of  $\gamma$ , we conclude that  $\gamma \leq 1/4$  and therefore

$$C_1 = \frac{4}{3} \tilde{C}_1(\bar{x}) = \frac{4}{3 \log\left(\frac{3}{2}\right)}.$$

Finally, call  $L$  the distance separating both the zero-waves  $q_0$  and  $\sigma_0$ : one can pass to the limit  $L \rightarrow 0$ . By compactness, it converges to a non-interacting Riemann fan endowed with a zero-wave of size  $q_0 + \sigma_0$ . The size of the reflected wave reads:

$$|q_{-1}^+| = \sum_{n=1}^{\infty} |\sigma_{2n}^+| \leq C_1 |q_1^-| \cdot |\sigma_0|,$$

and the estimate (49) follows after taking (52) into account. □

**Remark 6** *It is important to stress that there exists more direct manners to establish a quadratic estimate for the interaction of approaching waves like in Proposition 2. Indeed, let's consider for instance the proof of (48): one may proceed by just recalling the definition of sizes (17), so that  $\sigma_{-1} = J_m - J_r = f_r^- - f_m^-$  and then the second identity in (50) becomes*

$$q_1^+ - q_1^- = (J_*^- - J_*^+) - \sigma_{-1}.$$

Recalling the definition of  $\tilde{J}$ , see (25), the quantities  $J_*^+$ ,  $J_*^-$  are given by

$$J_*^+ = \tilde{J}(q_0, f_\ell^+, f_r^-) = \tilde{J}(q_0, f_\ell^+, f_m^- + \sigma_{-1}), \quad J_*^- = \tilde{J}(q_0, f_\ell^+, f_m^-),$$

therefore, by the mean-value theorem, one derives:

$$J_*^- - J_*^+ = -\frac{\partial \tilde{J}}{\partial(f^-)}(q_0, f_\ell^+, f_m^- + \theta\sigma_{-1})\sigma_{-1}, \quad \theta \in (0, 1).$$

Notice that for  $q_0 = 0$  one has  $\tilde{J}(0, f^\pm) = f^+ - f^-$ , so we substitute

$$\frac{\partial \tilde{J}}{\partial(f^-)}(0, f^\pm) \equiv -1, \quad \text{for all } f^-,$$

into the former expression. Accordingly we obtain:

$$q_1^+ - q_1^- = -\sigma_{-1} \left[ \frac{\partial \tilde{J}}{\partial(f^-)}(q_0, f_\ell^+, f_m^- + \theta\sigma_{-1}) - \frac{\partial \tilde{J}}{\partial(f^-)}(0, f_\ell^+, f_m^- + \theta\sigma_{-1}) \right],$$

which finally furnishes,

$$|q_1^+ - q_1^-| \leq |\sigma_{-1}| |q_0| \cdot \sup \left| \frac{\partial^2 \tilde{J}}{\partial \delta \partial (f^-)}(\delta, f^\pm) \right|.$$

However, the main issue with such a computation lies in the fact that the resulting interaction constant cannot be easily expressed (see also Appendix A).

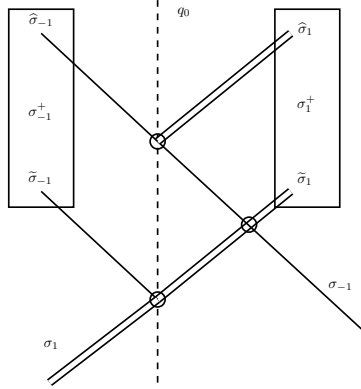


Figure 7. Schematic representation of the triple interaction

The following Proposition establishes a fundamental decay property:

**Proposition 3 (Multiple interaction)** *Assume that a 1-wave, a 0-wave and a -1-wave interact. Let  $\sigma_{-1}^-, \sigma_{-1}^+$  be the sizes of the incoming waves and  $\sigma_{-1}^+, \sigma_{-1}^-$  be the ones of the outgoing waves. Then*

$$|\sigma_{-1}^+| + |\sigma_{-1}^-| \leq |\sigma_{-1}^-| + |\sigma_{-1}^+|. \quad (53)$$

Besides, for  $\delta = a_r - a_\ell$ , one has

$$\begin{cases} |\sigma_{-1}^+| - |\sigma_{-1}^-| & \leq C_1 \delta (|\sigma_{-1}^-| + |\sigma_{-1}^+|) \\ |\sigma_{-1}^-| - |\sigma_{-1}^+| & \leq C_1 \delta (|\sigma_{-1}^-| + |\sigma_{-1}^+|) \end{cases} \quad (54)$$

*Proof.* One proceeds by letting interactions occur two at a time, and then collect the result: see Fig. 7.

- The first step is identical to the situation described in Lemma 2. Accordingly, the conclusion (38) holds for the present case, too. After the former interaction occurred, the wave of size  $\tilde{\sigma}_{-1}$  will cross the (-1)-wave of size  $\sigma_{-1}^-$  without changing size by linearity. The interaction between this last wave and the 0-wave produces two new waves,  $\hat{\sigma}_{\pm 1}$ . Analogously, they satisfy

$$|\hat{\sigma}_1| + |\hat{\sigma}_{-1}| = |\sigma_{-1}^-|. \quad (55)$$

- Due to the linearity of  $\pm 1$ -waves, no other interaction can occur. The sizes of the outgoing waves  $\sigma_{-1}^+, \sigma_{-1}^-$  must satisfy

$$\sigma_{-1}^+ = \tilde{\sigma}_{-1} + \hat{\sigma}_{-1}, \quad \sigma_{-1}^- = \tilde{\sigma}_{-1} + \hat{\sigma}_1.$$

Therefore, collecting (38) and (55), we finally get (53):

$$\begin{aligned} |\sigma_{-1}^+| + |\sigma_{-1}^-| & \leq |\tilde{\sigma}_{-1}| + |\hat{\sigma}_{-1}| + |\tilde{\sigma}_{-1}| + |\hat{\sigma}_1| \\ & = |\sigma_{-1}^-| + |\sigma_{-1}^+|. \end{aligned}$$

- Finally let us prove (54) for the 1-family, the other one being analogous. From the construction above and Prop. 2, it is easy to deduce that

$$|\tilde{\sigma}_1 - \sigma_1^-| \leq C_1 |\sigma_1^-| \delta, \quad |\hat{\sigma}_1| \leq C_1 |\sigma_{-1}^-| \delta.$$

One has

$$|\sigma_1^+| - |\sigma_1^-| \leq |\tilde{\sigma}_1| + |\hat{\sigma}_1| - |\sigma_1^-| \leq |\tilde{\sigma}_1 - \sigma_1^-| + |\hat{\sigma}_1|,$$

therefore we conclude thanks to the above estimates on  $|\tilde{\sigma}_1 - \sigma_1^-|$  and on  $|\hat{\sigma}_1|$ .

□

### 3.3 Lyapunov functional and error estimate for weak relaxation

Here, the main objective is to quantify the gap between 2 WB-approximations obtained with 2 different grid parameters  $(\Delta x)_1, (\Delta x)_2$ . Two approximations  $f_1^\pm, b_1(x)$  and  $f_2^\pm, b_2(x)$  being given, at each point  $(t, x)$ , one considers the “transversal Riemann problem” for (12) with left/right data:

$$f_1^\pm(t, x), b_1(x), \quad f_2^\pm(t, x), b_2(x).$$

We assume that  $b_1, b_2$  are piecewise constant, non-decreasing, with jumps located in  $(\Delta x)_1\mathbb{Z}$  and  $(\Delta x)_2\mathbb{Z}$  respectively, and that they satisfy

$$\text{TV } b_1 \leq \text{TV } a, \quad \text{TV } b_2 \leq \text{TV } a.$$

On the approximate initial data, we assume that

$$\text{TV } f_i^-(0, \cdot) + \text{TV } f_i^+(0, \cdot) \leq \text{TV } (f_0^+) + \text{TV } (f_0^-), \quad i = 1, 2.$$

Let

$$q_{\pm 1}(t, x), \quad q_0(x) = b_2(x) - b_1(x)$$

stand for the corresponding “transversal wave-strengths”, and consider, for instance, that  $f_1^-$  has a jump of size  $\sigma$  at the point  $(t, x_\alpha)$ : see Figure 8. In order to correctly devise the weights involved in the Lyapunov functional, it is necessary to know how the “transversal wave-strengths” evolve according to all the jumps present in both  $f_1^\pm, b_1(x)$  and  $f_2^\pm, b_2(x)$ .

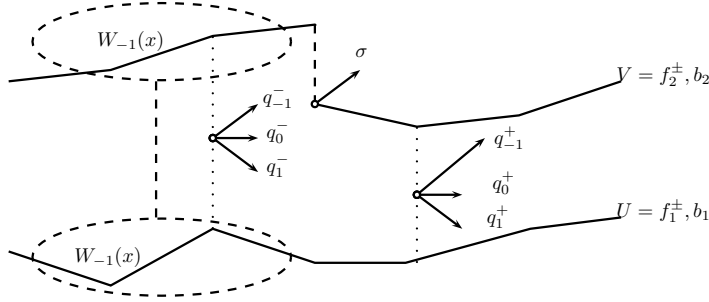


Figure 8. Interaction between a “transversal Riemann problem” (left) and a  $-1$ -wave resulting in the new Riemann problem (right) illustrating the simplest situation described by Prop. 2, Case (c).

In the sequel, we use all the standard notations by Bressan [8]; the only exception is that the characteristic families are numbered  $-1, 0, 1$  for obvious reasons. Let  $U, V$  stand for  $(f_1^-, f_1^+, b_1)$  and  $(f_2^-, f_2^+, b_2)$  respectively. We write  $\sigma_i^\alpha$  for the size a front located at  $x^\alpha$ , of the family  $i \in \{-1, 0, 1\}$ ; zero-waves are measured simply by the jump of  $b_1(x)$  or  $b_2(x)$ , respectively for  $U$  or  $V$ . Recall that all the  $\sigma_0^\alpha$  are positive, since  $b_1(x)$  and  $b_2(x)$  are assumed to be monotone, non-decreasing.

The Lyapunov functional  $\Phi[U, V]$  reads, for  $x_1 < x_2$  and  $t \leq T = (x_2 - x_1)/2$ :

$$t \mapsto \Phi[U, V](t) = \int_{x_1+t}^{x_2-t} |q_0(x)| W_0(t, x) dx + \sum_{i=\pm 1} \int_{x_1+t}^{x_2-t} |q_i(t, x)| W_i(x) dx, \quad (56)$$

where the weights  $W_i$  are defined as follows:

$$W_0(t, x) = 1 + \kappa_1 A_0(t, x) + \kappa_2 (Q(U) + Q(V)), \quad W_i(x) = 1 + \kappa_1 A_i(x), \quad i = -1, 1$$

and

$$\begin{aligned}
A_0(t, x) &= \sum_{x_\alpha < x} |\sigma_1^\alpha| + \sum_{x_\alpha > x} |\sigma_{-1}^\alpha|, \\
A_{-1}(x) &= \sum_{x_\alpha < x} \sigma_0^\alpha, \quad [0\text{-fronts on the left of } x], \\
A_1(x) &= \sum_{x_\alpha > x} \sigma_0^\alpha, \quad [0\text{-fronts on the right of } x].
\end{aligned}$$

The sums above extend over all jumps in  $U$  and  $V$ . An estimate for  $A_{\pm 1}$  reads:

$$A_{\pm 1}(x) \leq \text{TV } b_1 + \text{TV } b_2 \leq 2\text{TV } a.$$

On the other hand, an estimate on  $A_0$  goes as follows. By defining

$$\mathcal{A}_0 \doteq \text{TV } f_0^- + \text{TV } f_0^+ + 2C_0 \text{TV } a,$$

and recalling (30), one obtains

$$\begin{aligned}
A_0(t, x) &\leq L_\pm(t; U) + L_\pm(t; V) \\
&\leq \text{TV } f_1^-(0, \cdot) + \text{TV } f_1^+(0, \cdot) + \text{TV } f_2^-(0, \cdot) + \text{TV } f_2^+(0, \cdot) + 2C_0 (\text{TV } b_1 + \text{TV } b_2) \\
&\leq 2\mathcal{A}_0.
\end{aligned}$$

As usual,  $Q(U)$ ,  $Q(V)$  stand for interaction potentials between  $\pm 1$ -waves and 0-waves showing up in  $U$ ,  $V$  respectively:

$$Q(U)(t) = \sum_{\beta} \sigma_0^\beta \left[ \sum_{\alpha, x_\alpha < x_\beta} |\sigma_1^\alpha| + \sum_{\alpha, x_\alpha > x_\beta} |\sigma_{-1}^\alpha| \right]$$

where the sum runs over all jumps of  $U$  in  $(x_1 + t, x_2 - t)$ . Hence

$$Q(U)(t) \leq \text{TV } \{b_1\} L_\pm(t; U) \leq \text{TV } \{a\} L_\pm(0+, U) \leq \text{TV } \{a\} \mathcal{A}_0.$$

The situation is analogous for  $V$ . Therefore we estimate the sum of the  $Q$  as follows:

$$Q(U) + Q(V) \leq 2\text{TV } \{a\} \mathcal{A}_0.$$

In order to control the size of these weights, one must manage the bounds:

$$W_{\pm 1}(x) \leq 1 + 2\kappa_1 \text{TV } a, \tag{57}$$

$$W_0(t, x) \leq 1 + 2\mathcal{A}_0 (\kappa_1 + \kappa_2 \text{TV } \{a\}). \tag{58}$$

The constants  $\kappa_1$ ,  $\kappa_2$  still have to be determined. Here we are going to specialize the analysis presented in [8,9] for more general systems and avoid the smallness conditions on the initial data. Let us present the main steps of the analysis:

- (1) show that the functional decreases outside interaction times: see Lemma 3. A natural bound on  $\text{TV } a$  follows and  $\kappa_1$  is suitably chosen, see Remark 7.
- (2) show that the functional decreases at interaction times: see Lemma 4. The constant  $\kappa_2$  is chosen at this step.
- (3) quantify the relation between  $\Phi[U, V](t)$  and the  $L^1$  difference between the two approximate solutions, done in Lemma 5.

The next two lemmas state that  $t \mapsto \Phi[U, V](t)$  decreases both outside interaction times (Lemma 3) and at interaction times (Lemma 4).

**Lemma 3** Let  $U(t, \cdot)$  and  $V(t, \cdot)$  be two approximate solutions generated by the Well-Balanced algorithm, out of the initial data

$$U_0 = (f_1^\pm(t=0, \cdot), a(\cdot)), \quad V_0 = (f_2^\pm(t=0, \cdot), b(\cdot)).$$

Let  $K > 0$  such that the weights  $W_{\pm 1}$  satisfy a uniform bound of the following type:

$$\forall t \geq 0, \quad 1 \leq W_{\pm 1}(t, \cdot) \leq K. \quad (59)$$

and assume that

$$\kappa_1 \geq 2KC_1 \quad (60)$$

with  $C_1$  as in (47). Then, outside interaction times, one has

$$\frac{d\Phi[U, V]}{dt} \leq 0.$$

**Remark 7** From (57), one can choose  $K = 1 + 2\kappa_1 \text{TV } a$  and rewrite (60) as

$$\kappa_1 \geq 2C_1 [1 + 2\kappa_1 \text{TV } a].$$

The above inequality is possible whenever (see (47))

$$4C_1 \text{TV } a < 1, \quad \Leftrightarrow \quad 16\text{TV } a \leq 3 \log(3/2). \quad (61)$$

Therefore, provided that (61) holds, we can operate the following choice:

$$\kappa_1 = \frac{2C_1}{1 - 4C_1 \text{TV } a}, \quad K = \frac{\kappa_1}{2C_1} = \frac{1}{1 - 4C_1 \text{TV } a}. \quad (62)$$

*Proof.* Now we prove Lemma 3. Following Bressan (see [8, p.155]), outside interaction times it is convenient to write the time-derivative of  $\Phi$  as follows:

$$\begin{aligned} \frac{d\Phi[U, V]}{dt} &= \sum_{i=-1}^1 |q_i(x)| W_i(x) (-1 + \lambda_i) \Big|_{x=x_1+t} \\ &\quad + \sum_{i=-1}^1 |q_i(x)| W_i(x) (-1 - \lambda_i) \Big|_{x=x_2-t} + \sum_{\alpha} \sum_{i=-1}^1 E_{\alpha, i}, \end{aligned}$$

being

$$\begin{aligned} E_{\alpha, i} &= |q_i^{\alpha+}| W_i^{\alpha+} (\lambda_i^{\alpha+} - \dot{x}^\alpha) - |q_i^{\alpha-}| W_i^{\alpha-} (\lambda_i^{\alpha-} - \dot{x}^\alpha) \\ &= [ |q_i^{\alpha+}| W_i^{\alpha+} - |q_i^{\alpha-}| W_i^{\alpha-} ] (\lambda_i^\alpha - \dot{x}^\alpha) \end{aligned}$$

where we used that the  $\lambda_i$ 's are constant. Since  $|\lambda_i| \leq 1$ , the contribution from the boundaries is non-positive and then:

$$\frac{d\Phi[U, V]}{dt} \leq \sum_{\alpha} \sum_{i=-1}^1 E_{\alpha, i}.$$

Thanks to the linear structure of families  $\pm 1$ , lots of simplification occur in the sum above. For instance, if  $i = k_\alpha$  then the corresponding speeds coincide,  $\lambda_i^\alpha = \dot{x}^\alpha$ , and thus  $E_{\alpha, i} = 0$ . We shall analyze the jumps that occur in the  $V = (f_2^\pm, b_2)$  vector of unknowns; the analysis for the jumps in  $U$  is completely similar (see also [8, p.160]). Such a framework exactly meets with the interaction estimates given in Prop. 2. Accordingly, let  $k_\alpha \in \{\pm 1, 0\}$  denote the characteristic family of the jump present at the abscissa  $x_\alpha$ . To carry on, one distinguishes between each value of  $k_\alpha$ . For simplicity, in the following we will often omit the dependence on  $\alpha$ .

- If  $k_\alpha = -1 = \dot{x}^\alpha$ , an easy computation shows that  $E_{-1} = 0$  and that

$$E_0 = |q_0^+|W_0^+ - |q_0^-|W_0^-, \quad E_1 = 2[|q_1^+|W_1^+ - |q_1^-|W_1^-].$$

Moreover we have

$$q_0^+ = q_0^-, \quad W_1^+ = W_1^-, \quad W_0^+ - W_0^- = -\kappa_1|\sigma_{-1}|$$

and hence

$$\sum_{i=-1}^1 E_i = E_0 + E_1 = -\kappa_1|\sigma_{-1}||q_0^-| + 2\{|q_1^+| - |q_1^-|\}W_1^-.$$

From (48), Case (a) of Proposition 2, it follows that  $|q_1^+| \leq |q_1^-| + C_1|q_0^-||\sigma_{-1}|$ . Also, recalling (59), the weight  $W_1^-$  is supposed to be smaller than  $K$  and one gets

$$\sum_{i=-1}^1 E_i \leq |q_0^-||\sigma_{-1}|(-\kappa_1 + 2KC_1) \leq 0.$$

- If  $k_\alpha = 1 = \dot{x}^\alpha$ , this is the simple Case (c), and

$$\begin{aligned} \sum_{i=-1}^1 E_i &= E_{-1} + E_0 \\ &= -2\{|q_{-1}^+|W_{-1}^+ - |q_{-1}^-|W_{-1}^-\} - \{|q_0^+|W_0^+ - |q_0^-|W_0^-\}. \end{aligned}$$

Here  $q_0, q_{-1}, W_{-1}$  do not change, while

$$W_0^+ - W_0^- = +\kappa_1\sigma_1.$$

Hence one gets a negative sign for every  $\kappa_1 > 0$ :

$$\sum_{i=-1}^1 E_i = -|q_0|\{W_0^+ - W_0^-\} = -\kappa_1|q_0|\sigma_1 \leq 0.$$

- If  $k_\alpha = 0 = \dot{x}^\alpha$ , this is Case (b), depicted in Fig. 6, with  $\dot{x} = \lambda_0 = 0$  and thus  $E_0 = 0$ .

$$\begin{aligned} \sum_{i=-1}^1 E_i &= E_{-1} + E_1 \\ &= -\{|q_{-1}^+|W_{-1}^+ - |q_{-1}^-|W_{-1}^-\} + \{|q_1^+|W_1^+ - |q_1^-|W_1^-\}. \end{aligned}$$

The weights  $W_i^\pm$ ,  $i = \pm 1$  jump as follows:

$$W_{-1}^+ - W_{-1}^- = +\kappa_1|\sigma_0| \geq 0, \quad W_1^+ - W_1^- = -\kappa_1|\sigma_0|.$$

Hence, by means of (49), we find that

$$\begin{aligned} E_{-1} &= -|q_{-1}^+|\{W_{-1}^+ - W_{-1}^-\} - W_{-1}^-\{|q_{-1}^+| - |q_{-1}^-|\} \\ &\leq -W_{-1}^-\{|q_{-1}^+| - |q_{-1}^-|\} \\ &\leq K|q_{-1}^- - q_{-1}^+| \leq KC_1\sigma_0|q_1^-| \end{aligned}$$

while, in a quite similar way,

$$\begin{aligned} E_1 &= |q_1^-|(W_1^+ - W_1^-) + (|q_1^+| - |q_1^-|)W_1^+ \\ &\leq -\kappa_1\sigma_0|q_1^-| + K|q_1^+ - q_1^-| \\ &\leq -\kappa_1\sigma_0|q_1^-| + KC_1\sigma_0|q_1^-| \\ &\leq \sigma_0|q_1^-|(KC_1 - \kappa_1) \end{aligned}$$



At this point, having  $\kappa_1 \geq 2KC_1$  again ensures  $E_{-1} + E_1 \leq 0$ . □

**Lemma 4** *In the assumptions of Lemma 3, assume that (61) holds and that*

$$\kappa_2 \geq \frac{\kappa_1 C_1}{1 - C_1 \text{TV } a}. \quad (63)$$

Then  $\Phi[U, V](t)$  decreases at interaction times.

*Proof.* Assume that at a certain time  $t$  interactions occur for the approximate solution  $U$ . Recalling the definition (56) of  $\Phi$ , we notice that the  $|q_{\pm 1}(t, x)|$  change continuously in  $L^1_{loc}$ .

The only term that can change in a discontinuous way across the interaction time  $t$  is the weight  $W_0(t, x)$ :

$$\Delta W_0(t, x) = \kappa_1 \Delta A_0(t, x) + \kappa_2 \Delta Q(U)(t)$$

The term  $\Delta A_0$  can increase across interaction times, while  $\Delta Q(U)(t)$  decreases, as follows. For each  $x_\beta$  where a 0-wave is located, let

$$\sigma_0^\beta, \quad \sigma_{-1}^{\beta\pm}, \quad \sigma_1^{\beta\pm}$$

the waves involved in the interaction (with obvious notation). Thanks to Proposition 3, one of the two terms

$$|\sigma_1^{\beta+}| - |\sigma_1^{\beta-}|, \quad |\sigma_{-1}^{\beta+}| - |\sigma_{-1}^{\beta-}|$$

is negative, while the other one is possibly bounded by  $C_1 \sigma_0^\beta (|\sigma_1^{\beta-}| + |\sigma_{-1}^{\beta-}|)$ . Hence,

$$\begin{aligned} \Delta Q &= - \sum_{\beta} \sigma_0^\beta (|\sigma_1^{\beta-}| + |\sigma_{-1}^{\beta-}|) \\ &\quad + \sum_{\beta} (|\sigma_1^{\beta+}| - |\sigma_1^{\beta-}|) \text{TV} \{a; (x_\beta, \infty)\} + \sum_{\beta} (|\sigma_{-1}^{\beta+}| - |\sigma_{-1}^{\beta-}|) \text{TV} \{a; (-\infty, x_\beta)\} \\ &\leq (-1 + C_1 \text{TV } a) \sum_{\beta} \sigma_0^\beta (|\sigma_1^{\beta-}| + |\sigma_{-1}^{\beta-}|). \end{aligned}$$

On the other hand, thanks to (54) in Proposition 3, the possible increase of  $A_0$  is bounded uniformly in  $x$  as follows:

$$\Delta A_0(t, x) \leq C_1 \sum_{\beta} \sigma_0^\beta (|\sigma_1^{\beta-}| + |\sigma_{-1}^{\beta-}|).$$

Therefore

$$\Delta W_0 \leq (\kappa_1 C_1 - \kappa_2 (1 - C_1 \text{TV } a)) \sum_{\beta} \sigma_0^\beta (|\sigma_1^{\beta-}| + |\sigma_{-1}^{\beta-}|).$$

The above quantity is  $\leq 0$  whenever  $1 - C_1 \text{TV } a > 0$ , which is guaranteed by (61), and when  $\kappa_2$  satisfies condition (63). □

**Remark 8** *Following Remark 7, here we summarize the choice of  $\kappa_1$ ,  $\kappa_2$  and the bounds on  $W_i$  obtained so far. Thanks to Lemma 3, we have*

$$W_{\pm 1} \leq 1 + 2\kappa_1 \text{TV } a \leq K = \frac{\kappa_1}{2C_1};$$

this is possible if (61) holds, that is  $4C_1 \text{TV } a < 1$ . Then  $\kappa_1$  can be set as (62). Therefore a bound for  $W_{\pm 1}$  in terms of the data is:

$$W_{\pm 1} \leq \frac{1}{1 - 4C_1 \text{TV } a} = K. \quad (64)$$

Recalling (58) and thanks to Lemma 4, we get

$$\begin{aligned} W_0(t, x) &\leq 1 + 2\mathcal{A}_0 (\kappa_1 + \kappa_2 \text{TV} \{a\}) \leq 1 + 2 \frac{\kappa_1 \mathcal{A}_0}{1 - C_1 \text{TV} \{a\}} \\ &= 1 + \frac{4C_1 \mathcal{A}_0}{(1 - 4C_1 \text{TV } a)(1 - C_1 \text{TV} \{a\})} \doteq K_0. \end{aligned} \quad (65)$$

Now we take advantage of the equivalence of  $\Phi[U, V](t)$  and the  $L^1$  difference between any two approximate solutions.

**Lemma 5** For

$$I(t) = \int_{x_1+t}^{x_2-t} |f_1^+(t, x) - f_2^+(t, x)| + |f_1^-(t, x) - f_2^-(t, x)| dx$$

we get the estimate

$$I(t) \leq K \cdot I(0) + (2KC_0 + K_0) \int_{x_1}^{x_2} |b_1 - b_2| dx + (2C_0 - 1) \int_{x_1+t}^{x_2-t} |b_1 - b_2| dx. \quad (66)$$

**Remark 9** According to (65),  $K_0 = 1 + K \cdot \frac{4C_1 \mathcal{A}_0}{1 - C_1 \text{TV} a}$ , but simultaneously,

$$K = \frac{1}{1 - 4C_1 \text{TV} a}, \quad 1 - C_1 \text{TV} a = \frac{3K + 1}{4K}, \quad C_1 \leq \frac{14}{3}.$$

So,  $K_0 = 1 + \frac{16}{3} \frac{K^2 C_1 \mathcal{A}_0}{3K + 1}$  and for instance, if  $-x_1, x_2 \rightarrow +\infty$ , (66) rewrites,

$$I_{\mathbb{R}}(t) \leq K \left[ I_{\mathbb{R}}(0) + 2 \left( C_0 + \frac{8KC_1 \mathcal{A}_0}{3K + 1} \right) \int_{\mathbb{R}} |b_1 - b_2| dx \right] + 2C_0 \int_{\mathbb{R}} |b_1 - b_2| dx. \quad (67)$$

Notice also that the quantity  $(2C_0 - 1)$  in (66) can be negative.

*Proof.* Recalling (22) and using  $W_{\pm 1} \geq 1$ , one gets

$$\begin{aligned} I(t) &\leq \int_{x_1+t}^{x_2-t} \{|q_1| + |q_{-1}| + |b_1 - b_2|\} dx + (2C_0 - 1) \int_{x_1+t}^{x_2-t} |b_1 - b_2| dx \\ &\leq \Phi[U, V](t) + (2C_0 - 1) \int_{x_1+t}^{x_2-t} |b_1 - b_2| dx \end{aligned}$$

and also, always taking advantage of (22),

$$\begin{aligned} \Phi[U, V](t) &\leq K \sum_{i=-1,1} \int_{x_1+t}^{x_2-t} |q_i| dx + K_0 \int_{x_1+t}^{x_2-t} |b_1 - b_2| dx \\ &\leq KI(t) + (2KC_0 + K_0) \int_{x_1+t}^{x_2-t} |b_1 - b_2| dx. \end{aligned}$$

Altogether, since  $t \mapsto \Phi[U, V](t)$  decreases, it comes that:

$$\begin{aligned} I(t) &\leq \Phi[U, V](t) + (2C_0 - 1) \int_{x_1+t}^{x_2-t} |b_1 - b_2| dx \\ &\leq \Phi[U, V](0) + (2C_0 - 1) \int_{x_1+t}^{x_2-t} |b_1 - b_2| dx \\ &\leq KI(0) + (2KC_0 + K_0) \int_{x_1}^{x_2} |b_1 - b_2| dx + (2C_0 - 1) \int_{x_1+t}^{x_2-t} |b_1 - b_2| dx \end{aligned}$$

which is precisely (66). □

Since we now have the time-decay of  $\Phi[U, V]$  at hand, by just selecting  $\Delta x = (\Delta x)_1$ ,

$$b_1 = P^{(\Delta x)} a, \quad \partial_x a(x) = k(x), \quad (68)$$

$V(t = 0, \cdot) = P^{\Delta x}U(t = 0, \cdot)$  and sending  $(\Delta x)_2 \rightarrow 0$ , one obtains that the  $L^1$  error of the WB scheme at time  $t > 0$  is bounded by

$$\begin{aligned} \int_{x_1+t}^{x_2-t} |f_{\Delta x}^{\pm}(t, x) - f^{\pm}(t, x)| dx &\leq K \int_{x_1}^{x_2} |f_{\Delta x}^{\pm}(0, x) - f^{\pm}(0, x)| dx \\ &+ (2KC_0 + K_0)\Delta x \text{TV}\{a; (x_1, x_2)\} + (2C_0 - 1)\Delta x \text{TV}\{a; (x_1 + t, x_2 - t)\}, \end{aligned}$$

where  $K, K_0$  are given by (64), (65) respectively. Taking advantage of (67), we get that on the whole real line, the global  $L^1$  error is bounded uniformly in time by the quantity,

$$\boxed{\frac{1}{\Delta x} \int_{\mathbb{R}} |f_{\Delta x}^{\pm}(t, x) - f^{\pm}(t, x)| dx \leq K \text{TV}(f^{\pm}(0, \cdot)) + 2 \left[ (K + 1)C_0 + \frac{8K^2 C_1 \mathcal{A}_0}{3K + 1} \right] \|k\|_{L^1(\mathbb{R})}}$$

which blows up as  $C_1 \|k\|_{L^1(\mathbb{R})} \rightarrow \frac{1}{4}$ , since the constant  $K$  does (see (62)). This was to be expected, as for stiff relaxation regimes and well-prepared initial data, one expects  $\rho = f^+ + f^-$  to match the entropy solution of the conservation law  $\partial_t \rho + \partial_x A(\rho) = 0$ , and one cannot have order 1 convergence as  $\Delta x \rightarrow 0$ . This completes the proof of the first estimate,  $\mathcal{E}_1$ , in Theorem 1.

#### 4 Complementary $L^1$ error estimate through entropy dissipation

The former error estimate suits well the non-stiff case for (1). However, one may feel the need for a study of the complementary situation, where typically  $|k(x)|\Delta x$  can become (locally) big. In order to quantify the  $L^1$  error of WB schemes in this context too, we adapt a method of [28] (see also [22]) based on entropy dissipation and inspired by the seminal ideas of Kuznetsov [25].

##### 4.1 Quasi-monotonicity and entropy inequalities

Let us first describe what type of entropy inequalities are satisfied by the exact solution and by the WB approximation. On one hand, the exact solution of (8) is such that, for any constant values  $k_{\pm} \in \mathbb{R}^2$  and any test-function  $0 \leq \varphi(t, x) \in C_0^{\infty}(\mathbb{R}^+ \times \mathbb{R})$ ,

$$\begin{aligned} & - \int_0^T \int_{\mathbb{R}} (|f^+ - k_+| + |f^- - k_-|) \partial_t \varphi + (|f^+ - k_+| - |f^- - k_-|) \partial_x \varphi \cdot dx \cdot dt \\ & \quad + \int_{\mathbb{R}} (|f^+(T, x) - k_+| + |f^-(T, x) - k_-|) \varphi(T, x) \cdot dx \\ & \quad - \int_{\mathbb{R}} (|f^+(0, x) - k_+| + |f^-(0, x) - k_-|) \varphi(0, x) \cdot dx \tag{69} \\ & \leq \int_0^T \int_{\mathbb{R}} k(x) (\text{sgn}(f^+ - k_+) - \text{sgn}(f^- - k_-)) G(f^+, f^-) \varphi \cdot dx \cdot dt. \end{aligned}$$

On the other hand, the WB approximation is the exact solution of (12) with piecewise-constant initial data fitted to the length separating 2 zero-waves (see again Fig. 3), in particular there is no projection at each time-step. We have the following Lemma.

**Lemma 6** For any test-function  $\varphi(t, x) \geq 0$  compactly supported on  $(0, T) \times \mathbb{R}$ , one has

$$\begin{aligned}
& - \int_0^T \int_{\mathbb{R}} (|f^+ - k_+| + |f^- - k_-|) \partial_t \varphi + (|f^+ - k_+| - |f^- - k_-|) \partial_x \varphi \cdot dx \cdot dt \\
& \quad - \int_0^T \int_{\mathbb{R}} k(x) (\operatorname{sgn}(f^+ - k_+) - \operatorname{sgn}(f^- - k_-)) G(f^+, f^-) \varphi \cdot dx \cdot dt \\
& \leq C_\alpha \sum_{n,j} \operatorname{TV} (f^\pm(t^n, \cdot); \{x_{j-1}, x_j\}) \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \varphi(t, x_{j-\frac{1}{2}}) dx \cdot dt \\
& \quad + C_\beta \sum_{n,j} \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \left| \varphi(t, x) - \varphi(t, x_{j-\frac{1}{2}}) \right| dx \cdot dt, \tag{70}
\end{aligned}$$

where  $C_\alpha = \operatorname{Lip}(G)$  and  $C_\beta = 2C_0$ , the Maxwellian gap defined in (23).

*Proof.* The proof is divided into several steps.

- Using the standard notation,  $t^n = n\Delta t$ ,  $C_j = (x_{j-\frac{1}{2}}, x_{j+\frac{1}{2}})$  with  $x_{j-\frac{1}{2}} = (j - \frac{1}{2})\Delta x$  the locus of the zero-waves, comes in each “cell”  $C_j \times (t^n, t^{n+1})$ ,

$$\begin{aligned}
& - \int_{t^n}^{t^{n+1}} \int_{C_j} (\eta_+(f^+) + \eta_-(f^-)) \partial_t \varphi + (\eta_+(f^+) - \eta_-(f^-)) \partial_x \varphi \cdot dx \cdot dt \\
& \quad + \int_{C_j} (\eta_+(f^+(t^{n+1}, x)) + \eta_-(f^-(t^{n+1}, x))) \varphi(t^{n+1}, x) \cdot dx \\
& \quad - \int_{C_j} (\eta_+(f^+(t^n, x)) + \eta_-(f^-(t^n, x))) \varphi(t^n, x) \cdot dx \\
& \quad + \int_{t^n}^{t^{n+1}} [(\eta_+(f^+) - \eta_-(f^-)) \varphi] (t, x_{j+\frac{1}{2}} - 0) \cdot dt \\
& \quad - \int_{t^n}^{t^{n+1}} [(\eta_+(f^+) - \eta_-(f^-)) \varphi] (t, x_{j-\frac{1}{2}} + 0) \cdot dt \leq 0
\end{aligned}$$

for any couple of smooth, convex functions  $\eta_\pm \in C^2(\mathbb{R})$  and  $j, n \in \mathbb{Z} \times \mathbb{N}$ . Clearly, as the Courant number is 1, there is no need for a projection step so the summation on  $j, n$  is rather straightforward:

$$\begin{aligned}
& - \sum_{j,n \in \mathbb{Z} \times \mathbb{N}} \int_{t^n}^{t^{n+1}} \int_{C_j} (\eta_+(f^+) + \eta_-(f^-)) \partial_t \varphi + (\eta_+(f^+) - \eta_-(f^-)) \partial_x \varphi \cdot dx \cdot dt \\
& \leq \sum_{j \in \mathbb{Z}, n \in \mathbb{N}} \left( \mathcal{I}_{n,j-\frac{1}{2}}^+ - \mathcal{I}_{n,j-\frac{1}{2}}^- \right) \int_{t^n}^{t^{n+1}} \varphi(t, x_{j-\frac{1}{2}}) \cdot dt, \tag{71}
\end{aligned}$$

because  $\varphi(t, \cdot)$  is continuous in  $x = x_{j-\frac{1}{2}}$  and  $\varphi(t, \cdot) = 0$  for  $t = 0, T$ . We used the following notation,

$$\mathcal{I}_{n,j-\frac{1}{2}}^\pm = \eta_\pm \left( f^\pm(t^n, x_{j-\frac{1}{2}} + 0) \right) - \eta_\pm \left( f^\pm(t^n, x_{j-\frac{1}{2}} - 0) \right).$$

These terms  $\mathcal{I}_{n,j-\frac{1}{2}}^\pm$  stand for the jump of entropy flux across each zero-wave, located at the grid’s interface. They are independent of  $t$  thanks to the CFL condition, which ensures that linear waves propagate exactly  $\Delta x$  during  $\Delta t$ .

- One needs to recover, up to  $\Delta x$ , the source term which appears in the entropy inequality for the exact solution, and which seems to be missing here. By definition of the stationary equations, see (14), at any time-step  $t^n = n\Delta t$ , the corresponding smooth profiles  $\tilde{f}_n^\pm$  satisfy modified ODE’s too,

$$\partial_x \left( \eta_\pm(\tilde{f}_n^\pm) \right) = k(x) G^\pm(\tilde{f}_n^+, \tilde{f}_n^-), \quad G^\pm(\tilde{f}^+, \tilde{f}^-) := \eta'_\pm(\tilde{f}^\pm) G(\tilde{f}^+, \tilde{f}^-).$$

Accordingly, the entropy jumps rewrite:

$$\eta_{\pm}(f^{\pm})(t^n, x_{j-\frac{1}{2}} + 0) = \eta_{\pm}(f^{\pm})(t^n, x_{j-\frac{1}{2}} - 0) + \int_{x_{j-1}}^{x_j} k(s) G^{\pm}(\tilde{f}_n^+(s), \tilde{f}_n^-(s)) \cdot ds,$$

therefore, the former jumps are amended as follows,

$$\mathcal{I}_{n,j-\frac{1}{2}}^+ - \mathcal{I}_{n,j-\frac{1}{2}}^- = \int_{x_{j-1}}^{x_j} k(s) \left( \eta'_+(\tilde{f}_n^+(s)) - \eta'_-(\tilde{f}_n^-(s)) \right) G(\tilde{f}_n^+(s), \tilde{f}_n^-(s)) \cdot ds.$$

So the contribution of the source term can be reconstructed:

$$\begin{aligned} & \left( \mathcal{I}_{n,j-\frac{1}{2}}^+ - \mathcal{I}_{n,j-\frac{1}{2}}^- \right) \int_{t^n}^{t^{n+1}} \varphi(t, x_{j-\frac{1}{2}}) dt \\ &= \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \left[ \eta'_+(f^+) - \eta'_-(f^-) \right] G(f^+, f^-) \varphi(t, x) dx \cdot dt \\ & \quad - \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \underbrace{\left[ \eta'_+(f^+) - \eta'_-(f^-) \right] G(f^+, f^-)}_{G^+(f^{\pm}) - G^-(f^{\pm}) = \beta(t,x)} \left[ \varphi(t, x) - \varphi(t, x_{j-\frac{1}{2}}) \right] dx \cdot dt \\ & \quad - \int_{t^n}^{t^{n+1}} \varphi(t, x_{j-\frac{1}{2}}) \int_{x_{j-1}}^{x_j} k(x) \left[ \left( \eta'_+(f^+) - \eta'_-(f^-) \right) G(f^+, f^-) \right. \\ & \quad \quad \left. - \underbrace{\left( \eta'_+(\tilde{f}_n^+(s)) - \eta'_-(\tilde{f}_n^-(s)) \right) G(\tilde{f}_n^+(x), \tilde{f}_n^-(x))}_{G^+(f^{\pm}) - G^+(\tilde{f}_n^{\pm}) - G^-(f^{\pm}) + G^-(\tilde{f}_n^{\pm}) = \alpha(t,x)} \right] dx \cdot dt. \end{aligned} \tag{72}$$

The above terms  $\alpha$ ,  $\beta$  are bounded as follows:

$$|\alpha(t, x)| \leq Lip(G) TV \left( \tilde{f}_n^{\pm}(\cdot); \{x_{j-1}, x_j\} \right), \quad |\beta(t, x)| \leq |G^+ - G^-| \leq 2C_0.$$

- For any  $\ell \in \mathbb{R}$ , we approximate a weak Kruřkov entropy  $u \mapsto |u - \ell|$  by means of a smooth function  $E \in C^2(\mathbb{R})$  such that  $E'' \geq 0$ ,  $E(v) = |v|$  for  $|v| \geq 1$ ,  $E'(0) = 0$  and  $|E'| \leq 1$ . It is rescaled like  $\eta_{\delta}(v) = \delta E\left(\frac{v-\ell}{\delta}\right)$ , and therefore  $\eta'_{\delta}(v) \rightarrow \text{sgn}(v - \ell)$  as  $\delta \rightarrow 0$ , for all  $v \neq 0$ . Using (71), (72) and thanks to the bound above on  $\alpha$  and  $\beta$ , we pass to the limit as  $\delta \rightarrow 0$  by means of the dominated convergence theorem and finally recover (70). □

#### 4.2 Derivation of the complementary $L^1$ error estimate

Hereafter we shall denote  $f_{\Delta x}^{\pm}$  the piecewise-constant numerical approximations delivered by the WB algorithm described in the former sections, and keep  $f^{\pm}$  for the corresponding exact solution. Each one satisfies a specific entropy dissipation inequality, (69) and (70). An error estimate can be derived by taking advantage of the simple fact that (weak) *Kruřkov entropies are symmetric*, together with a specific choice of nonnegative test-functions. Indeed, adopting the notations of [7,16,22], let us consider,

$$\mathbb{R}^+ \times \mathbb{R} \times \mathbb{R}^+ \times \mathbb{R} \rightarrow \mathbb{R}^+, \quad 0 \leq \phi(t, x, s, y) = \varphi(t, x) \zeta(t - s, x - y).$$

The choice of  $\zeta$  corresponds to a smooth approximation of the Dirac mass, namely for  $\Delta, \delta > 0$ :

$$\zeta(t, x) = \zeta_t(t) \zeta_x(x) = \frac{1}{\delta} \zeta_t^1\left(\frac{t}{\delta}\right) \cdot \frac{1}{\Delta} \zeta_x^1\left(\frac{x}{\Delta}\right), \quad 0 \leq \zeta_t^1, \zeta_x^1 \in C_0^{\infty}(\mathbb{R}).$$

Moreover, one can ensure that they are symmetric and:

$$\|\zeta_t\|_{L^1(\mathbb{R})} = \|\zeta_x\|_{L^1(\mathbb{R})} = 1, \quad \zeta_t^1(\cdot) \zeta_x^1(\cdot) \text{ supported in } (-1, 0) \times \left(-\frac{1}{4}, \frac{1}{4}\right).$$

Now, thanks to entropies' symmetry, it is possible to consider (70) with  $k_{\pm} = f^{\pm}(s, y)$ , for any  $s, y \in \mathbb{R}^+ \times \mathbb{R}$  and reciprocally. By double integration, and usual simplifications, one arrives at:

$$\begin{aligned}
0 \leq & \int \int \int \int \zeta(t-s, x-y) \left\{ [|f_{\Delta x}^+(t, x) - f^+(s, y)| + |f_{\Delta x}^-(t, x) - f^-(s, y)|] \partial_t \varphi(t, x) \right. \\
& + [|f_{\Delta x}^+(t, x) - f^+(s, y)| - |f_{\Delta x}^-(t, x) - f^-(s, y)|] \partial_x \varphi(t, x) \\
& + (\operatorname{sgn}(f_{\Delta x}^+(t, x) - f^+(s, y)) - \operatorname{sgn}(f_{\Delta x}^-(t, x) - f^-(s, y))) \times \\
& \left. [k(x)G(f_{\Delta x}^+, f_{\Delta x}^-)(t, x) - k(y)G(f^+, f^-)(s, y)] \varphi(t, x) \right\} ds dy dt dx \quad (73) \\
& + C_{\alpha} \sum_{n,j} \operatorname{TV} \left( \tilde{f}_n^{\pm}(\cdot); \{x_{j-1}, x_j\} \right) \int dy \int ds \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \phi(t, x_{j-\frac{1}{2}}, s, y) dt dx \\
& + 2C_0 \sum_{n,j} \int dy \int ds \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \left| \phi(t, x, s, y) - \phi(t, x_{j-\frac{1}{2}}, s, y) \right| dt dx.
\end{aligned}$$

By imposing that  $\varphi(t, x)$  is a regularized characteristic function as in [7,22] with  $\nu = 0$ ,  $\delta = \Delta$ ,  $L = 1$  and  $\theta = \Delta/4$ , space and time derivatives simplify each other in order to produce

$$\begin{aligned}
& \int_{x_1 - \frac{\Delta}{2}}^{x_2 + \frac{\Delta}{2}} |f_{\Delta x}^{\pm}(T, x) - f^{\pm}(T, x)| dx \\
& \leq \int_{x_1 - \frac{\Delta}{2} - T}^{x_2 + \frac{\Delta}{2} + T} |f_{\Delta x}^{\pm}(0, x) - f^{\pm}(0, x)| dx + 4C \operatorname{TV}(f^{\pm}(0, \cdot)) \Delta + [\dots].
\end{aligned}$$

Now, in contrast with the similar computation in [1], one can get rid of the contribution of  $G$  in the term (73) by taking advantage of its quasi-monotonicity: in fact, since  $\pm \frac{\partial G}{\partial f^{\pm}} \leq 0$  (see (16)) and  $\operatorname{sgn}(b)a - |a| \leq 0$  for any  $a, b \in \mathbb{R}^2$ , we have

$$[\operatorname{sgn}(f_{\Delta x}^+(t, x) - f^+(s, y)) - \operatorname{sgn}(f_{\Delta x}^-(t, x) - f^-(s, y))] [G(f_{\Delta x}^+, f_{\Delta x}^-)(t, x) - G(f^+, f^-)(s, y)] \leq 0.$$

Since  $k(x) \geq 0$ , from the integrand in (73) we get a negative term, while the remaining term comes from the difference  $k(x) - k(y)$  and is smaller than:

$$\begin{aligned}
& \int \int \int \int \varphi(t, x) |k(x) - k(y)| \zeta(t-s, x-y) ds dy dt dx \\
& \leq T \int_x \int_y |k(x) - k(y)| \zeta_x(x-y) dx dy \\
& \leq \frac{T}{\Delta} \int_x \int_{-\frac{\Delta}{4}}^{\frac{\Delta}{4}} |k(x) - k(x+\xi)| d\xi dx \\
& \leq \frac{T}{\Delta} \operatorname{TV}(k) \int_{-\frac{\Delta}{4}}^{\frac{\Delta}{4}} |\xi| d\xi = \operatorname{TV}(k) \cdot \frac{\Delta}{16} \cdot T \sup |\varphi|.
\end{aligned}$$

Above, we used that  $|\varphi| \leq 1$  by construction. It is necessary to derive suitable bounds for the error terms:

- following the construction of [7],  $|\partial_x \varphi| \leq C/\Delta$  and this affects the term:

$$\begin{aligned}
& \int dy \int ds \sum_{n,j} \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \left| \phi(t, x, s, y) - \phi(t, x_{j-\frac{1}{2}}, s, y) \right| dt dx \\
& \leq \sum_{n,j} \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \left| \varphi(t, x) - \varphi(t, x_{j-\frac{1}{2}}) \right| dt dx \leq CT \frac{\Delta x}{\Delta} \|k\|_{L^1(\mathbb{R})}.
\end{aligned}$$

- The other term depends on  $\text{TV}(\tilde{f}_n^\pm; x_{j-1}, x_j)$ , which is bounded by  $C_0 \Delta x \|k\|_{L^\infty(\mathbb{R})}$ , so one gets:

$$\sum_{n,j} \text{TV} \left\{ \tilde{f}_n^\pm(\cdot); (x_{j-1}, x_j) \right\} \int_{t^n}^{t^{n+1}} \int_{x_{j-1}}^{x_j} k(x) \varphi(t, x_{j-\frac{1}{2}}) dt dx \leq T \cdot C_0 \Delta x \|k\|_{L^\infty} \|k\|_{L^1}.$$

Since  $C_\alpha \leq \text{Lip}(G) \leq 2$  (within the assumptions (6), (5)) the inequality reduces to:

$$\begin{aligned} \int_{x_1 - \frac{\Delta}{2}}^{x_2 + \frac{\Delta}{2}} |f_{\Delta x}^\pm(T, x) - f^\pm(T, x)| \cdot dx &\leq \int_{x_1 - \frac{\Delta}{2} - T}^{x_2 + \frac{\Delta}{2} + T} |f_{\Delta x}^\pm(0, x) - f^\pm(0, x)| \cdot dx \\ &+ \frac{4CT \Delta x \|k\|_{L^1}}{\Delta} + 2C_0(T \Delta x) \|k\|_{L^1} \|k\|_{L^\infty} \\ &+ \Delta [4C \text{TV}(f^\pm(0, \cdot)) + C_0 \text{TV}(k) T / 8]. \end{aligned}$$

The optimal value for  $\Delta$  can be computed by standard ways, and one finds:

$$\begin{aligned} \int_{x_1 - \frac{\Delta}{2}}^{x_2 + \frac{\Delta}{2}} |f_{\Delta x}^\pm(T, x) - f^\pm(T, x)| \cdot dx &\leq \int_{x_1 - \frac{\Delta}{2} - T}^{x_2 + \frac{\Delta}{2} + T} |f_{\Delta x}^\pm(0, x) - f^\pm(0, x)| \cdot dx \\ + 2T \left\{ 2 \sqrt{2 \Delta x C_0 \|k\|_{L^1} \left[ \frac{4}{T \cdot C_0} \text{TV}(f^\pm(0, \cdot)) + \frac{\text{TV}(k)}{8} \right]} + \Delta x C_0 \|k\|_{L^1} \|k\|_{L^\infty} \right\}. \end{aligned}$$

The absolute constant  $C$  which is used in [7] is fixed here at 2, based on [18, Theorem 2]. We have established the second estimate,  $\mathcal{E}_2$ : the proof of Theorem 1 is yet complete.

## A An elementary example

The implicit flux-function  $\tilde{J}(\delta, f^\pm)$ , as derived in Proposition 1, is not very convenient for computing sharp interaction estimates; in particular, we have the usual derivation rule,

$$\frac{\partial \tilde{J}}{\partial \delta} = \frac{1}{\frac{\partial F}{\partial J}(\tilde{J}, \delta, f^\pm)}, \quad \frac{\partial \tilde{J}}{\partial f^\pm} = - \frac{\frac{\partial F}{\partial f^\pm}(\tilde{J}, \delta, f^\pm)}{\frac{\partial F}{\partial J}(\tilde{J}, \delta, f^\pm)}.$$

Accordingly, mixed derivatives of  $\tilde{J}$  have an intricate expression:

$$\frac{\partial^2 \tilde{J}}{\partial f^\pm \partial \delta}(\delta, f^\pm) = - \frac{\left( \frac{\partial^2 F}{\partial f^\pm \partial J}(\tilde{J}, \delta, f^\pm) + \frac{\partial \tilde{J}}{\partial f^\pm}(\delta, f^\pm) \cdot \frac{\partial^2 F}{\partial J^2}(\tilde{J}, \delta, f^\pm) \right)}{\left( \frac{\partial F}{\partial J}(\tilde{J}, \delta, f^\pm) \right)^2}.$$

Assume that the relaxation is just  $g(\rho, J) = \alpha \rho - J$ ,  $0 \leq \alpha < 1$ . An elementary computation yields that

$$B(\rho, J) = \frac{1}{2\alpha} \log \left| \rho - \frac{J}{\alpha} \right|,$$

and (24) rewrites as

$$F(J, \delta, f^\pm) = B(2f^+ - J, J) - B(2f^- + J, J) - \delta = \frac{1}{2\alpha} \log \left| \frac{2f^+ - J(1 + \frac{1}{\alpha})}{2f^- + J(1 - \frac{1}{\alpha})} \right| - \delta.$$

The function  $\tilde{J}$  is, for this simple case,

$$\tilde{J}(\delta, f^\pm) = \frac{2\alpha(f^+ - f^- \exp(2\alpha\delta))}{1 + \alpha - \exp(2\alpha\delta)(1 - \alpha)},$$

so its partial derivative in  $\delta$  reads:

$$\frac{\partial \tilde{J}}{\partial \delta} = 4\alpha^2 \exp(2\alpha\delta) \frac{f^+(1 - \alpha) - f^-(1 + \alpha)}{(1 + \alpha - \exp(2\alpha\delta)(1 - \alpha))^2} = \frac{4\alpha^2 \exp(2\alpha\delta)}{(1 + \alpha - \exp(2\alpha\delta)(1 - \alpha))^2} (J - \alpha\rho),$$

which clearly changes sign when the equilibrium curve is crossed. Its partial derivatives in  $f^\pm$  read,

$$\frac{\partial \tilde{J}}{\partial f^+} = \frac{2\alpha}{1 + \alpha - \exp(2\alpha\delta)(1 - \alpha)}, \quad \frac{\partial \tilde{J}}{\partial f^-} = \frac{-2\alpha}{(1 + \alpha) \exp(-2\alpha\delta) - (1 - \alpha)},$$

so  $\frac{\partial \tilde{J}}{\partial f^\pm} \simeq \frac{\pm 1}{1 - \delta(1 \mp \alpha)}$  for small  $\delta > 0$ . Consequently, we get second-order (mixed) derivatives as follows:

$$\begin{aligned} \frac{\partial^2 \tilde{J}}{\partial f^+ \partial \delta}(\delta, f^\pm) &= 4\alpha^2 \frac{(1 - \alpha) \exp(2\alpha\delta)}{[1 + \alpha - \exp(2\alpha\delta)(1 - \alpha)]^2} \geq 0, \\ \frac{\partial^2 \tilde{J}}{\partial f^- \partial \delta}(\delta, f^\pm) &= -4\alpha^2 \frac{(1 + \alpha) \exp(-2\alpha\delta)}{[(1 + \alpha) \exp(-2\alpha\delta) - (1 - \alpha)]^2} \leq 0. \end{aligned} \tag{A.1}$$

## References

- [1] D. Amadori, L. Gosse. Transient  $L^1$  error estimates for well-balanced schemes on non-resonant scalar balance laws, *J. Differential Equations* **255** (2013), 469–502
- [2] D. Amadori, L. Gosse. Error estimates for well-balanced and time-split schemes on a damped semilinear wave equation, Preprint (2013)
- [3] D. Amadori, G. Guerra. Global BV solutions and relaxation limit for a system of conservation laws, *Proc. Roy. Soc. Edinburgh Sect. A* **131** (2001), 1–26
- [4] S. Bianchini. A Glimm-type functional for a special Jin–Xin relaxation model, *Ann. Inst. H. Poincaré Anal. Non-linéaire* **18**(1) (2001), 19–42
- [5] S. Bianchini. Relaxation limit of the Jin–Xin relaxation model, *Comm. Pure Appl. Math.* **59**(5) (2006), 688–753
- [6] S. Bianchini, B. Hanouzet, R. Natalini. Asymptotic behavior of smooth solutions for partially dissipative hyperbolic systems with a convex entropy, *Comm. Pure Appl. Math.* **60** (2007), 1559–1622
- [7] F. Bouchut and B. Perthame. Kružkov’s estimates for scalar conservation laws revisited, *Trans. Amer. Math. Soc.* **350** (1998), no. 7, 2847–2870
- [8] A. Bressan. **Hyperbolic Systems of Conservation Laws. The one-dimensional Cauchy problem**, Oxford Lecture Series in Mathematics and its Applications **20**, Oxford University Press, Oxford, 2000
- [9] A. Bressan, T.-P. Liu, T. Yang.  $L^1$  stability estimates for  $n \times n$  conservation laws. *Arch. Ration. Mech. Anal.* **149** (1999), 1–22



- [10] A. Edwards, B. Perthame, N. Seguin, M. Tournus. Analysis of a simplified model of the urine concentration mechanism, *Netw. Heterog. Media* **7** (2012), 989–1018
- [11] B. Cockburn, H. Gau. A posteriori error estimates for general numerical schemes for conservations laws, *Mat. Apl. Comput.* **14** (1995), 37–47
- [12] F. Coquel, S. Jin, J.-G. Liu, L. Wang. Well-posedness and singular limit of a semilinear hyperbolic relaxation system with a two-scale discontinuous relaxation rate. Preprint (2013)
- [13] W. E. Homogenization of scalar conservation laws with oscillatory forcing terms, *SIAM J. Appl. Math* **52** (1992), 959–972
- [14] J. Glimm, D.H. Sharp. An  $S$ -matrix theory for classical nonlinear physics, *Found. Phys.* **16** (1986), 125–141
- [15] L. Gosse, *Time-splitting schemes and measure source terms for a quasilinear relaxing system*, M3AS **13**(8) (2003), 1081–1101
- [16] L. Gosse. **Computing Qualitatively Correct Approximations of Balance Laws**, SIMAI Springer Series, Vol. 2, Springer (2013)
- [17] L. Gosse, Redheffer products and numerical approximation of currents in one-dimensional semiconductor kinetic models, *SIAM MMS* (2014), to appear
- [18] L. Gosse, Ch. Makridakis. Two a Posteriori Error Estimates for One-Dimensional Scalar Conservation Laws, *SIAM J. Numer. Anal.* **30** (2000), 964–988
- [19] L. Gosse, G. Toscani. Space localization and well-balanced schemes for discrete kinetic models in diffusive regimes, *SIAM J. Numer. Anal.* **41** (2004), 641–658
- [20] L. Gosse, A. Tzavaras. Convergence of relaxation schemes to the equations of elastodynamics, *Math. Comp.* **70**(234) (2001), 555–577
- [21] S. Jin, Z. Xin. The relaxation schemes for systems of conservation laws in arbitrary space dimension, *Comm. Pure Appl. Math.* **48** (1995), 235–276
- [22] M.A. Katsoulakis, G. Kossioris, Ch. Makridakis. Convergence and error estimates of relaxation schemes for multidimensional conservation laws, *Comm. Partial Differential Equations* **24** (3-4) (1999), 395–422
- [23] M.A. Katsoulakis, A.E. Tzavaras. Contractive relaxation systems and the scalar multidimensional conservation law, *Comm. Partial Differential Equations* **22** (1997), 195–233
- [24] S.N. Kružkov. First order quasilinear equations in several independant space variables, *Mat. USSR Sbornik* **81** (1970), 228–255
- [25] N.N. Kuznetsov. Accuracy of some approximate methods for computing the weak solutions of a first-order quasilinear equation, *Zh. Vychisl. Mat. i Mat. Fiz.*, 16 (1976), pp. 1489–1502; English transl. in *USSR Comp. Math. and Math. Phys.*, 16 (1976), pp. 105–119
- [26] M. Laforest. A Posteriori Error Estimate for Front-Tracking: System of Conservation Laws, *SIAM J. Math. Anal.* **35** (2004), 1347–1370
- [27] H. Liu, R. Natalini. Long-Time Diffusive Behavior of Solutions to a Hyperbolic Relaxation System, *Asympt. Analysis* **25** (2001), 21–38

- [28] H.L. Liu, G. Warnecke. Convergence rates for relaxation schemes approximating conservation laws, *SIAM J. Numer. Anal.* **37** (4) (2000), 1316–1337
- [29] C. Mascia, K. Zumbrun. Pointwise Green’s function bounds and stability of relaxation shocks. *Indiana Univ. Math. Journal* **51**(4) (2001), 773–904
- [30] R. Natalini. Recent results on hyperbolic relaxation problems, in **Analysis of Systems of Conservation Laws** (Aachen, 1997), Monogr. Surv. Pure Appl. Math. **99** (Chapman & Hall/CRC, 1999)
- [31] B. Perthame, N. Seguin, M. Tournus, A simple derivation of BV bounds for inhomogeneous relaxation systems, Preprint ArXiv (1014)
- [32] F. Sabac, *The optimal convergence rate of monotone finite difference methods for hyperbolic conservation laws*, *SIAM J. Numer. Anal.* **34** (1997), 2306–2318
- [33] D. Serre. Relaxation semi-linéaire et cinétique des systèmes de lois de conservation, *Ann. Inst. Henri Poincaré, Analyse non linéaire* **17**(2) (2000), 169–192
- [34] E. Tadmor, T. Tang. Pointwise error estimates for relaxation approximations to conservation laws. *SIAM J. Math. Anal.* **32**(4) (2000), 870–886
- [35] Z.H. Teng. First-order  $L^1$  convergence for relaxation approximations to conservation laws. *Comm. Pure Appl. Math.* **51**(8) (1998), 857–895