

# On Convergence of Minmod-Type Schemes

Sergei Konyagin\*, Bojan Popov† and Ognian Trifonov‡

April 16, 2004

## Abstract

A class of non-oscillatory numerical methods for solving nonlinear scalar conservation laws in one space dimension is considered. This class of methods contains the classical Lax-Friedrichs and the second order Nessyahu-Tadmor scheme. In the case of linear flux, new  $l_2$  stability results and error estimates for the methods are proved. Numerical experiments confirm that these methods are one-sided  $l_2$  stable for convex flux instead of the usual Lip+ stability.

## 1 Introduction

We are interested in the scalar hyperbolic conservation law

$$(1) \quad \begin{cases} u_t + f(u)_x = 0, & (x, t) \in \mathbb{R} \times (0, \infty), \\ u(x, 0) = u^0(x), & x \in \mathbb{R}, \end{cases}$$

where  $f$  is a given flux function. In recent years, there has been enormous activity in the development of the mathematical theory and in the construction of numerical methods for (1). Even though the existence-uniqueness theory of weak solutions is complete [12], there are many numerically efficient methods for which the questions of convergence and error estimates are still open. For example, there are many non-oscillatory schemes based on the minmod limiter which are numerically robust, at least in many numerical tests, but theoretical results about convergence and error estimates are still missing [3, 6, 7, 18].

In this paper, we consider a class of the so-called Godunov-type schemes for solving (1). There are two main steps in such schemes: evolution and projection. In the original Godunov scheme, the projection is onto piecewise constant functions – the cell averages. In the general Godunov-type method, the projection is onto piecewise polynomials. To determine the properties of these schemes it is necessary to study the

---

\*Department of Mathematics, Moscow State University, Moscow, Russia, [konyagin@ok.ru](mailto:konyagin@ok.ru) Supported by the the grant NSh - 304.2003.1.

†Department of Mathematics, Texas A&M University, College Station, TX 77845, USA, [popov@math.tamu.edu](mailto:popov@math.tamu.edu). Supported by the ONR Grant No. N00014-91-J-1076.

‡Department of Mathematics, University of South Carolina, Columbia, SC 29208, USA, [trifonov@math.sc.edu](mailto:trifonov@math.sc.edu). Supported by the NSF DMS Grant No. 9970455.

properties of the projection operator. We limit our attention to the case of piecewise linear projection based on cell averages using minmod limiters for the slope reconstruction and we call such a scheme *minmod-type*. For example, the Nessyahu-Tadmor scheme [15] is of minmod-type and it is based on staggered evolution, other examples include the second order non-oscillatory central schemes with non-staggered grids given in [8, 9], and the UNO and TVD2 schemes in [6]. Theoretical results about convergence of such schemes to the entropy solution, or error estimates, are still missing. In most cases the authors give a variation bound for such a scheme which is enough to conclude that the method converges to a weak solution, see [10]. The only paper which has a convergence result is the Nessyahu-Tadmor paper [15] in which the authors prove a single cell entropy inequality for a minor modification of the original MinMod scheme. A single entropy inequality is enough to conclude that the scheme is convergent to the unique entropy solution but does not give any rate of convergence. In order to get a rate, one has to have a family of entropy inequalities (see [1, 2, 11, 14]). Alternatively, for a convex flux, one can impose Lip+ stability on the projection and then prove convergence via Tadmor's Lip' theory [16, 19]. Unfortunately, it is well known that minmod-type schemes are incompatible with the Lip+ condition – the Lip+ seminorm is not preserved by a minmod-type projection. It is easy to think about minmod-type schemes in terms of new/old cell averages. That is, we start with a sequence of cell averages  $\{w_j\}$  and after one time step (projection and evolution) we get a new sequence  $\{w'_j\}$ . A scheme is *total variation diminishing* (TVD) if the variation of the new sequence  $\sum_j |w'_j - w'_{j-1}|$  is not bigger than the variation of the old one  $\sum_j |w_j - w_{j-1}|$ , i.e., the  $l_1$  norm of the jumps does not increase in time. In the Lip+ case (for convex flux) the condition on the jumps is that the biggest non-negative jump does not increase in time

$$\sup_j (w'_j - w'_{j-1})_+ \leq \sup_j (w_j - w_{j-1})_+.$$

In Section 3 of this paper, we prove that for linear flux the  $l_2$  norm of the jumps for some minmod-type schemes does not increase in time. This class of schemes include the NT scheme and the TVD2 scheme considered in [6]. Based on that, we use the dual approach (see [16, 19]) to derive a new error estimate in  $L_2$  in Section 4. The rate of convergence we prove is 1/2 in  $L_2$  which improves the known result of 1/4 (see [19]). In Section 5, we present numerical examples in the case of linear and convex flux, and discuss the non-convex case. Our numerical tests show that for convex flux the minmod schemes preserve the one-sided analog

$$\sum_j (w'_j - w'_{j-1})_+^2 \leq \sum_j (w_j - w_{j-1})_+^2$$

which suggests a different approach to prove convergence and error estimates for such schemes in the convex case. The  $l_2$  norm of the jumps is a natural candidate norm for the analysis of high-order schemes, such as central or ENO [7] type, due to its numerical viscosity. We view the results of this paper as a step toward obtaining convergence results and estimates for the rate of convergence of minmod-type schemes for solving (1) in the case of convex nonlinear flux.

## 2 Non-Oscillatory Central Schemes

In this section, we are concerned with non-oscillatory central differencing approximations to the scalar conservation law

$$(2) \quad u_t + f(u)_x = 0.$$

The prototypes of all central schemes are the staggered form of the Lax- Friedrichs (LxF) scheme and its second order extension, the Nessyahu-Tadmor (NT) scheme [15]. For an introduction in central schemes see [8, 9, 13, 15]. For simplicity, we limit our attention to the staggered NT scheme described below. Let  $v(x, t)$  be an approximate solution to (2), and assume that the space mesh  $\Delta x$  and the time mesh  $\Delta t$  are uniform. Let  $x_j := j\Delta x$ ,  $j \in \mathbb{Z}$ ,  $\lambda := \frac{\Delta t}{\Delta x}$  and

$$(3) \quad v_j(t) := \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} v(x, t) dx$$

be the average of  $v$  at time  $t$  over  $(x_{j-1/2}, x_{j+1/2})$ . Let us assume that  $v(\cdot, t)$  is a piecewise linear function, and it is linear on the intervals  $(x_{j-1/2}, x_{j+1/2})$ ,  $j \in \mathbb{Z}$ , of the form

$$(4) \quad v(x, t) = L_j(x, t) := v_j(t) + (x - x_j) \frac{1}{\Delta x} v'_j, \quad x_{j-1/2} < x < x_{j+1/2},$$

where  $\frac{1}{\Delta x} v'_j$  is the numerical derivative of  $v$  which is yet to be determined. Integration of (2) over the staggered space-time cell  $(x_j, x_{j+1}) \times (t, t + \Delta t)$  yields

$$(5) \quad v_{j+1/2}(t + \Delta t) = \frac{1}{\Delta x} \left( \int_{x_j}^{x_{j+1/2}} L_j(x, t) dx + \int_{x_{j+1/2}}^{x_{j+1}} L_{j+1}(x, t) dx \right) - \frac{1}{\Delta x} \left( \int_t^{t+\Delta t} f(v(x_{j+1}, \tau)) d\tau - \int_t^{t+\Delta t} f(v(x_j, \tau)) d\tau \right).$$

The first two integrals on the right of (5) can be evaluated exactly. Moreover, if the CFL condition

$$(6) \quad \lambda \max_{x_j \leq x \leq x_{j+1}} |f'(v(x, t))| \leq \frac{1}{2}, \quad j \in \mathbb{Z},$$

is met, then the last two integrands on the right of (5) are smooth functions of  $\tau$ . Hence, they can be integrated approximately by the midpoint rule with third order local truncation error. Note that, in the case of zero slopes  $\frac{1}{\Delta x} v'_j$  and  $\frac{1}{\Delta x} v'_{j+1}$ , the time integration is exact for any flux  $f$ , and even for nonzero slopes the time integration can be exact for a low degree polynomial flux if a higher order quadrature rule is used. Thus, following [15], we arrive at

$$(7) \quad v_{j+1/2}(t + \Delta t) = \frac{1}{2}(v_j(t) + v_{j+1}(t)) + \frac{1}{8}(v'_j - v'_{j+1}) - \lambda(f(v(x_{j+1}, t + \Delta t/2)) - f(v(x_j, t + \Delta t/2))).$$

By Taylor expansion and the conservation law (2), we obtain

$$(8) \quad v(x_j, t + \Delta t/2) = v_j(t) - \frac{1}{2} \lambda f'_j,$$

where  $\frac{1}{\Delta x} f'_j$  stand for an approximate numerical derivative of the flux  $f(v(x = x_j, t))$ . The following choices are widely used as approximations of the numerical derivatives (we drop  $t$  to simplify the notation)

$$(9) \quad \begin{aligned} v'_j &= \text{m}(v_{j+1} - v_j, v_j - v_{j-1}), \\ f'_j &= \text{m}(f(v_{j+1}) - f(v_j), f(v_j) - f(v_{j-1})), \end{aligned}$$

where  $\text{m}(a, b)$  stands for the minmod limiter

$$(10) \quad \text{m}(a, b) \equiv \text{MinMod}(a, b) := \frac{1}{2}(\text{sgn}(a) + \text{sgn}(b)) \cdot \min(|a|, |b|)$$

with the usual generalization

$$(11) \quad \text{m}(E) := \begin{cases} \inf(E) & \text{if } E \subset \mathbb{R}_+ \\ \sup(E) & \text{if } E \subset \mathbb{R}_- \\ 0 & \text{otherwise} \end{cases}.$$

A generalization of this numerical approximation is based the so-called minmod- $\theta$  limiters

$$(12) \quad \begin{aligned} v'_j &= \text{m}\left(\theta(v_{j+1} - v_j), \frac{1}{2}(v_{j+1} - v_{j-1}), \theta(v_j - v_{j-1})\right), \\ f'_j &= \text{m}\left(\theta(f(v_{j+1}) - f(v_j)), \frac{1}{2}(f(v_{j+1}) - f(v_{j-1})), \theta(f(v_j) - f(v_{j-1}))\right). \end{aligned}$$

Given the approximate slopes and flux derivatives (12), we have a family of central schemes in the predictor-corrector form

$$(13) \quad \begin{aligned} v(x_j, t + \Delta t/2) &= v_j(t) - \frac{1}{2}\lambda f'_j, \\ v_{j+1/2}(t + \Delta t) &= \frac{1}{2}(v_j(t) + v_{j+1}(t)) + \frac{1}{8}(v'_j - v'_{j+1}) \\ &\quad - \lambda(f(v(x_{j+1}, t + \Delta t/2)) - f(v(x_j, t + \Delta t/2))), \end{aligned}$$

where we start with  $v_j(0) := \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u_0(x) dx$ . Note that we alternate between two grids uniform partition of the real line: all intervals with integer end points for  $t = 2k\Delta t$ ,  $k \in \mathbb{Z}$ , and half integers for  $t = (2k + 1)\Delta t$ ,  $k \in \mathbb{Z}$ . As a special case, we recover the staggered LxF scheme for  $\theta = 0$  and the basic MinMod scheme for  $\theta = 1$  (the middle slope in the minmod limiter (12) drops if  $\theta \leq 1$ ).

### 3 $l_2$ Stability for Linear Flux

In this section we will prove that the central scheme given in (13) is  $l_2$ -stable for any  $\theta$  in the interval  $[0, 1]$ . Based on this stability we will also derive a new error estimate in  $L_2$  instead of the usual  $L_1$  estimates in the conservation laws. Note that even for linear flux  $f$ , the minmod-type schemes are not linear and the only global property known was that

the total variation does not increase in time under an appropriate CFL condition, see [15]. The class of minmod-type schemes is also not Lip+ stable except for the obvious choice  $\theta = 0$ . Let us consider a linear flux  $f(u) = au$ , uniform time steps  $t_n = n\Delta t$ , and restrict the minmod limiter to  $\theta \leq 1$ . We denote  $v_j^n := v_j(t_n)$ ,  $\delta_j^n := v_j^n - v_{j-1}^n$ . The minmod scheme (13) reduces to

$$(14) \quad v_{j+1/2}^{n+1} = \frac{1}{2}(v_j^n + v_{j+1}^n) + \frac{\theta}{8}(\mathfrak{m}(\delta_j^n, \delta_{j+1}^n) - \mathfrak{m}(\delta_{j+1}^n, \delta_{j+2}^n)) \\ - \frac{a\Delta t}{\Delta x} \left( v_{j+1}^n - \frac{a\Delta t}{2\Delta x} \theta \mathfrak{m}(\delta_{j+1}^n, \delta_{j+2}^n) - v_j^n + \frac{a\Delta t}{2\Delta x} \theta \mathfrak{m}(\delta_j^n, \delta_{j+1}^n) \right).$$

Hence, we have an explicit formula for the new cell averages (at time  $t_{n+1}$ ) on staggered grid in terms of the old ones (at time  $t_n$ ) on regular grid. In order to simplify the notation, we drop the time dependence and denote  $w_j := v_j^n$ ,  $w'_{j+1} := v_{j+1/2}^{n+1}$ ,  $\delta'_j := w'_j - w'_{j-1}$ ,  $\alpha := \frac{1}{2} + \frac{a\Delta t}{\Delta x}$ , and  $\beta := \frac{1}{2}\alpha(1 - \alpha)$ . With this notation, we have the following relation between the sequence of the new averages  $\{w'_j\}$  and the old ones  $\{w_j\}$

$$(15) \quad w'_j = \alpha w_{j-1} + (1 - \alpha)w_j + \theta\beta(\mathfrak{m}(\delta_{j-1}, \delta_j) - \mathfrak{m}(\delta_j, \delta_{j+1})).$$

Using that  $\delta_j = w_j - w_{j-1}$ , we derive the formula for the sequence of new jumps in terms of the old ones

$$(16) \quad \delta'_j = \alpha\delta_{j-1} + (1 - \alpha)\delta_j - \theta\beta\mathfrak{m}(\delta_{j-2}, \delta_{j-1}) + 2\theta\beta\mathfrak{m}(\delta_{j-1}, \delta_j) - \theta\beta\mathfrak{m}(\delta_j, \delta_{j+1}).$$

The CFL condition (6) reduces to  $0 \leq \alpha \leq 1$  because  $\alpha = 1/2 + \frac{a\Delta t}{\Delta x}$  and  $|\frac{a\Delta t}{\Delta x}| \leq 1/2$ . The main result in this section is the following stability result.

**Theorem 1.** *If the initial condition  $u_0 \in L^1_{loc}(\mathbb{R})$ , then the  $l_2$  norm of the jumps of the approximate solution  $v(\cdot, t)$  is non-increasing in time. That is*

$$(17) \quad \|\{\delta'_j\}\|_{l_2} \equiv \|\{v_j^{n+1} - v_{j-1}^{n+1}\}\|_{l_2} \leq \|\{\delta_j\}\|_{l_2} \equiv \|\{v_j^n - v_{j-1}^n\}\|_{l_2}$$

for all  $n \geq 1$ .

*Proof.* It is clear that we have to prove the result for one step assuming that  $\|\{\delta_j\}\|_{l_2} < \infty$ . We split the proof in two parts. First, we prove the stability for a monotone sequence  $\{w_j\}$ . By symmetry, it is sufficient to consider the case  $\delta_j \geq 0$  for all  $j \in \mathbb{Z}$ . Then we apply that result locally to derive the  $l_2$ -stability for a general sequence.

**Theorem 2.** *Let us assume that  $\delta_j \geq 0$ ,  $j \in \mathbb{Z}$ , and  $\delta'_j$  be given by (16). Then*

$$(18) \quad \|\{\delta'_j\}\|_{l_2} \leq \|\{\delta_j\}\|_{l_2}.$$

*Proof.* Let us recall that  $\{\delta_j\}_{-\infty}^{\infty} \in l_2$ , and  $\delta_j \geq 0$  for all  $j$ . It is enough to prove Theorem 2 only for  $0 < \alpha < 1$ . Let  $\beta_1 := \theta\beta$ , then  $0 < \beta_1 \leq \beta$ . We construct the new sequence  $\{\delta'_j\}$  by using the rule

$$(19) \quad \delta'_j = (1 - \alpha)\delta_j + \alpha\delta_{j-1} - \beta_1 \min(\delta_{j-2}, \delta_{j-1}) + 2\beta_1 \min(\delta_{j-1}, \delta_j) - \beta_1 \min(\delta_j, \delta_{j+1}),$$

for each  $j$ . First we assume that  $\{\delta_j\}$  has finite support. It is easy to see how to modify the proof in case the support is not finite. Therefore we assume  $\delta_j = 0$  for  $j \leq 3$  and for  $j \geq N - 3$  for some integer  $N$ . Then  $\delta'_j = 0$  for  $j \leq 3$  and  $j \geq N - 2$ . Thus it suffices to prove

$$(20) \quad \sum_{j=1}^N \delta_j^2 \geq \sum_{j=1}^N (\delta'_j)^2.$$

Let us introduce some notation. Let  $y_j = \min(\delta_j, \delta_{j+1})$ ,  $\Delta\delta_j = \delta_j - \delta_{j-1}$ ,  $\Delta y_j = y_j - y_{j-1}$ ,  $\Delta^2\delta_j = \delta_j - 2\delta_{j-1} + \delta_{j-2}$ , and  $\Delta^2 y_j = y_j - 2y_{j-1} + y_{j-2}$ . Then (19) becomes

$$\delta'_j = ((1 - \alpha)\delta_j + \alpha\delta_{j-1}) - \beta_1\Delta^2 y_j, \text{ and}$$

$$\sum_{j=1}^N (\delta'_j)^2 = \sum_{j=1}^N \left( ((1 - \alpha)\delta_j + \alpha\delta_{j-1})^2 - 2\beta_1((1 - \alpha)\delta_j + \alpha\delta_{j-1})\Delta^2 y_j + \beta_1^2(\Delta^2 y_j)^2 \right).$$

Note that since  $\delta_0 = \delta_1 = 0$  and  $\delta_{N-1} = \delta_N = 0$ , we have

$$\begin{aligned} \sum_{j=1}^N \delta_j^2 - ((1 - \alpha)\delta_j + \alpha\delta_{j-1})^2 &= \sum_{j=1}^N (1 - (1 - \alpha)^2 - \alpha^2)\delta_j^2 - 2\alpha(1 - \alpha)\delta_j\delta_{j-1} \\ &= 2\beta \sum_{j=1}^N (2\delta_j^2 - 2\delta_j\delta_{j-1}) = 2\beta \sum_{j=1}^N (\delta_j - \delta_{j-1})^2 = 2\beta \sum_{j=1}^N (\Delta\delta_j)^2. \end{aligned}$$

Thus we get

$$(21) \quad \sum_{j=1}^N \delta_j^2 - \sum_{j=1}^N (\delta'_j)^2 = \sum_{j=1}^N (2\beta(\Delta\delta_j)^2 + 2\beta_1((1 - \alpha)\delta_j + \alpha\delta_{j-1})\Delta^2 y_j - \beta_1^2(\Delta^2 y_j)^2).$$

To prove Theorem 2, we need to prove

$$\sum_{j=1}^N (2\beta(\Delta\delta_j)^2 + 2\beta_1((1 - \alpha)\delta_j + \alpha\delta_{j-1})\Delta^2 y_j - \beta_1^2(\Delta^2 y_j)^2) \geq 0.$$

Note that

$$\begin{aligned} &\sum_{j=1}^N (2\beta(\Delta\delta_j)^2 + 2\beta_1((1 - \alpha)\delta_j + \alpha\delta_{j-1})\Delta^2 y_j - \beta_1^2(\Delta^2 y_j)^2) \\ &= \beta_1 \left( \sum_{j=1}^N (2(\beta/\beta_1)(\Delta\delta_j)^2 + 2((1 - \alpha)\delta_j + \alpha\delta_{j-1})\Delta^2 y_j - \beta_1(\Delta^2 y_j)^2) \right) \\ &\geq \beta_1 \left( \sum_{j=1}^N (2(\Delta\delta_j)^2 + 2((1 - \alpha)\delta_j + \alpha\delta_{j-1})\Delta^2 y_j - \beta(\Delta^2 y_j)^2) \right) \\ &= (\beta_1/\beta) \sum_{j=1}^N (2\beta(\Delta\delta_j)^2 + 2\beta((1 - \alpha)\delta_j + \alpha\delta_{j-1})\Delta^2 y_j - \beta^2(\Delta^2 y_j)^2). \end{aligned}$$

Therefore it is sufficient to prove the theorem in the case  $\beta_1 = \beta$ , it is the worst case in certain sense. Now we use  $\Delta^2 y_j = \Delta y_j - \Delta y_{j-1}$ ,  $\Delta y_j = 0$ ,  $\delta_j = 0$  for  $j \leq 1$ ,  $j \geq N - 1$ , and Abel's transform to obtain

$$\sum_{j=1}^N \delta_j \Delta^2 y_j = \sum_{j=1}^N (\delta_j - \delta_{j+1}) \Delta y_j, \quad \text{and} \quad \sum_{j=1}^N \delta_{j-1} \Delta^2 y_j = \sum_{j=1}^N (\delta_{j-1} - \delta_j) \Delta y_j.$$

So, (21) becomes

$$\begin{aligned} \sum_{j=1}^N \delta_j^2 - \sum_{j=1}^N (\delta'_j)^2 &= 2\beta \left( \sum_{j=1}^N (\Delta \delta_j)^2 - (1 - \alpha) \sum_{j=1}^N \Delta \delta_{j+1} \Delta y_j \right. \\ &\quad \left. - \alpha \sum_{j=1}^N \Delta \delta_j \Delta y_j - \frac{\beta}{2} \sum_{j=1}^N (\Delta^2 y_j)^2 \right). \end{aligned}$$

Recall that  $y_j = \min(\delta_j, \delta_{j+1})$ ,  $\Delta \delta_j = \delta_j - \delta_{j-1}$ ,  $\Delta y_j = y_j - y_{j-1}$ ,  $\Delta^2 \delta_j = \delta_j - 2\delta_{j-1} + \delta_{j-2}$ , and  $\Delta^2 y_j = y_j - 2y_{j-1} + y_{j-2}$ . To finish the proof of Theorem 2, it is sufficient to prove the following two Lemmas:

**Lemma 1.**

$$\sum_{j=1}^N (\Delta \delta_j)^2 - (1 - \alpha) \sum_{j=1}^N \Delta \delta_{j+1} \Delta y_j - \alpha \sum_{j=1}^N \Delta \delta_j \Delta y_j - \beta \sum_{j=1}^N (\Delta^2 \delta_j)^2 \geq 0.$$

**Lemma 2.**  $2 \sum_{j=1}^N (\Delta^2 \delta_j)^2 \geq \sum_{j=1}^N (\Delta^2 y_j)^2.$

The proof of Lemma 1. We consider that  $\sum_j$  denotes  $\sum_{j=1}^N$ . Denote

$$A = \sum_j \Delta \delta_{j+1} \Delta y_j, \quad \text{and} \quad B = \sum_j \Delta \delta_j \Delta y_j.$$

Our aim is to prove that

$$(22) \quad \sum_j (\Delta \delta_j)^2 - (1 - \alpha)A - \alpha B - \beta \sum_j (\Delta^2 \delta_j)^2 \geq 0.$$

Let  $u_+ = \max(u, 0)$ ,  $u_- = \min(u, 0)$ . It is easy to check that

$$(23) \quad \Delta y_j = (\Delta \delta_j)_+ + (\Delta \delta_{j+1})_-.$$

We can transform  $A$  as follows:

$$\begin{aligned} (24) \quad A &= \sum_j \Delta \delta_{j+1} ((\Delta \delta_j)_+ + (\Delta \delta_{j+1})_-) \\ &= \sum_j \Delta \delta_{j+1} (\Delta \delta_j)_+ + \sum_j \Delta \delta_j (\Delta \delta_j)_- = \sum_{\Delta \delta_j \leq 0} (\Delta \delta_j)^2 + \sum_{\Delta \delta_j \geq 0} \Delta \delta_j \Delta \delta_{j+1} \\ &= \sum_{\Delta \delta_j \leq 0} (\Delta \delta_j)^2 + \sum_{\Delta \delta_j \geq 0, \Delta \delta_{j+1} \leq 0} \Delta \delta_j \Delta \delta_{j+1} + D, \end{aligned}$$

where  $D = \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j+1} \geq 0} \Delta\delta_j \Delta\delta_{j+1}$ . Further,

$$\begin{aligned}
D &= \frac{1}{2} \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j+1} \geq 0} ((\Delta\delta_j)^2 + (\Delta\delta_{j+1})^2 - (\Delta^2\delta_{j+1})^2) \\
&= \frac{1}{2} \sum_{\Delta\delta_{j-1} \geq 0, \Delta\delta_j \geq 0} ((\Delta\delta_{j-1})^2 + (\Delta\delta_j)^2 - (\Delta^2\delta_j)^2) \\
&= \frac{1}{2} \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j+1} \geq 0} (\Delta\delta_j)^2 + \frac{1}{2} \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j-1} \geq 0} (\Delta\delta_j)^2 - \frac{1}{2} \sum_{\Delta\delta_{j-1} \geq 0, \Delta\delta_j \geq 0} (\Delta^2\delta_j)^2.
\end{aligned}$$

By (24),

$$\begin{aligned}
(25) \quad A &= \sum_j (\Delta\delta_j)^2 - \frac{1}{2} \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j+1} < 0} (\Delta\delta_j)^2 - \frac{1}{2} \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j-1} < 0} (\Delta\delta_j)^2 \\
&\quad - \frac{1}{2} \sum_{\Delta\delta_{j-1} \geq 0, \Delta\delta_j \geq 0} (\Delta^2\delta_j)^2 + \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j+1} \leq 0} \Delta\delta_j \Delta\delta_{j+1}.
\end{aligned}$$

Transform  $B$  in the same way as  $A$ :

$$(26) \quad B = \sum_{\Delta\delta_j \geq 0} (\Delta\delta_j)^2 + \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j+1} \leq 0} \Delta\delta_j \Delta\delta_{j+1} + E,$$

where  $E = \sum_{\Delta\delta_j \leq 0, \Delta\delta_{j+1} \leq 0} \Delta\delta_j \Delta\delta_{j+1}$ . The quantity  $E$  can also be rewritten in the same way as  $D$ :

$$E = \frac{1}{2} \sum_{\Delta\delta_j \leq 0, \Delta\delta_{j+1} \leq 0} (\Delta\delta_j)^2 + \frac{1}{2} \sum_{\Delta\delta_{j-1} \leq 0, \Delta\delta_j \leq 0} (\Delta\delta_j)^2 - \frac{1}{2} \sum_{\Delta\delta_{j-1} \leq 0, \Delta\delta_j \leq 0} (\Delta^2\delta_j)^2.$$

Combining this equality with (26) we get

$$\begin{aligned}
(27) \quad B &= \sum_j (\Delta\delta_j)^2 - \frac{1}{2} \sum_{\Delta\delta_j \leq 0, \Delta\delta_{j+1} > 0} (\Delta\delta_j)^2 - \frac{1}{2} \sum_{\Delta\delta_j \leq 0, \Delta\delta_{j-1} > 0} (\Delta\delta_j)^2 \\
&\quad - \frac{1}{2} \sum_{\Delta\delta_{j-1} \leq 0, \Delta\delta_j \leq 0} (\Delta^2\delta_j)^2 + \sum_{\Delta\delta_{j+1} \leq 0, \Delta\delta_j \geq 0} \Delta\delta_j \Delta\delta_{j+1}.
\end{aligned}$$

By (25) and (27),

$$(28) \quad \sum_j (\Delta\delta_j)^2 - (1 - \alpha)A - \alpha B - \beta \sum_j (\Delta^2\delta_j)^2 = F + G + H + I + J + K + L,$$

where

$$F = \frac{1 - \alpha}{2} \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j+1} < 0} (\Delta\delta_j)^2, \quad G = \frac{1 - \alpha}{2} \sum_{\Delta\delta_j \geq 0, \Delta\delta_{j-1} < 0} (\Delta\delta_j)^2,$$

$$H = \frac{\alpha}{2} \sum_{\Delta\delta_j \leq 0, \Delta\delta_{j+1} > 0} (\Delta\delta_j)^2, \quad I = \frac{\alpha}{2} \sum_{\Delta\delta_j \leq 0, \Delta\delta_{j-1} > 0} (\Delta\delta_j)^2,$$

$$\begin{aligned} J &= - \sum_{\Delta\delta_{j+1} \leq 0, \Delta\delta_j \geq 0} \Delta\delta_j \Delta\delta_{j+1} + \left( \frac{1-\alpha}{2} - \beta \right) \sum_{\Delta\delta_{j-1} \geq 0, \Delta\delta_j \geq 0} (\Delta^2\delta_j)^2 \\ &= + \left( \frac{\alpha}{2} - \beta \right) \sum_{\Delta\delta_{j-1} \leq 0, \Delta\delta_j \leq 0} (\Delta^2\delta_j)^2, \end{aligned}$$

$$K = -\beta \sum_{\Delta\delta_{j-1} > 0, \Delta\delta_j < 0} (\Delta^2\delta_j)^2, \quad \text{and } L = -\beta \sum_{\Delta\delta_{j-1} < 0, \Delta\delta_j > 0} (\Delta^2\delta_j)^2.$$

We have to prove that  $F+G+H+I+J+K+L \geq 0$ . Among the sums  $F, G, H, I, J, K, L$  only the two last sums might be negative; we will show that and

$$(29) \quad F + I + K \geq 0$$

$$(30) \quad G + H + L \geq 0.$$

Indeed,

$$\begin{aligned} \frac{1-\alpha}{2}(\Delta\delta_{j-1})^2 + \frac{\alpha}{2}(\Delta\delta_j)^2 - \beta(\Delta^2\delta_j)^2 &= \frac{1-\alpha}{2}(\Delta\delta_{j-1})^2 + \frac{\alpha}{2}(\Delta\delta_j)^2 \\ &\quad - \frac{(1-\alpha)\alpha}{2}(\Delta\delta_j - \Delta\delta_{j-1})^2 = \frac{1}{2}((1-\alpha)\Delta\delta_{j-1} + \alpha\Delta\delta_j)^2 \geq 0. \end{aligned}$$

Summing the last inequality over all  $j$  with  $\Delta\delta_{j-1} > 0, \Delta\delta_j < 0$ , we get (29). The inequality (30) can be proven in similarly.

Also, we have

$$(31) \quad J \geq 0$$

Finally, plugging (29), (30), and (31) into (28), we obtain the required (22). This completes the proof of Lemma 1.

Proof of Lemma 2.

First, recall that  $\Delta^2 y_j = 0$  for  $j \leq 1$  and  $j \geq N$ . Also, from the proof of Lemma 1 we have:

$$\Delta y_j = \begin{cases} \Delta\delta_{j+1} & \text{if } \Delta\delta_{j+1} \leq 0, \Delta\delta_j \leq 0 \\ \Delta\delta_{j+1} + \Delta\delta_j & \text{if } \Delta\delta_{j+1} \leq 0, \Delta\delta_j \geq 0 \\ \Delta\delta_j & \text{if } \Delta\delta_{j+1} \geq 0, \Delta\delta_j \geq 0 \\ 0 & \text{if } \Delta\delta_{j+1} \geq 0, \Delta\delta_j \leq 0 \end{cases}.$$

Similarly,

$$\Delta y_{j-1} = \begin{cases} \Delta\delta_j & \text{if } \Delta\delta_j \leq 0, \Delta\delta_{j-1} \leq 0 \\ \Delta\delta_j + \Delta\delta_{j-1} & \text{if } \Delta\delta_j \leq 0, \Delta\delta_{j-1} \geq 0 \\ \Delta\delta_{j-1} & \text{if } \Delta\delta_j \geq 0, \Delta\delta_{j-1} \geq 0 \\ 0 & \text{if } \Delta\delta_j \geq 0, \Delta\delta_{j-1} \leq 0 \end{cases}.$$

Therefore  $\Delta\delta_{j-1}, \Delta\delta_j, \Delta\delta_{j+1}$  and their signs determine uniquely  $\Delta^2 y_j$ . We have eight cases depending on what the signs of  $\Delta\delta_{j-1}, \Delta\delta_j, \Delta\delta_{j+1}$  are.

- Case I. (+, +, +), that is  $\Delta\delta_{j-1} \geq 0, \Delta\delta_j \geq 0, \Delta\delta_{j+1} \geq 0$ .  
Then,  $\Delta y_j = \Delta\delta_j, \Delta y_{j-1} = \Delta\delta_{j-1}$ , so  $\Delta^2 y_j = \Delta^2 \delta_j$ , thus,  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_j)^2$  in this case.
- Case II. (+, +, -), that is  $\Delta\delta_{j-1} \geq 0, \Delta\delta_j \geq 0, \Delta\delta_{j+1} < 0$ .  
Then,  $\Delta y_j = \Delta\delta_{j+1} + \Delta\delta_j, \Delta y_{j-1} = \Delta\delta_{j-1}$ , so  $\Delta^2 y_j = \Delta\delta_{j+1} + \Delta\delta_j - \Delta\delta_{j-1}$ .
- Case III. (+, -, +), that is  $\Delta\delta_{j-1} \geq 0, \Delta\delta_j < 0, \Delta\delta_{j+1} \geq 0$ .  
Then,  $\Delta y_j = 0, \Delta y_{j-1} = \Delta\delta_j + \Delta\delta_{j-1}$ , so  $\Delta^2 y_j = -\Delta\delta_j - \Delta\delta_{j-1}$ . In this case  $(\Delta^2 y_j)^2 - (\Delta^2 \delta_j)^2 = 4\Delta\delta_j \Delta\delta_{j-1} \leq 0$ , and  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_j)^2$ .
- Case IV. (+, -, -), that is  $\Delta\delta_{j-1} \geq 0, \Delta\delta_j < 0, \Delta\delta_{j+1} < 0$ .  
Then,  $\Delta y_j = \Delta\delta_{j+1}, \Delta y_{j-1} = \Delta\delta_j + \Delta\delta_{j-1}$ , so  $\Delta^2 y_j = \Delta\delta_{j+1} - \Delta\delta_j - \Delta\delta_{j-1}$ .
- Case V. (-, +, +), that is  $\Delta\delta_{j-1} < 0, \Delta\delta_j \geq 0, \Delta\delta_{j+1} \geq 0$ .  
Then,  $\Delta y_j = \Delta\delta_j, \Delta y_{j-1} = 0$ , so  $\Delta^2 y_j = \Delta\delta_j$ . In this case  $0 \leq \Delta\delta_j < \Delta\delta_j - \Delta\delta_{j-1}$ , and  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_j)^2$ .
- Case VI. (-, +, -), that is  $\Delta\delta_{j-1} < 0, \Delta\delta_j \geq 0, \Delta\delta_{j+1} < 0$ .  
Then,  $\Delta y_j = \Delta\delta_{j+1} + \Delta\delta_j, \Delta y_{j-1} = 0$ , so  $\Delta^2 y_j = \Delta\delta_{j+1} + \Delta\delta_j$ . In this case  $(\Delta^2 y_j)^2 - (\Delta^2 \delta_{j+1})^2 = 4\Delta\delta_{j+1} \Delta\delta_j \leq 0$ , and  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_{j+1})^2$ .
- Case VII. (-, -, +), that is  $\Delta\delta_{j-1} < 0, \Delta\delta_j < 0, \Delta\delta_{j+1} \geq 0$ .  
Then,  $\Delta y_j = 0, \Delta y_{j-1} = \Delta\delta_j$ , so  $\Delta^2 y_j = -\Delta\delta_j$ . In this case  $0 < -\Delta\delta_j \leq \Delta\delta_{j+1} - \Delta\delta_j$ , and  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_{j+1})^2$ .
- Case VIII. (-, -, -), that is  $\Delta\delta_{j-1} < 0, \Delta\delta_j < 0, \Delta\delta_{j+1} < 0$ .  
Then,  $\Delta y_j = \Delta\delta_{j+1}, \Delta y_{j-1} = \Delta\delta_j$ , so  $\Delta^2 y_j = \Delta^2 \delta_{j+1}$ . In this case  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_{j+1})^2$ .

Therefore in cases I (+, +, +), III (+, -, +), and V (-, +, +),  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_j)^2$ , and in cases VI (-, +, -), VII (-, -, +), and VIII (-, -, -),  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_{j+1})^2$ . There are only two “bad” cases: II (+, +, -) and IV (+, -, -) which need a special treatment.

Next, we define a sequence of + and - signs  $\{s_j\}$  where  $s_j = +$  if  $\Delta\delta_j \geq 0$  and  $s_j = -$  if  $\Delta\delta_j < 0$ . Note that  $s_j = +$  for  $j \leq 3$  and  $j \geq N - 2$ . There are three types of “bad” quadruples.

- Type A quadruple: (+, +, -, -), that is  $\Delta\delta_{j-1} \geq 0, \Delta\delta_j \geq 0, \Delta\delta_{j+1} < 0, \Delta\delta_{j+2} < 0$  for some  $j$ . We claim that in this case the following inequality holds:

$$(32) \quad (\Delta^2 y_j)^2 + (\Delta^2 y_{j+1})^2 \leq 2(\Delta^2 \delta_j)^2 + 2(\Delta^2 \delta_{j+1})^2 + 2(\Delta^2 \delta_{j+2})^2.$$

In this case  $\Delta^2 y_j = \Delta\delta_{j+1} + \Delta\delta_j - \Delta\delta_{j-1}$  and  $\Delta^2 y_{j+1} = \Delta\delta_{j+2} - \Delta\delta_{j+1} - \Delta\delta_j$ . If we denote  $\Delta\delta_{j-1}$  by  $a$ ,  $\Delta\delta_j$  by  $b$ ,  $\Delta\delta_{j+1}$  by  $c$ , and  $\Delta\delta_{j+2}$  by  $d$  the above inequality becomes

$$(c+b-a)^2 + (d-c-b)^2 \leq 2(b-a)^2 + 2(c-b)^2 + 2(d-c)^2 \text{ for } a \geq 0, b \geq 0, c < 0, d < 0,$$

which is equivalent to  $a^2 + 2b^2 + 2c^2 + d^2 - 2ab + 2ac - 8bc + 2bd - 2cd \geq 0$ , or

$$(a - b + c)^2 + (b - c + d)^2 - 4bc \geq 0$$

which holds since  $b \geq 0, c < 0$ .

- Type B quadruple:  $(+, +, -, +)$ , that is  $\Delta\delta_{j-1} \geq 0, \Delta\delta_j \geq 0, \Delta\delta_{j+1} < 0, \Delta\delta_{j+2} \geq 0$  for some  $j$ . We claim that in this case the following inequality holds:

$$(33) \quad (\Delta^2 y_j)^2 + (\Delta^2 y_{j+1})^2 \leq 2(\Delta^2 \delta_j)^2 + 2(\Delta^2 \delta_{j+1})^2 + (\Delta^2 \delta_{j+2})^2.$$

In this case  $\Delta^2 y_j = \Delta\delta_{j+1} + \Delta\delta_j - \Delta\delta_{j-1}$  and  $\Delta^2 y_{j+1} = -\Delta\delta_{j+1} - \Delta\delta_j$ . Using the notation we just introduced, the inequality becomes:

$$(c + b - a)^2 + (c + b)^2 \leq 2(b - a)^2 + 2(c - b)^2 + (d - c)^2 \text{ for } a \geq 0, b \geq 0, c < 0, d \geq 0.$$

Since  $(d - c)^2 \geq c^2$  it is sufficient to prove

$$(c + b - a)^2 + (c + b)^2 \leq 2(b - a)^2 + 2(c - b)^2 + c^2, \text{ or}$$

$$a^2 + 2b^2 + c^2 - 2ab + 2ac - 8bc \geq 0, \text{ or } (a - b + c)^2 + b^2 - 6bc \geq 0$$

which holds for  $b \geq 0, c < 0$ .

- Type C quadruple:  $(-, +, -, -)$ , that is  $\Delta\delta_{j-1} < 0, \Delta\delta_j \geq 0, \Delta\delta_{j+1} < 0, \Delta\delta_{j+2} < 0$  for some  $j$ . We claim that in this case the following inequality holds:

$$(34) \quad (\Delta^2 y_j)^2 + (\Delta^2 y_{j+1})^2 \leq (\Delta^2 \delta_j)^2 + 2(\Delta^2 \delta_{j+1})^2 + 2(\Delta^2 \delta_{j+2})^2.$$

In this case  $\Delta^2 y_j = \Delta\delta_{j+1} + \Delta\delta_j$  and  $\Delta^2 y_{j+1} = \Delta\delta_{j+2} - \Delta\delta_{j+1} - \Delta\delta_j$ . Using the notation we just introduced the inequality becomes:

$$(c + b)^2 + (d - c - b)^2 \leq (b - a)^2 + 2(c - b)^2 + 2(d - c)^2 \text{ for } a < 0, b \geq 0, c < 0, d < 0.$$

Since  $(b - a)^2 \geq b^2$  it is sufficient to prove

$$(c + b)^2 + (d + c - b)^2 \leq b^2 + 2(c - b)^2 + 2(d - c)^2, \text{ or}$$

$$b^2 + 2c^2 + d^2 - 8bc + 2bd - 2cd \geq 0, \text{ or } (b - c + d)^2 + c^2 - 6bc \geq 0$$

which holds for  $b \geq 0, c < 0$ .

We will call the type A (32), type B (33), and type C (34) inequalities - “long” inequalities; we call the inequalities of type  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_j)^2$ ,  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_{j+1})^2$  - “short” inequalities.

Now, we construct a set of inequalities. We identify all “bad” quadruples and include the corresponding inequality (type A, type B, or type C) in the set. Next, for all  $j \in [1, N]$  such that  $s_j$  is not a middle element of a “bad” quadruple, and such that  $j$  does not belong to the “bad” cases II and IV we include the corresponding “short” inequality in the set. Finally, we add all inequalities in the set. Taking into account that  $\Delta^2 \delta_j = 0$  and  $\Delta^2 y_j = 0$  for  $j > N$  the resulting inequality is

$$(35) \quad \sum_{j=1}^N a_j (\Delta^2 y_j)^2 \leq \sum_{j=1}^N b_j (\Delta^2 \delta_j)^2,$$

where the  $a_j$ s and  $b_j$ s are non-negative integers. To finish the proof of the lemma we need to show  $a_j \geq 1$  and  $b_j \leq 2$  for all  $j \in [1, N]$ .

Note that all “long” inequalities have the form  $(\Delta^2 y_j)^2 + (\Delta^2 y_{j+1})^2 \leq \dots$ , where  $s_j$  and  $s_{j+1}$  are the middle elements of a “bad” quadruple. Then  $a_j \geq 1$  if  $s_j$  is a middle element of a “bad” quadruple. (By middle element of a quadruple we mean second or third element of the quadruple.)

Now, suppose  $s_j$  is not a middle element of a “bad” quadruple. Then  $j$  does not belong to the “bad” cases II and IV. Indeed if  $j$  is in case II:  $(s_{j-1}, s_j, s_{j+1}) = (+, +, -)$  then  $s_j$  is a middle element of type B quadruple if  $s_{j+2} = +$  and a middle element of type A quadruple if  $s_{j+2} = -$ . Similarly, if  $j$  is in case IV:  $(s_{j-1}, s_j, s_{j+1}) = (+, -, -)$  then  $s_j$  is a middle element of type A quadruple if  $s_{j-1} = +$  and a middle element of type C quadruple if  $s_{j-1} = -$ . Therefore a “short” inequality for  $(\Delta^2 y_j)^2$  has been included in the set of inequalities. Thus  $a_j \geq 1$  in this case as well. We proved  $a_j \geq 1$  for all  $j \in [1, N]$ .

Now, we prove  $b_j \leq 2$  for all  $j \in [1, N]$ . Note that  $(\Delta^2 \delta_j)^2$  can appear in only two “short” inequalities:  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_j)^2$  and  $(\Delta^2 y_{j-1})^2 \leq (\Delta^2 \delta_j)^2$ . Therefore  $b_j \leq 2$  if  $(\Delta^2 \delta_j)^2$  does not appear in any “long” inequalities, that if  $s_j$  is not a second, third, or fourth element of a “bad” quadruple.

The case when  $s_j$  is a second, third, or fourth element of a “bad” quadruple requires more work. First, note that two distinct “bad” quadruples have at most two common elements. Indeed all “bad” quadruples are of the form  $(*, +, -, *)$  where  $*$  denotes  $+$  or  $-$ , and no “bad” quadruple has  $(+, -)$  as its first two or last two elements. Next, the only case when two “bad” quadruples have two common elements is the following configuration:

$$(36) \quad (s_{j-1}, s_j, s_{j+1}, s_{j+2}, s_{j+3}, s_{j+4}) = (+, +, -, +, -, -)$$

Indeed, Type A quadruple can not share exactly two elements with another “bad” quadruple because no “bad” quadruple has  $(-, -)$  as first two elements, or  $(+, +)$  as last two elements. Similar analysis shows that the only way a type B or type C quadruple can share exactly two elements with another type B or type C quadruple is when the configuration (36) occurs.

Let us analyze the configuration (36). The “long” inequalities which correspond to the two “bad” quadruples in this configuration are:

$$(37) \quad (\Delta^2 y_j)^2 + (\Delta^2 y_{j+1})^2 \leq 2(\Delta^2 \delta_j)^2 + 2(\Delta^2 \delta_{j+1})^2 + (\Delta^2 \delta_{j+2})^2, \text{ and}$$

$$(38) \quad (\Delta^2 y_{j+2})^2 + (\Delta^2 y_{j+3})^2 \leq (\Delta^2 \delta_{j+2})^2 + 2(\Delta^2 \delta_{j+3})^2 + 2(\Delta^2 \delta_{j+4})^2.$$

Their sum is

$$\begin{aligned} & (\Delta^2 y_j)^2 + (\Delta^2 y_{j+1})^2 + (\Delta^2 y_{j+2})^2 + (\Delta^2 y_{j+3})^2 \leq \\ & 2(\Delta^2 \delta_j)^2 + 2(\Delta^2 \delta_{j+1})^2 + 2(\Delta^2 \delta_{j+2})^2 + 2(\Delta^2 \delta_{j+3})^2 + 2(\Delta^2 \delta_{j+4})^2. \end{aligned}$$

In this case  $s_j, s_{j+1}, s_{j+2}, s_{j+3}$ , and  $s_{j+4}$  appear as second, third, or fourth element of a “bad” quadruple. Since the configuration (36) starts with  $(+, +)$  and end with  $(-, -)$  it cannot share two elements with a “bad” quadruple outside the configuration. This means that none of  $s_j, s_{j+1}, s_{j+2}, s_{j+3}$ , and  $s_{j+4}$  can be a second, third, or fourth element of a “bad” quadruple outside the configuration. In other words, none of  $(\Delta^2 \delta_j)^2, (\Delta^2 \delta_{j+1})^2, (\Delta^2 \delta_{j+2})^2, (\Delta^2 \delta_{j+3})^2$ , and  $(\Delta^2 \delta_{j+4})^2$  can appear in a “long” inequality other than (37) and (38). Since,  $s_j, s_{j+1}, s_{j+2}$ , and  $s_{j+3}$  are middle elements of “bad” quadruples  $(\Delta^2 \delta_{j+1})^2, (\Delta^2 \delta_{j+2})^2$ , and  $(\Delta^2 \delta_{j+3})^2$  can not appear in “short” inequalities as well. Thus  $b_{j+1} = b_{j+2} = b_{j+3} = 2$ . Also,  $(\Delta^2 \delta_j)^2$  can not appear in a “short” inequality. The only way this could happen is  $(\Delta^2 y_{j-1})^2 \leq (\Delta^2 \delta_j)^2$  which is impossible since  $j - 1$  is either in case I  $(+, +, +)$  or case V  $(-, +, +)$  depending on what  $s_{j-2}$  is, and in both cases the short inequality is  $(\Delta^2 y_{j-1})^2 \leq (\Delta^2 \delta_{j-1})^2$ . Thus  $b_j = 2$ . Similarly,  $(\Delta^2 \delta_{j+4})^2$  can not appear in a “short” inequality. The only way this could happen is  $(\Delta^2 y_{j+4})^2 \leq (\Delta^2 \delta_{j+4})^2$  which is impossible since  $j + 4$  is either in case VII  $(-, -, +)$  or case VIII  $(-, -, -)$  depending on what  $s_{j+5}$  is, and in both cases the short inequality is  $(\Delta^2 y_{j+4})^2 \leq (\Delta^2 \delta_{j+5})^2$ . Thus  $b_{j+4} = 2$ . This concludes the analysis of the configuration (36).

Now, let  $s_j$  be a second, third or fourth element of a “bad” quadruple but not an element of a configuration (36). This means  $(\Delta^2 \delta_j)^2$  appears in exactly one “long” inequality (it can not be a second, third, or fourth element of two distinct “bad” quadruples.) If  $s_j$  is a third element of a “bad” quadruple then  $s_{j-1}$  and  $s_j$  are the middle elements of the quadruple and  $(\Delta^2 \delta_j)^2$  does not appear in a short inequality. Thus,  $b_j \leq 2$  in this case. The cases when  $s_j$  is a second or fourth element of a “bad” quadruple need separate consideration.

1.  $s_j$  is a second element of type A quadruple  $(+, +, -, -)$ . The only way  $(\Delta^2 \delta_j)^2$  could appear in a “short” inequality is  $(\Delta^2 y_{j-1})^2 \leq (\Delta^2 \delta_j)^2$  which is impossible since  $j - 1$  is either in case I  $(+, +, +)$  or case V  $(-, +, +)$  depending on what  $s_{j-2}$  is, and in both cases the short inequality is  $(\Delta^2 y_{j-1})^2 \leq (\Delta^2 \delta_{j-1})^2$ . Thus  $b_j = 2$ .
2.  $s_j$  is a fourth element of type A quadruple  $(+, +, -, -)$ . The only way  $(\Delta^2 \delta_j)^2$  could appear in a “short” inequality is  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_j)^2$  which is impossible since  $j$  is either in case VII  $(-, -, +)$  or case VIII  $(-, -, -)$  depending on what  $s_{j+1}$  is, and in both cases the short inequality is  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_{j+1})^2$ . Thus  $b_j = 2$ .

3.  $s_j$  is a second element of type B quadruple  $(+, +, -, +)$ . Here the argument is word by word like in 1). The only way  $(\Delta^2 \delta_j)^2$  could appear in a “short” inequality is  $(\Delta^2 y_{j-1})^2 \leq (\Delta^2 \delta_j)^2$  which is impossible since  $j - 1$  is either in case I  $(+, +, +)$  or case V  $(-, +, +)$  depending on what  $s_{j-2}$  is, and in both cases the short inequality is  $(\Delta^2 y_{j-1})^2 \leq (\Delta^2 \delta_{j-1})^2$ . Thus  $b_j = 2$ .
4.  $s_j$  is a fourth element of type B quadruple  $(+, +, -, +)$ . Since the coefficient of  $(\Delta^2 \delta_j)^2$  in the corresponding “long” inequality (33) is 1 and  $(\Delta^2 \delta_j)^2$  could appear in at most 1 “short” inequality  $b_j \leq 2$ .
5.  $s_j$  is a second element of type C quadruple  $(-, +, -, -)$ . Since the coefficient of  $(\Delta^2 \delta_j)^2$  in the corresponding “long” inequality (34) is 1 and  $(\Delta^2 \delta_j)^2$  could appear in at most 1 “short” inequality  $b_j \leq 2$ .
6.  $s_j$  is a fourth element of type C quadruple  $(-, +, -, -)$ . Here the argument is word by word like in case 2. The only way  $(\Delta^2 \delta_j)^2$  could appear in a “short” inequality is  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_j)^2$  which is impossible since  $j$  is either in case VII  $(-, -, +)$  or case VIII  $(-, -, -)$  depending on what  $s_{j+1}$  is, and in both cases the short inequality is  $(\Delta^2 y_j)^2 \leq (\Delta^2 \delta_{j+1})^2$ . Thus  $b_j = 2$ .

We have shown that in all six cases  $b_j \leq 2$  for  $j \in [1, N]$  which completes the proof of Lemma 2 and Theorem 2.  $\square$

Now, we continue with the general case, the proof of Theorem 1. That is, we want to show that the  $l_2$  norms inequality

$$(39) \quad \|\{\delta'_j\}\|_{l_2} \leq \|\{\delta_j\}\|_{l_2}$$

holds for *any* initial sequence  $\{\delta_j\}$  with finite  $l_2$  norm. We consider the sequence  $\{w_j\}$  and restrict the index  $j$  to a maximal subset  $\Lambda_m$  on which the piecewise constant function  $w$  is monotone, recall that  $\delta_j = w_j - w_{j-1}$ . Given a sequence  $\{w_j\}$ , we can decompose it into monotone subsequences. This decomposition also gives a decomposition of the sequence  $\{\delta_j\}$  into subsequences such that in each subsequence all jumps have the same sign (non-negative or non-positive). Without any limitations, we assume that the jumps  $\{\delta_j\}$  are non-negative for all  $l \leq j \leq r$ ,  $\delta_{l-1} < 0$  and  $\delta_{r+1} < 0$ . That is,  $w_{l-1}$  is a local minimum and  $w_r$  is a local maximum of the piecewise constant function  $w$ . Let  $w^m$  be the following piecewise constant correction of  $w$

$$(40) \quad w_j^m := \begin{cases} w_j, & \text{if } l \leq j \leq r, \\ w_{l-1}, & \text{if } j < l, \\ w_r, & \text{if } j > r. \end{cases}$$

Note that  $\Lambda_m = \{j : l \leq j \leq r + 1\}$  and the jumps sequence  $\delta^m := \{\delta_j^m\}$  of  $w^m$  is given by

$$(41) \quad \delta_j^m := \begin{cases} w_j - w_{j-1}, & \text{if } l \leq j \leq r, \\ 0, & \text{otherwise} \end{cases}$$

Hence, we have a sequence of monotone functions  $\{w^m\}$  and the corresponding jump sequences  $\{\delta_j^m\}_j$  such that

$$\sum_m \sum_{j \in \Lambda_m} \|\{\delta_j^m\}\|_{l_2}^2 = \sum_{m, j \in \mathbb{Z}} \|\{\delta_j^m\}\|_{l_2}^2 = \|\{\delta_j\}\|_{l_2}^2$$

because the sequence of the jumps of  $\{\delta_j\}$  is decomposed into disjoint jump subsequences  $\{\delta_j^m\}$ . There are two types of jumps  $\delta_j^l$ . A jump  $\delta_j^l$  is of *type 1* if it is equal to the jump  $\delta_j^l(\delta^m)$ , that is the jump generated with the starting sequence  $\{\delta_j^m\}$ , where the index  $m$  such that  $j \in \Lambda_m$ . A jump is of *type 2* if it is not of *type 1*. Note that a *type 2* jump  $\delta_{j^*}^l$  occurs only inside an interval which contains a strict local extremum. Near a local extremum we have two nonzero jumps, say  $\delta_{j^*}^l$  and  $\delta_{j^*}^r$ , generated by the two monotone  $w^m$ -s with index sets finishing/starting with  $j^*$ . It is easy to verify that

$$|\delta_{j^*}^l| = \left| |\delta_{j^*}^l| - |\delta_{j^*}^r| \right|.$$

Hence, we have that

$$(\delta_{j^*}^l)^2 < (\delta_{j^*}^l)^2 + (\delta_{j^*}^r)^2,$$

and we conclude that

$$\sum_j (\delta_j^l)^2 \leq \sum_m \sum_{j \in \Lambda_m} (\delta_j^l(\delta^m))^2 \leq \sum_m \sum_{j \in \Lambda_m} (\delta_j^m)^2 = \sum_n (\delta_j)^2,$$

where we use the notation  $\delta_j^l(\delta^m)$  for the new jumps generated by  $\delta^m$ . It is also easy to prove a local inequality but with index set for  $\delta_j^l$  starting from an interval right after an extremum and finishing right before one.  $\square$

## 4 Error Estimates for Linear Flux

Recall that  $u$  is the entropy solution to the conservation law  $u_t + f(u)_x = 0$  with initial condition  $u^0$ , and  $v$  is the numerical solution described in (13). In the case of linear flux and  $0 \leq \theta \leq 1$ , the formula for the new averages of the minmod scheme is given in (14) and the conservation laws (2) reduces to

$$(42) \quad u_t + au_x = 0.$$

Let  $S_\tau$  be the shift operator defined by  $S_\tau g(\cdot) := g(\cdot - \tau)$ . Then the exact solution of (42) at time  $t$  for any initial data  $u^0$  is  $u(\cdot, t) = S_{at}u^0$ . Let  $A_h$  be the averaging operator defined on a uniform partition by  $A_h g|_I := \frac{1}{h} \int_I g(s) ds$ , where  $|I| = h$ . It will be useful to define a global approximate solution  $v$ . We first define the approximate solution at discrete times by  $v^n := v(\cdot, n\Delta t)$ ,  $n = 0, 1, \dots, N$ , in the following way: (i)  $v^0 := u^0$ ; (ii)  $v^n := S_{a\Delta t} P_h v^{n-1}$ , where for odd  $n$ ,  $1 \leq n \leq N$ ,  $P_h v$  is the linear function on  $I_j := (x_{j-1/2}, x_{j+1/2})$  defined in (4) with the minmod slopes (9), and for even  $n$  we have the analog definition of  $P_h$  on the shifted partition  $\{I_{j+1/2} | j \in \mathbb{Z}\}$ . Note that,  $P_h v^n = P_h A_h v^n$  because the piecewise linear projection  $P_h$ , defined in (4) and (9), is

based only on the averages of  $v^n$  on the corresponding partition. The formula (14) for the new cell averages can be written as

$$v_{j+1/2}^n = A_h(v^n)|_{I_{j+1/2}} = A_h(S_{a\Delta t}P_h v^{n-1})|_{I_{j+1/2}}$$

for odd  $n$ , with  $A_h$  based on the staggered partition  $\{I_{j+1/2}|j \in \mathbb{Z}\}$  and  $P_h$  based on regular partition  $\{I_j|j \in \mathbb{Z}\}$ . For even  $n$ , we have the same sequence of operators but on the reversed partitions. The global approximate solution  $v$  is defined by  $v(\cdot, n\Delta t) = v^n$  and  $v(\cdot, t) = S_{a(t-n\Delta t)}(P_h v^n)$  for  $n\Delta t < t \leq (n+1)\Delta t$  and  $n = 0, 1, \dots, N-1$ . That is  $v$  solves exactly (42) for  $n\Delta t < t \leq (n+1)\Delta t$  with initial data  $P_h v^n$ ,  $n = 0, 1, \dots, N-1$ .

In order to describe the next result, we need to introduce some notation. A function  $g$  is of bounded variation, i.e.,  $g \in \text{BV}(\mathbb{R})$ , if

$$|g|_{\text{BV}(\mathbb{R})} := \sup \sum_{i=1}^n |g(x_{i+1}) - g(x_i)| < \infty,$$

where the supremum is taken over all finite sequences  $x_1 < \dots < x_n$  in  $\mathbb{R}$ . Functions of bounded variation have at most countable many discontinuities, and their left and right limits  $g(x^-)$  and  $g(x^+)$  exist at each point  $x \in \mathbb{R}$ . Since the values of the initial condition  $u^0$  on a set of measure zero have no influence on the numerical solution  $v$  and the entropy solution  $u$ , it is desirable to replace the seminorm  $|\cdot|_{\text{BV}(\mathbb{R})}$  by a similar quantity independent of the function values on sets of measure zero. The standard approach in conservation laws is to consider the space  $\text{Lip}(1, L^1(\mathbb{R}))$  of all functions  $g \in L^1(\mathbb{R})$  such that the seminorm

$$(43) \quad |g|_{\text{Lip}(1, L^1(\mathbb{R}))} := \limsup_{s>0} \frac{1}{s} \int_{\mathbb{R}} |g(x+s) - g(x)| dx$$

is finite. It is clear that  $|g|_{\text{Lip}(1, L^1(\mathbb{R}))}$  will not change if  $g$  is modified on a set of measure zero. At the same time the above two seminorms are equal for functions  $g \in \text{BV}(\mathbb{R})$  such that the value of  $g$  at a point of discontinuity lies between  $g(x^-)$  and  $g(x^+)$  (see Theorem 9.3 in [5]). Similarly, we define the space  $\text{Lip}(1, L^p(\mathbb{R}))$ ,  $1 \leq p \leq \infty$ , which is the set of all functions  $g \in L^p(\mathbb{R})$  for which

$$(44) \quad \|g(\cdot - s) - g(\cdot)\|_{L^p(\mathbb{R})} \leq Ms, \quad s > 0.$$

The smallest  $M \geq 0$  for which (44) holds is  $|g|_{\text{Lip}(1, L^p(\mathbb{R}))}$ . It is easy to see that in the case  $p = 1$  the seminorm given in (44) is the same as the one in (43). In the case  $p > 1$ , the space  $\text{Lip}(1, L^p(\mathbb{R}))$  is essentially the same as  $W^1(L^p(\mathbb{R}))$ , see [5] for details. With this notation, we have the following result.

**Theorem 3.** *Let  $u(x, t) = u(x - at, 0)$  be the solution to (2) with linear flux  $f(z) = az$  and  $v$  be the numerical solution described in (14) with  $0 \leq \theta \leq 1$ . If the CFL condition (6) is satisfied,  $t_n = n\Delta t$ ,  $0 \leq n \leq N$ , and  $T = N\Delta t$ , we have*

$$(45) \quad \|u(\cdot, T) - v(\cdot, T)\|_{L^p(\mathbb{R})} \leq C(Nh)^{1/2} h^{1/2} |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))},$$

for  $p = 1, 2$  where  $C$  is an absolute constant.

*Proof.* The  $L_1$  estimate is based on the TVD property of the numerical solution  $v$  and the  $L_2$  estimate is based on the  $l_2$  stability of the jumps proved in Theorem 1. Both estimates use a dual argument similar to the one in [19] and in the proof we use an index  $p$ , where  $p \in \{1, 2\}$ . Note that we consider the case of linear flux and the usual Lip+ stability requirement is not needed in the dual approach because the the negative norm stability (47) holds for any initial data (not just Lip+). In the proof,  $C$  will be an absolute constant that can be different at different places.

Let  $e(x, t) := u(x, t) - v(x, t)$ , and  $E(x, t) := \int_{-\infty}^x e(s, t) ds$ , where we assume that  $u^0 \in L^1(\mathbb{R})$  to guarantee that  $E$  is well defined for all  $(x, t) \in \mathbb{R} \times (0, T)$ . We have that  $E$  also satisfies (42) for  $n\Delta t < t \leq (n+1)\Delta t$  with initial data  $\int_{-\infty}^x u(s, n\Delta t) - P_h v^n(s) ds$ ,  $n = 0, 1, \dots, N-1$ . For a function  $g \in L^1(\mathbb{R})$  and  $1 \leq p \leq \infty$ , we define a *minus one* norm in the following way

$$(46) \quad \|g\|_{-1,p} := \left\| \int_{-\infty}^{\cdot} g(s) ds \right\|_{L^p(\mathbb{R})}.$$

It is easy to verify that for any  $\tau \in \mathbb{R}$

$$(47) \quad \|S_\tau g\|_{-1,p} = \|g\|_{-1,p}.$$

Recall that  $T = N\Delta t$ . Then, we have the representations  $u(\cdot, T) = (S_{a\Delta t})^N u^0$  and  $v(\cdot, T) = (S_{a\Delta t} P_h)^N u^0$ . Using (47), we have

$$\|e(\cdot, T)\|_{-1,p} = \|(S_{a\Delta t})^N u^0 - (S_{a\Delta t} P_h)^N u^0\|_{-1,p} = \|(S_{a\Delta t})^{N-1} u^0 - P_h (S_{a\Delta t} P_h)^{N-1} u^0\|_{-1,p},$$

and by the triangle inequality we obtain

$$(48) \quad \begin{aligned} \|e(\cdot, T)\|_{-1,p} &\leq \|(S_{a\Delta t})^{N-1} u^0 - (S_{a\Delta t} P_h)^{N-1} u^0\|_{-1,p} \\ &\quad + \|P_h (S_{a\Delta t} P_h)^{N-1} u^0 - (S_{a\Delta t} P_h)^{N-1} u^0\|_{-1,p}. \end{aligned}$$

Let  $e^n = ((S_{a\Delta t})^n - (S_{a\Delta t} P_h)^n) u^0$ ,  $n = 0, 1, \dots, N$ . Then (48) is equivalent to

$$(49) \quad \|e^N\|_{-1,p} \leq \|e^{N-1}\|_{-1,p} + \|P_h v^{N-1} - v^{N-1}\|_{-1,p},$$

and applying (49) for  $n = N, N-1, \dots, 1$ , we get

$$(50) \quad \|e^N\|_{-1,p} \leq \sum_{n=1}^{N-1} \|P_h v^n - v^n\|_{-1,p}$$

because  $e^0 \equiv 0$ . To prove the error estimates, we need the following technical lemma.

**Lemma 4.** *For any  $p \in \{1, 2\}$  and any  $n = 0, 1, \dots, N$ , we have*

- (i)  $\|\{\delta_j^n\}\|_{l_p} \leq h^{1-\frac{1}{p}} |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))}$ ,
- (ii)  $\|P_h v^n - A_h v^n\|_{-1,p} \leq \left(\frac{2}{p+1}\right)^{\frac{1}{p}} h^{1+\frac{1}{p}} \|\{\delta_j\}\|_{l_p}$ ,
- (iii)  $\|A_h v^n - v^n\|_{-1,p} \leq \left(\frac{4}{p+1}\right)^{\frac{1}{p}} h^2 |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))}$ .

*Proof.* The inequalities (i) and (ii) follow by standard arguments, therefore we only prove (i) in the case  $p = 2$  and omit the rest because their proofs are similar. Recall that  $\delta_j^n = v_j^n - v_{j-1}^n$ , and by Theorem 1 we have

$$\left( \sum_j (\delta_j^n)^2 \right)^{1/2} \leq \left( \sum_j (\delta_j^0)^2 \right)^{1/2},$$

where  $\delta_j^0 = u_j^0 - u_{j-1}^0$ ,  $u_j^0 := \frac{1}{h} \int_{I_j} u^0(s) ds$ . Hence, to prove (i) for  $p = 2$ , we need to prove

$$\sum_j (\delta_j^0)^2 \leq h |u^0|_{\text{Lip}(1, L^2(\mathbb{R}))}^2.$$

Since

$$\sum_j (\delta_j^0)^2 = \sum_j \left( \frac{1}{h} \int_{I_j} (u^0(s+h) - u^0(s)) ds \right)^2 \leq h^{-2} \sum_j \left( \int_{I_j} |u^0(s+h) - u^0(s)| ds \right)^2,$$

and since by Cauchy-Schwartz  $\left( \int_{I_j} |u^0(s+h) - u^0(s)| ds \right)^2 \leq h \int_{I_j} |u^0(s+h) - u^0(s)|^2 ds$ , we obtain

$$(51) \quad \sum_j (\delta_j^0)^2 \leq h^{-1} \int_{\mathbb{R}} |u^0(s+h) - u^0(s)|^2 ds.$$

From (44), we have  $\int_{\mathbb{R}} |u^0(s+h) - u^0(s)|^2 ds \leq h^2 |u^0|_{\text{Lip}(1, L^2(\mathbb{R}))}^2$  and using that in (51), we conclude

$$\sum_j (\delta_j^0)^2 \leq h |u^0|_{\text{Lip}(1, L^2(\mathbb{R}))}^2$$

which proves (i) for  $p = 2$ . To prove (iii), we note that

$$(52) \quad |A_h v^n - v^n|_{I_j} \leq \max_{x \in I_j} v^n(x) - \min_{x \in I_j} v^n(x),$$

and because  $v^n = S_{a\Delta t} v^{n-1}$  we have that

$$\max_{x \in I_j} v^n(x) - \min_{x \in I_j} v^n(x) \leq 2 \max(|\delta_{j-1}^{n-1}|, |\delta_j^{n-1}|).$$

The rest of the proof of (iii) is analogous to the proof of (i).  $\square$

Combining (i)-(iii), we have  $\|P_h v^n - v^n\|_{-1, p} \leq Ch^2 |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))}$  and after applying the above inequality in (50), we derive the following estimate

$$(53) \quad \|e^N\|_{-1, p} \leq CNh^2 |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))}.$$

Because  $v^N \notin \text{Lip}(1, L^2(\mathbb{R}))$ , we approximate  $v^N$  by

$$\tilde{v} := \frac{1}{h} \int_{x-h/2}^{x+h/2} A_h v^N(s) ds.$$

Similar to Lemma 4, it is easy to verify that for  $p \in \{1, 2\}$  we have

$$(54) \quad \|\tilde{v} - v^N\|_{-1,p} \leq Ch^2 |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))},$$

$$(55) \quad \|\tilde{v} - v^N\|_{L^p(\mathbb{R})} \leq Ch |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))},$$

and

$$(56) \quad \|\tilde{v}\|_{\text{Lip}(1, L^p(\mathbb{R}))} \leq |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))}.$$

Let  $\tilde{e} := u(\cdot, T) - \tilde{v}$ . Then  $\|\tilde{e}\|_{-1,p} \leq \|e^N\|_{-1,p} + \|\tilde{v} - v^N\|_{-1,p}$ , and combining the estimates (53) and (54), we have

$$(57) \quad \|\tilde{e}\|_{-1,p} \leq CNh^2 |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))}.$$

Kolmogorov-Landau inequalities in  $L^p(\mathbb{R})$  (page 156 in [5]) for the functions  $\tilde{E}(x) := \int_{-\infty}^x \tilde{e}(s) ds$ ,  $\tilde{E}'$ , and  $\tilde{E}''$  give

$$\|\tilde{e}\|_{L^p} = \|\tilde{E}'\|_{L^p} \leq \sqrt{2} \|\tilde{E}\|_{L^p}^{1/2} \|\tilde{E}''\|_{L^p}^{1/2} = \sqrt{2} \|\tilde{e}\|_{-1,p}^{1/2} |\tilde{e}|_{\text{Lip}(1, L^p(\mathbb{R}))}^{1/2}.$$

Using (57) and (56), we arrive at

$$(58) \quad \|\tilde{e}\|_{L^p} \leq C(Nh)^{1/2} h^{1/2} |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))}.$$

Finally, by the triangle inequality

$$\|e\|_{L^p} \leq \|\tilde{e}\|_{L^p} + \|\tilde{v} - v^N\|_{L^p(\mathbb{R})} \leq C(Nh)^{1/2} h^{1/2} |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))},$$

and we combine (55) and (58) to conclude

$$\|u(\cdot, T) - v(\cdot, T)\|_{L^p(\mathbb{R})} = \|e\|_{L^p} \leq C(Nh)^{1/2} h^{1/2} |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))}.$$

Note that  $C$  can be computed explicitly and it is not very big ( $C < 20$ ). In the case  $p = 2$  and  $u^0 \notin L^1(\mathbb{R})$ , we get the same error estimate via an approximation procedure because the estimate is independent of the  $L^1$  norm.  $\square$

**Corollary 5.** *In the case of  $Nh \leq C$ , we get the convergence rate*

$$\|u(\cdot, T) - v(\cdot, T)\|_{L^p(\mathbb{R})} \leq Ch^{1/2} |u^0|_{\text{Lip}(1, L^p(\mathbb{R}))},$$

for  $p = 1$  and  $p = 2$ .

The  $L^1$  estimate is not new, it follows from the arguments in [19], but the  $1/2$  rate in  $L^2$  is new. Note that, using the  $L^1$  estimate, by interpolation arguments, we get only  $1/4$  rate in  $L^2$ . The rate  $1/2$  is optimal for the case  $\theta = 0$  because the numerical method in that case reduces to the LxF scheme, a special case of a monotone scheme. In the case  $p = 1$ , the sharpness of the  $1/2$  bound is given in [20] with an extension to the nonlinear case in [17]. The sharpness in the case  $p = 2$  follows from the more general result for formal first order linear schemes, see [4]. The case  $\theta > 0$  is more complicated because the schemes are nonlinear and it will be addressed elsewhere.

## 5 Numerical Examples

In this section, we present numerical evidence for the new  $l_2$ -stability result we proved in Section 3. Our numerical tests suggest that in the case of linear flux the NT schemes do not increase the  $l_2$  norm of the jumps not only for  $\theta \leq 1$  (as proved in Theorem 1) but also for  $1 < \theta \leq 2$ . In the case of convex flux, we numerically observe the one-sided analog of this property. We now give generic examples for this  $l_2$ -stability in the linear and convex case.

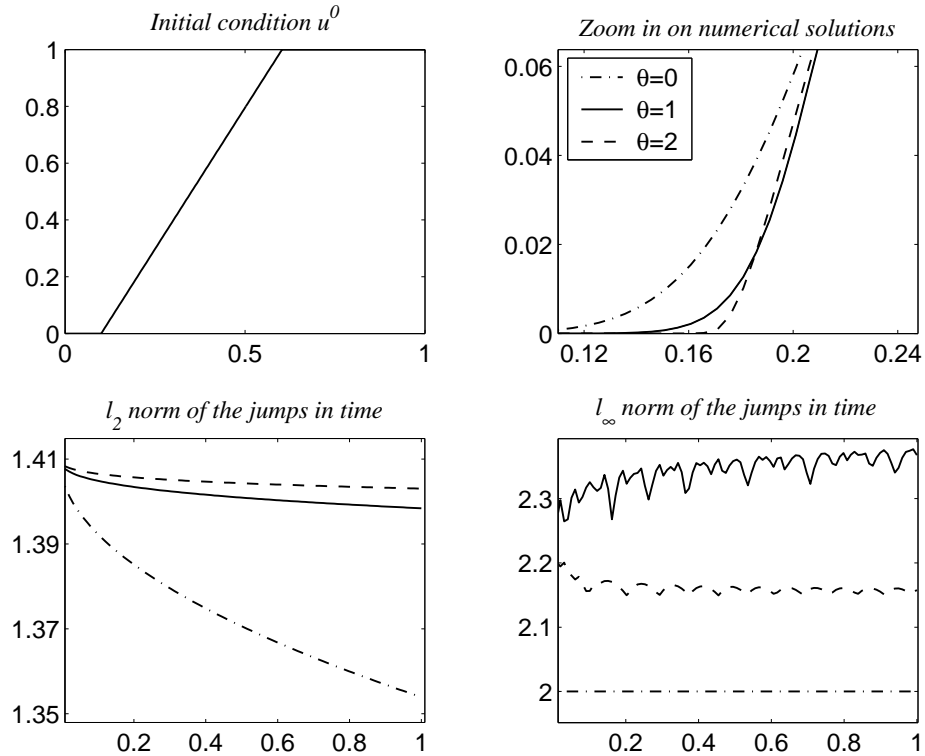


Figure 1:  $u_t + 0.5 u_x = 0$

*Example 1.* Linear equation. We take a piecewise linear initial condition  $u^0$  (see top left on Fig. 1) and compare three different approximate solutions. The solid line represents  $\theta = 1$ , the dashed line represents  $\theta = 2$ , and the dash-dotted line stands for  $\theta = 0$  – the staggered LxF scheme. The values we used are:  $\Delta x = 0.005$ ,  $\lambda = 0.15$ , final time  $T = 0.15$ , and the flux is  $f(u) = 0.5u$ . It is easy to see that for a bigger value of  $\theta$  we get smaller numerical diffusion (see top right on Fig. 1). The other two plots on Fig. 1 give the behaviors of the  $l_2$  and the  $l_\infty$  norm of the jumps in time where the time is re-scaled from  $[0, 0.15]$  to  $[0, 1]$  and the  $l_2$  norm is also re-scaled. Note the oscillatory behavior of the  $l_\infty$  norm and the monotonicity of the  $l_2$  norm for  $\theta = 1, 2$ .

The presence of shocks or local extrema in the initial data will only make the decrease of the  $l_2$  norm of the jumps faster in the beginning and then for large time the  $l_2$  norm will decrease very slowly again. In some sense, the total amount of numerical diffusion is given in the decrease of the  $l_2$  norm. In the so-called second order methods (like  $\theta = 1, 2$ ),

the amount of diffusion is much smaller than in a first order methods represented here

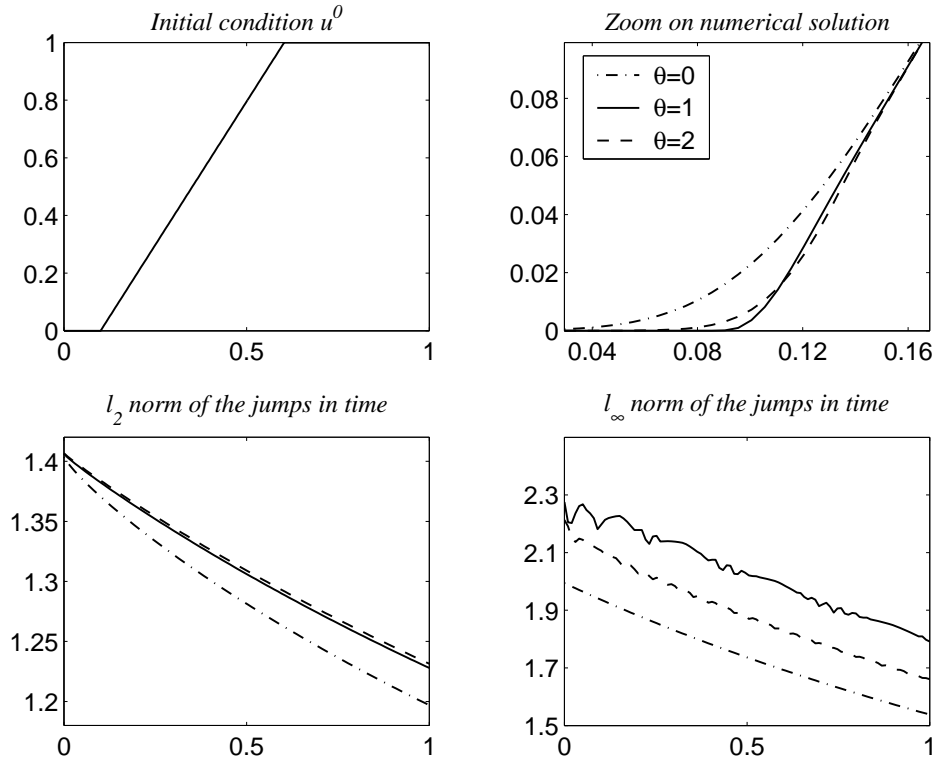


Figure 2:  $u_t + (0.5 u^2)_x = 0$

by the LxF scheme ( $\theta = 0$ ). We will address this issue in a different paper and use it to improve the error estimate for  $\theta = 1$ .

*Example 2.* Burgers' equation. We consider the same initial data,  $\Delta x$ ,  $\lambda$ , and  $T$ , as in the first example. Note again the oscillatory behavior of the  $l_\infty$  norm and the monotonicity of the  $l_2$  norm for  $\theta = 1, 2$ . This is the case because we have non-decreasing initial data which corresponds to a region of spreading of the characteristics.

The nonlinearity of the flux in such regions helps to decrease overall any norm of the jumps. In the case of a general initial condition, the  $l_2$  norm of the jumps decreases in every region of rarefaction. That is, for convex flux the numerical schemes decrease the *one-sided*  $l_2$  norm of the jumps

$$\sum_j (v_j^{n+1} - v_{j-1}^{n+1})_+^2 \leq \sum_j (v_j^n - v_{j-1}^n)_+^2.$$

It is important to note that in the case of convex/concave flux the extreme values separate the regions of rarefactions from the regions of shocks and we observe numerically that the  $l_2$  norm of the jumps decreases in every interval where the numerical solution is non-decreasing/non-increasing.

In the non-convex case (at least one inflection point), the situation is quite different. In one interval of monotonicity we can have both shocks and rarefaction waves. It that

case, the NT scheme with  $\theta = 2$  scheme converges to a wrong weak solution even for the Buckley-Leverett problem, see Example 3 in [8]. Our numerical tests show that the NT scheme give a wrong solution to that problem for any value of  $\theta \geq 1.2$  and in general it looks like the biggest reliable value of  $\theta$  for a non-convex flux is  $\theta = 1$ .

## 6 Acknowledgments

The authors are grateful to Ronald DeVore for his inspiring discussions and constant support.

The authors also thank the anonymous referees. Their comments and suggestions helped improve the paper.

## References

- [1] F. Bouchut, Ch. Bourdarias and B. Perthame, A MUSCL method satisfying all entropy inequalities, *Math. Comp.*, volume 65: 1439–1461, 1996.
- [2] F. Bouchut and B. Perthame, Kruřkov’s estimates for scalar conservation laws revisited, *Trans. AMS*, volume 350, #7: 2847–2870, 1998.
- [3] Y. Brenier and S. Osher, The one-sided Lipschitz condition for convex scalar conservation laws, *SIAM J. Numer. Anal.*, 25: 8–23, 1988.
- [4] P. Brenner, V. Thomée and L. B. Wahlbin, Besov spaces and applications to difference methods for initial value problems, (A. Dold and B. Eckmann, eds.) *Lecture Notes in Math.*, vol. 434, Springer-Verlag, Berlin and New York, 1975.
- [5] R. A. DeVore and G. G. Lorentz, *Constructive Approximation*, Springer-Verlag, Berlin, 1993.
- [6] A. Harten and S. Osher, Uniformly high order accurate non-oscillatory schemes, I, *J. Appl. Num. Math.*, volume 71, 2: 279–309, 1987.
- [7] A. Harten, B. Enquist, S. Osher and S.R. Chakravarthy, Uniformly high order accurate essentially non-oscillatory schemes, III, *J. Comp. Phys.*, volume 71, 2: 231–303, 1987.
- [8] G.-S. Jiang, D. Levi, C.-T. Lin, S. Osher and E. Tadmor High-resolution non-oscillatory central schemes with nonstaggered grids for hyperbolic conservation laws. *SIAM J. Numer. Anal.*, 35, 6: 2147–2169, 1998.
- [9] G.-S. Jiang and E. Tadmor Nonoscillatory central schemes for hyperbolic conservation laws. *SIAM J. Sci. Comput.*, 19, 6: 1892–1917, 1998.
- [10] P. Lax and B. Wendroff, Systems of conservation laws, *Comm. Pure Appl. Math.* 13: 217–237, 1960.

- [11] K. Kopotun, M. Neamtu and B. Popov, Weakly Non-Oscillatory Schemes for Scalar Conservation Laws, to appear in *Math. Comp.*,
- [12] S.N. Kruzhkov, First order quasi-linear equations in several independent variables, *Math. USSR Sbornik*, Vol. 10, #2: 217–243, 1970.
- [13] A. Kurganov and E. Tadmor, New high-resolution central schemes for nonlinear conservation laws and convection-diffusion equations, *J. Comp. Phys.*, volume 160: 241–282, 2000.
- [14] N.N. Kuznetsov, Accuracy of some approximate methods for computing the weak solutions of a first order quasi-linear equations. *USSR Comput. Math. and Math. Phys.*, 16: 105–119, 1976.
- [15] H. Nessyahu and E. Tadmor, Non-oscillatory central differencing for hyperbolic conservation laws, *J. Comp. Phys.*, volume 87, 2: 408–463, 1990.
- [16] H. Nessyahu and E. Tadmor, The convergence rate of nonlinear scalar conservation laws, *SIAM J. Numer. Anal.*, 29: 1505–1519, 1992.
- [17] F. Sabac, The optimal convergence rate of monotone finite difference methods for hyperbolic conservation laws, *SIAM J. Numer. Anal.*, 34: 2306–2318, 1997.
- [18] C.-W. Shu, Numerical experiments on the accuracy of ENO and modified ENO schemes, *J. Comp. Phys.*, 5: 127–149, 1990.
- [19] E. Tadmor, Local error estimates for discontinuous solutions of nonlinear hyperbolic equations, *SIAM J. Numer. Anal.*, 28: 891–906, 1991.
- [20] T. Tang and Z.-H. Teng, The sharpness of Kuznetsov’s  $O(\sqrt{\Delta x})$   $L^1$ -error estimate for monotone difference scheme, *Math. Comp.*, 64: 581–589, 1995.