# Advanced Numerical Methods for Hyperbolic Equations and Applications

## Lecture notes by Michael Dumbser

LABORATORY OF APPLIED MATHEMATICS,
UNIVERSITY OF TRENTO
VIA MESIANO 77, I-38040 TRENTO, ITALY
*E-mail address*: michael.dumbser@ing.unitn.it

# Contents

# Preface

These lecture notes summarize in a very compact form the material presented during the second course of the *Winter School on Advanced Numerical Methods for Free Surface Flows*, held in February 2011 in Trento, Italy. The reader is introduced to the elementary theory and basic concepts of modern numerical methods for the discretization of hyperbolic conservation laws in one and multiple space dimensions. These lecture notes do not claim to be complete and should be used together with additional standard literature available on the topic, such as the books of Toro [77], LeVeque [55], Hirsch [43, 44], Kröner [51] and the classical book of Godlewski–Raviart [36].

Trento, 13/02/2011
Prof. Dr.-Ing. Michael Dumbser

# Part 1

# Basic Algorithms for the Discretization of Hyperbolic Conservation Laws

# Linear Hyperbolic Equations

## 1. Linear scalar advection equation

We consider the following partial differential equation (PDE)

$$(1.1) \qquad \frac{\partial q}{\partial t} + a \frac{\partial q}{\partial x} = 0,$$

with $q = q(x,t)$, $x \in \mathbb{R}$, $t \in \mathbb{R}_0^+$, $a \in \mathbb{R}$ and the initial condition (IC)

$$(1.2) \qquad q(x,0) = h(x).$$

The exact solution of $(1.1)$ is

$$(1.3) \qquad q(x,t) = h(x - at),$$

which can be easily proven as follows:

$$(1.4) \qquad q_t = -ah', q_x = h' \Rightarrow -ah' + ah' = 0.$$

The solution $q(x,t)$ consists of the *transport* of the initial condition $h(x)$, moving with a constant velocity $a$.

### 1.1. Characteristic curves.
The so–called *characteristic curves* or *characteristics* are functions $x(t)$ that satisfy the following initial value problem of the ordinary differential equation (ODE)

$$(1.5) \qquad \frac{dx}{dt} = a,$$

with $x \in \mathbb{R}$, $t \in \mathbb{R}_0^+$, $a \in \mathbb{R}$ and the initial condition $x(0) = x_0$. The value $x_0$ is the initial position or the so–called *foot of the characteristic* and $a$ denotes the characteristic velocity in $(1.5)$. The solution of the ordinary differential equation (ODE) given by $(1.5)$ is

$$(1.6) \qquad x(t) = x_0 + at,$$

plotted in Fig. 1.

The *total derivative* or *material derivative* of $q(x,t)$ computed along the characteristic curve $x(t)$ is given by

$$(1.7) \qquad \frac{dq}{dt} = \frac{\partial q}{\partial t}\frac{dt}{dt} + \frac{\partial q}{\partial x}\frac{dx}{dt} = \frac{\partial q}{\partial t} + a\frac{\partial q}{\partial t} = 0.$$

Hence, along the characteristic $x(t) = x_0 + at$ the PDE $(1.1)$ reduces to the trivial ODE

$$(1.8) \qquad \frac{dq}{dt} = 0.$$

It follows that the solution $q = q(x,t)$ remains *constant* along the characteristic curve. Hence, the solution at a general point $x$ in space and at a general time $t$ can be computed by evaluating the initial condition at the *foot of the characteristic* passing through $x$ and $t$:

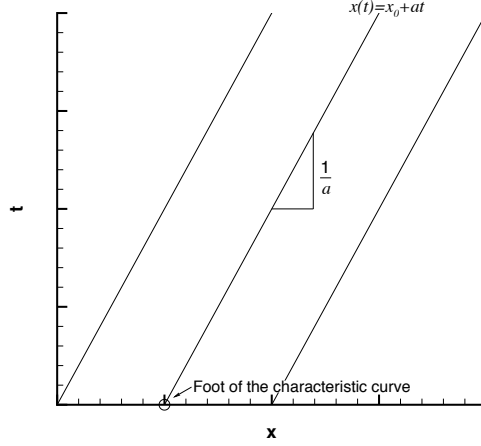$$(1.9) \qquad q(x,t) = q(x_0, 0) = h(x_0) = h(x - at).$$

FIGURE 1. The characteristic curves resulting from the initial value problem (1.5).

Note that for the construction of the solution in Eqn. (1.9), the differentiability of the solution $q = q(x,t)$ has not been used. In fact, the function $q = q(x,t)$ can also be discontinuous and discontinuities will propagate along characteristic curves.

Already with this very simple example, we can see the important difference between the *Eulerian* description of physical processes, where the observer looks at the process from a fixed laboratory frame, and the *Lagrangian* description of physical processes, where the observer moves along characteristic curves. In the Eulerian case, the physics of advection is described by a partial differential equation, here the PDE (1.1). In the Lagrangian case, the physics of advection is described by a system of two ordinary differential equations, namely the ODE (1.5) for the evolution of the observer's position and the ODE (1.8), which describes the temporal evolution of the state variable $q$ in the co–moving frame of the observer. While the Eulerian description uses *partial* derivatives to model the physics, the Lagrangian description uses *total* or *material* derivatives. However, both formulations are equivalent, even in the presence of discontinuities, as proven by Wagner in [**89**].

EXAMPLE 1. *Let us assume a unitary characteristic velocity $a = 1$ in the linear scalar advection equation* (1.1) *and the following initial condition*

$$(1.10) \qquad q(x,0) = h(x) = \begin{cases} 1, & \text{if } -1 \leq x \geq 1, \\ 0, & \text{otherwise.} \end{cases}$$

*Then, the exact solution reproduces the initial profile translated to the right with a velocity of one. The discontinuities remain sharp, since the PDE* (1.1) *does not contain any diffusion terms. The exact solution is plotted in Fig. 2 following the Eulerian description (on the left) and the Lagrangian approach (on the right).*

**1.2. Riemann problem.** The Riemann problem is a particular initial value problem (Cauchy problem) of a PDE. The initial condition consists in this case of two piecewise constant states $q_L$ and $q_R$, separated by a discontinuity:

$$(1.11) \qquad q_t + aq_x = 0,$$

$$(1.12) \qquad q(x,0) = h(x) = \begin{cases} q_L, & \text{if } x \leq 0, \\ q_R, & \text{if } x > 0, \end{cases}$$

FIGURE 2. The exact solution of the initial value problem (1.10) of the linear scalar advection equation (1.1) with $a = 1$ according to the Eulerian description (on the left) and the Lagrangian description (on the right).



FIGURE 3. Sketch of the exact solution of the Riemann problem for the linear scalar advection equation.

where the notation $q_t$, $q_x$ denotes the partial derivatives of the function $q(x, t)$ with respect to $t$ and $x$, respectively. The exact solution of the Riemann problem is

$$(1.13) \qquad q(x, t) = h(x - at) = \begin{cases} q_L, & \text{if } x - at \leq 0, \\ q_R, & \text{if } x - at > 0, \end{cases}$$

that is equivalent for $t > 0$ to

$$(1.14) \qquad q(x, t) = \begin{cases} q_L, & \text{if } \frac{x}{t} \leq a, \\ q_R, & \text{if } \frac{x}{t} > a. \end{cases}$$

An example of an exact solution of the Riemann problem for the linear scalar advection equation is plotted in Fig. 3.

**1.3. Initial boundary value problems.** Let us consider a general initial boundary value problem given by the partial differential equation (PDE) and the

corresponding initial condition (IC) and boundary conditions (BC), as follows

$$
\begin{aligned}
\text{PDE:} &\quad q_t + aq_x = 0, \quad x \in [x_L, x_R], t \in \mathbb{R}_0^+, \\
\text{IC:} &\quad q(x,0) = h(x), \\
\text{BC:} &\quad q(x_L, t) = b_L(t), \quad q(x_R, t) = b_R(t).
\end{aligned}
$$
(1.15)

The exact solution for $a \in \mathbb{R}$ is the following
(1.16)
$$
q(x,t) = \begin{cases}
h(x - at), & \text{if } \max(x_L, x_L + at) \le x \le \min(x_R, x_R + at), \\
b_L\left(t - \frac{x - x_L}{a}\right), & \text{if } x_L \le x < x_L + at, \\
b_R\left(t - \frac{x_R - x}{a}\right), & \text{if } x_R \ge x > x_R + at.
\end{cases}
$$

## 2. Linear scalar advection reaction equation

We consider a linear scalar advection reaction equation

$$
q_t + aq_x = \beta q, \quad \text{with } x \in \mathbb{R}, t \in \mathbb{R}_0^+, \beta \in \mathbb{R}_0^+,
$$
(1.17)

$$
q(x,0) = h(x),
$$
(1.18)

where $\beta$ in the linear source term is a constant. Solving the linear scalar advection reaction equation (1.17) along a characteristic curve $x = x_0 + at$, given again by the solution of the ODE (1.5) $\frac{dx}{dt} = a$, we obtain

$$
\frac{dq}{dt} = \frac{\partial q}{\partial t}\frac{dt}{dt} + \frac{\partial q}{\partial x}\frac{dx}{dt} = q_t + aq_x = \beta q.
$$
(1.19)

The exact solution along the characteristic curve is the solution of the previous ODE:

$$
q(t) = q(0) e^{\beta t}.
$$
(1.20)

Hence, the exact solution of the linear scalar advection reaction equation (1.17) is

$$
q(x,t) = h(x - at) e^{\beta t}.
$$
(1.21)

The fact that Eqn. (1.21) is the exact solution of the advection–reaction equation can be proven by computing the derivative in time of Eqn. (1.21)

$$
q_t = h\beta e^{\beta t} - ah' e^{\beta t},
$$
(1.22)

and the derivative in space

$$
q_x = h' e^{\beta t}.
$$
(1.23)

Replacing these expressions in (1.19) yields

$$
q_t + aq_x = h\beta e^{\beta t} - ah' e^{\beta t} + ah' e^{\beta t} = \beta h e^{\beta t} = \beta q.
$$
(1.24)

## 3. Linear advection with variable coefficient

We consider the following initial value problem

$$
q_t + a(x,t) q_x = 0,
$$
(1.25)

$$
q(x,0) = h(x).
$$
(1.26)

The characteristic curves are the solution of the problem

$$
\frac{dx}{dt} = a(x(t), t),
$$
(1.27)

$$
x(0) = x_0.
$$
(1.28)

In this case, the characteristics are in general no longer straight lines, but curves. The formal solution of the ODE (1.27) is

$$(1.29) \qquad x(t) = x_0 + \int_0^t a\left(x(\tau), \tau\right) d\tau,$$

and hence the exact solution of the advection equation with variable coefficient (1.25) can be written formally as

$$(1.30) \qquad q(x,t) = h(x_0) = h\left(x(t) - \int_0^t a\left(x(\tau), \tau\right) d\tau\right).$$

EXERCISE 1. *Prove that* (1.30) *actually is a solution of* (1.25).

EXAMPLE 2. *We consider the following problem*

$$(1.31) \qquad q_t + cxq_x = 0,$$

$$(1.32) \qquad q(x,0) = h(x) = \begin{cases} 1, & \text{if } |x| \le 1, \\ 0, & \text{if } |x| > 1. \end{cases}$$

*The corresponding characteristic curves are given by the following ODE*

$$(1.33) \qquad \frac{dx}{dt} = cx,$$

$$(1.34) \qquad x(0) = x_0,$$

*the solution of which leads to*

$$(1.35) \qquad x(t) = x_0 e^{ct}, \quad \Rightarrow \quad x_0 = xe^{-ct}.$$

*It follows that the exact solution of the problem* (1.32) *is*

$$(1.36) \qquad q(x,t) = h(x_0) = h\left(xe^{-ct}\right) = \begin{cases} 1, & \text{if } |xe^{-ct}| \le 1, \\ 0, & \text{if } |xe^{-ct}| > 1. \end{cases}$$

EXAMPLE 3. *Another linear advection equation with variable coefficient is given by*

$$(1.37) \qquad q_t + cxtq_x = 0,$$

$$(1.38) \qquad q(x,0) = h(x) \begin{cases} 1, & \text{if } |x| \le 1, \\ 0, & \text{if } |x| > 1. \end{cases}$$

*We apply the method of the characteristics, as follows*

$$(1.39) \qquad \frac{dx}{dt} = cxt,$$

$$(1.40) \qquad x(0) = x_0,$$

*with the solution*

$$(1.41) \qquad x(t) = x_0 e^{\frac{ct^2}{2}}, \quad \Rightarrow \quad x_0 = xe^{\frac{-ct^2}{2}}.$$

*This leads to the solution*

$$(1.42) \qquad q(x,t) = h(x_0) = h\left(xe^{-\frac{ct^2}{2}}\right) \begin{cases} 1, & \text{if } \left|xe^{-\frac{ct^2}{2}}\right| \le 1, \\ 0, & \text{if } \left|xe^{-\frac{ct^2}{2}}\right| > 1. \end{cases}$$

## 4. Linear hyperbolic systems

We consider the following system of equations

$$(1.43) \qquad \frac{\partial \vec{Q}}{\partial t} + \underline{\underline{A}} \cdot \frac{\partial \vec{Q}}{\partial x} = \vec{0}, \quad x \in \mathbb{R}, t \in \mathbb{R}_0^+, \underline{\underline{A}} \in \mathbb{R}^{m \times m},$$

where the vector of the unknowns $\vec{Q} = \vec{Q}(x,t)$ and the matrix $\underline{\underline{A}}$ are defined as

$$(1.44) \qquad \vec{Q} = \begin{pmatrix} q_1 \\ q_2 \\ ... \\ q_m \end{pmatrix}, \quad \underline{\underline{A}} = \begin{pmatrix} a_{11} & a_{12} & ... & a_{1m} \\ a_{21} & a_{22} & ... & a_{2m} \\ ... & ... & ... & ... \\ a_{m1} & a_{m2} & ... & a_{mm} \end{pmatrix},$$

with the following initial conditions

$$(1.45) \qquad \vec{Q}(x,0) = \vec{h}(x).$$

The system of equations (1.43) is called *hyperbolic* if all the eigenvalues of the matrix $\underline{A}$ are real numbers *and* if there exists a complete set of $m$ eigenvectors that are linearly independent. A system (1.43) is *strictly hyperbolic* if all the eigenvalues are real and *distinct*. A system that is strictly hyperbolic is necessarily hyperbolic because distinct eigenvalues automatically imply the existence of a complete set of linearly independent eigenvectors. A system (1.43) is called *weakly hyperbolic* if all the eigenvalues are real, but not all are distinct and there does *not exist* a complete set of linearly independent eigenvectors.

The properties of the right eigenvectors $\vec{r}$ are

$$(1.46) \qquad \underline{\underline{A}} \cdot \vec{r} = \lambda \vec{r},$$

where $\lambda$ denotes the eigenvalues, which are the roots of the characteristic polynomial

$$(1.47) \qquad P(\lambda) = \left| \underline{\underline{A}} - \lambda \underline{\underline{I}} \right| = 0,$$

with the identity matrix $\underline{\underline{I}}$. The left eigenvectors $\vec{l}$ satisfy the relation

$$(1.48) \qquad \vec{l}^T \cdot \underline{\underline{A}} = \lambda \vec{l}^T.$$

The eigenvalues $\lambda$ in (1.46) and (1.48) are the same, which can be seen easily by some simple manipulations:

$$(1.49) \qquad \vec{l}^T \underline{\underline{A}} = \left( \underline{\underline{A}}^T \vec{l} \right)^T,$$

hence with (1.48) it follows that

$$(1.50) \qquad \underline{\underline{A}}^T \vec{l} = \lambda \vec{l},$$

whose characteristic polynomial is

$$(1.51) \qquad P(\lambda) = |\underline{\underline{A}}^T - \lambda \underline{\underline{I}}| = |\underline{\underline{A}} - \lambda \underline{\underline{I}}|,$$

which is identical to the characteristic polynomial of the right eigenvectors.

We are furthermore interested in a relationship between the left and right eigenvectors. Let us define $\underline{\underline{R}}$, $\underline{\underline{L}}$ and $\underline{\underline{A}}$, as follows

$$(1.52) \qquad \underline{\underline{R}} = (\vec{r}_1, \vec{r}_2, ..., \vec{r}_m), \quad \underline{\underline{L}} = \begin{pmatrix} \vec{l}_1 \\ \vec{l}_2 \\ ... \\ \vec{l}_m \end{pmatrix}, \quad \underline{\underline{\Lambda}} = \begin{pmatrix} \lambda_1 & 0 & ... & 0 \\ 0 & \lambda_2 & ... & 0 \\ ... & ... & ... & ... \\ 0 & 0 & ... & \lambda_m \end{pmatrix}.$$

According to the definition of the right eigenvectors (1.46) and by multiplying with the matrix of the left eigenvectors $L$ from the left, we obtain

$$(1.53) \qquad \underline{\underline{L}}(\underline{\underline{A}}\,\underline{\underline{R}}) = \underline{\underline{L}}(\underline{\underline{R}}\,\underline{\underline{\Lambda}}),$$

and similarly for the left eigenvectors,

$$(\underline{L}\,\underline{\underline{A}})\,\underline{\underline{R}} = (\underline{\underline{\Lambda}}\,\underline{L})\,\underline{\underline{R}}. \tag{1.54}$$

This leads to

$$\underline{\underline{L}}\,\underline{\underline{R}}\,\underline{\underline{\Lambda}} = \underline{\underline{\Lambda}}\,\underline{\underline{L}}\,\underline{\underline{R}}. \tag{1.55}$$

Hence, $\left(\underline{\underline{L}}\,\underline{\underline{R}}\right)$ and $\underline{\underline{\Lambda}}$ are commutative matrices. The only matrix that is commutative with a diagonal matrix with distinct values, is also a *diagonal* matrix. It follows that the right and left eigenvectors are orthogonal, i.e.

$$\vec{l}_i^T \vec{r}_j \begin{cases} \neq 0, & \text{if } i = j, \\ = 0, & \text{if } i \neq j. \end{cases} \tag{1.56}$$

With a proper choice of the scaling factors of the eigenvectors, one can *always* obtain the result that the right and left eigenvectors become *orthonormal*, that is

$$\underline{\underline{L}}\,\underline{\underline{R}} = \underline{\underline{I}}, \quad \text{or} \quad \underline{\underline{L}} = \underline{\underline{R}}^{-1}. \tag{1.57}$$

Finally, from (1.46) and (1.57) we obtain

$$\underline{\underline{A}} = \underline{\underline{R}}\,\underline{\underline{\Lambda}}\,\underline{\underline{R}}^{-1} = \underline{\underline{R}}\,\underline{\underline{\Lambda}}\,\underline{\underline{L}}, \tag{1.58}$$

which is equivalent to

$$\underline{\underline{\Lambda}} = \underline{\underline{R}}^{-1}\underline{\underline{A}}\,\underline{\underline{R}} = \underline{\underline{L}}\underline{\underline{A}}\,\underline{\underline{R}}. \tag{1.59}$$

**4.1. Characteristic variables.** To solve the linear hyperbolic system of equations one introduces the so called characteristic variables $\vec{C}$, defined by

$$\vec{Q} = \underline{\underline{R}}\vec{C}, \quad \vec{C} = \underline{\underline{R}}^{-1}\vec{Q}. \tag{1.60}$$

Then, the system (1.43) becomes

$$\frac{\partial}{\partial t}\left(\underline{\underline{R}}\vec{C}\right) + \underline{\underline{A}}\frac{\partial}{\partial x}\left(\underline{\underline{R}}\vec{C}\right) = 0. \tag{1.61}$$

In the case of a linear hyperbolic system with constant coefficients, the matrix $\underline{\underline{A}}$ and thus also $\underline{\underline{R}}$, $\underline{\underline{R}}^{-1}$ and $\underline{\underline{\Lambda}}$ do not depend on $(x,t)$, hence we have

$$\underline{\underline{R}}\frac{\partial \vec{C}}{\partial t} + \underline{\underline{R}}\,\underline{\underline{A}}\frac{\partial \vec{C}}{\partial x} = 0. \tag{1.62}$$

According to (1.59) this leads to a new formulation of the linear system (1.43) written in characteristic variables

$$\frac{\partial \vec{C}}{\partial t} + \underline{\underline{\Lambda}}\frac{\partial \vec{C}}{\partial x} = 0. \tag{1.63}$$

The system (1.63) is a system of $m$ *decoupled* linear scalar advection equations,

$$\frac{\partial c_i}{\partial t} + \lambda_i \frac{\partial c_i}{\partial x} = 0, \quad \text{for } i = 1,..,m, \tag{1.64}$$

which can be easily solved componentwise. The initial condition for the system (1.63) is just obtained from (1.60), as follows

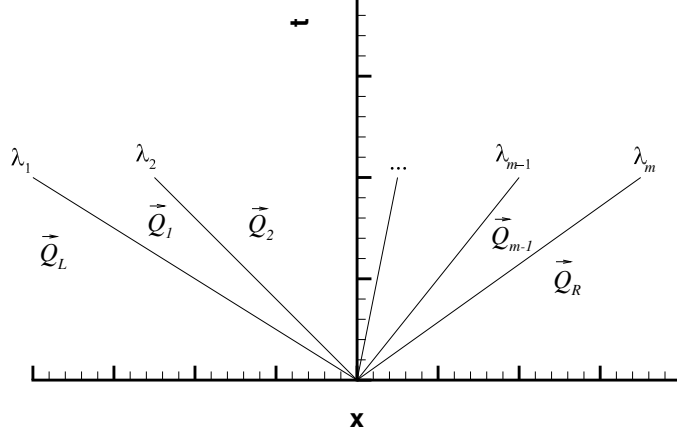$$\vec{C}(x,0) = \underline{\underline{R}}^{-1}\vec{h}(x). \tag{1.65}$$

FIGURE 4. Sketch of the exact solution of the Riemann problem for a linear hyperbolic system of equations.

**4.2. Riemann problem.** Now, we consider the Riemann problem

$$\frac{\partial \vec{Q}}{\partial t} + \underline{\underline{A}} \frac{\partial \vec{Q}}{\partial x} = 0, \tag{1.66}$$

$$\vec{Q}(x,0) = \vec{h}(x) = \begin{cases} \vec{Q}_L, & \text{if } x \leq 0, \\ \vec{Q}_R, & \text{if } x > 0. \end{cases} \tag{1.67}$$

Written in characteristic variables it becomes

$$\frac{\partial \vec{C}}{\partial t} + \underline{\underline{\Lambda}} \frac{\partial \vec{C}}{\partial x} = 0. \tag{1.68}$$

$$\vec{C}(x,0) = \begin{cases} \underline{\underline{R}}^{-1} \vec{Q}_L := \vec{C}_L, & \text{if } x \leq 0, \\ \underline{\underline{R}}^{-1} \vec{Q}_R := \vec{C}_R, & \text{if } x > 0. \end{cases} \tag{1.69}$$

Then, the exact solution of the Riemann problem (1.66) written in terms of the characteristic variables is

$$c_i(x,t) = \begin{cases} c_i^L, & \text{if } \frac{x}{t} \leq \lambda_i, \\ c_i^R, & \text{if } \frac{x}{t} > \lambda_i, \end{cases} \tag{1.70}$$

or, in more compact form, using the matrix sign function, which is applied componentwise for a diagonal matrix:

$$\vec{C}(x,t) = \frac{1}{2}\left(1 + \text{sign}(\underline{\underline{\Lambda}} - \frac{x}{t}\underline{\underline{I}})\right)\vec{C}_L + \frac{1}{2}\left(1 - \text{sign}(\underline{\underline{\Lambda}} - \frac{x}{t}\underline{\underline{I}})\right)\vec{C}_R. \tag{1.71}$$

The exact solution in physical variables results from (1.60). Hence, we obtain the following general expression for the solution in physical variables in terms of the initial condition (1.67):

$$\vec{Q}(x,t) = \frac{1}{2}\underline{\underline{R}}\left(\underline{\underline{I}} + \text{sign}(\underline{\underline{\Lambda}} - \frac{x}{t}\underline{\underline{I}})\right)\underline{\underline{R}}^{-1}\vec{Q}_L + \frac{1}{2}\underline{\underline{R}}\left(\underline{\underline{I}} - \text{sign}(\underline{\underline{\Lambda}} - \frac{x}{t}\underline{\underline{I}})\right)\underline{\underline{R}}^{-1}\vec{Q}_R. \tag{1.72}$$

The solution consists of $m + 1$ piecewise constant states, separated by the characteristics moving with speeds $\lambda_i$, see Fig. 4.

EXAMPLE 4. *Let us solve the following system of equations*

$$\frac{\partial \vec{Q}}{\partial t} + \underline{\underline{A}} \frac{\partial \vec{Q}}{\partial x} = \vec{0}, \quad \text{with } \underline{\underline{A}} = \begin{pmatrix} 0 & a \\ b & 0 \end{pmatrix}, \tag{1.73}$$

$$(1.74) \qquad \vec{Q} = h(x).$$

*The characteristic polynomial results*

$$(1.75) \qquad P(\lambda) = \left\| \begin{array}{cc} -\lambda & a \\ b & -\lambda \end{array} \right\| = \lambda^2 - ab = 0,$$

*and the corresponding eigenvalues are $\lambda_{1,2} = \mp\sqrt{ab}$. Now, we compute the right eigenvectors, solving the following linear systems related to the first and second eigenvalue, respectively*

$$(1.76) \qquad \left( \begin{array}{cc} +\sqrt{ab} & a \\ b & +\sqrt{ab} \end{array} \right) \left( \begin{array}{c} x_1 \\ x_2 \end{array} \right) = \left( \begin{array}{c} 0 \\ 0 \end{array} \right),$$

$$(1.77) \qquad \left( \begin{array}{cc} -\sqrt{ab} & a \\ b & -\sqrt{ab} \end{array} \right) \left( \begin{array}{c} x_1 \\ x_2 \end{array} \right) = \left( \begin{array}{c} 0 \\ 0 \end{array} \right).$$

*We obtain*
$$(1.78)$$
$$\underline{\underline{R}} = \left( \begin{array}{cc} 1 & 1 \\ -\sqrt{\frac{b}{a}} & +\sqrt{\frac{b}{a}} \end{array} \right), \quad \underline{\underline{\Lambda}} = \left( \begin{array}{cc} -\sqrt{ba} & 0 \\ 0 & +\sqrt{ba} \end{array} \right), \quad \underline{\underline{R}}^{-1} = \frac{1}{2} \left( \begin{array}{cc} 1 & -\sqrt{\frac{a}{b}} \\ 1 & +\sqrt{\frac{a}{b}} \end{array} \right).$$

EXAMPLE 5. *We consider a system of three equations, as follows*

$$(1.79) \qquad \frac{\partial \vec{Q}}{\partial t} + \underline{\underline{A}} \frac{\partial \vec{Q}}{\partial x} = \vec{0}, \quad \text{with } \underline{\underline{A}} = \left( \begin{array}{ccc} 1 & 2 & -2 \\ -1 & 1 & 1 \\ -1 & 2 & 0 \end{array} \right),$$

$$(1.80) \qquad \vec{Q} = h(x).$$

*We compute the eigenvalues from the characteristic polynomial*
$$(1.81)$$
$$P(\lambda) = \left\| \begin{array}{ccc} 1-\lambda & 2 & -2 \\ -1 & 1-\lambda & 1 \\ -1 & 2 & -\lambda \end{array} \right\| = (1-\lambda)(-\lambda(1-\lambda)-2)-2(\lambda+1)-2(-2+(1-\lambda)) = 0.$$

*They are $\lambda_{1,2} = \mp 1$ and $\lambda_3 = 2$. Using the first eigenvalue we obtain the following right eigenvector*

$$(1.82) \qquad \left( \begin{array}{ccc} 2 & 2 & -2 \\ -1 & 2 & 1 \\ -1 & 2 & 1 \end{array} \right) \left( \begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right) = \left( \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right), \quad \Rightarrow \quad \vec{r_1} = \left( \begin{array}{c} 1 \\ 0 \\ 1 \end{array} \right).$$

*Similarly, for the other eigenvectors we have*

$$(1.83) \qquad \left( \begin{array}{ccc} 0 & 2 & -2 \\ -1 & 0 & 1 \\ -1 & 2 & -1 \end{array} \right) \left( \begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right) = \left( \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right), \quad \Rightarrow \quad \vec{r_2} = \left( \begin{array}{c} 1 \\ 1 \\ 1 \end{array} \right),$$

$$(1.84) \qquad \left( \begin{array}{ccc} -1 & 2 & -2 \\ -1 & -1 & 1 \\ -1 & 2 & -2 \end{array} \right) \left( \begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right) = \left( \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right), \quad \Rightarrow \quad \vec{r_3} = \left( \begin{array}{c} 0 \\ 1 \\ 1 \end{array} \right).$$

*Finally, the matrices $\underline{\underline{R}}$, $\underline{\underline{R}}^{-1}$ and $\underline{\underline{\Lambda}}$ result as*
$$(1.85)$$
$$\underline{\underline{R}} = \left( \begin{array}{ccc} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 1 & 1 \end{array} \right), \quad \underline{\underline{\Lambda}} = \left( \begin{array}{ccc} -1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{array} \right), \quad \underline{\underline{R}}^{-1} = \left( \begin{array}{ccc} 0 & -1 & 1 \\ 1 & 1 & -1 \\ -1 & 0 & 1 \end{array} \right).$$

FIGURE 5. The mesh used to compute the numerical solution.

## 5. Numerical methods for the linear scalar advection equation

The aim of this section is to suggest some numerical methods to solve numerically the linear advection equation (PDE), given the initial (IC) and boundary (BC) conditions

$$
\begin{aligned}
&\text{PDE:} && q_t + aq_x = 0, \quad x \in [x_L; x_R], t \in \mathbb{R}_0^+, q \in \mathbb{R}, a \in \mathbb{R}, \\
&\text{IC:} && q(x, 0) = h(x), \\
&\text{BC:} && q(x_L, t) = b_L(t), \quad q(x_R, t) = b_R(t).
\end{aligned}
$$
(1.86)

For the numerical solution of (1.86) we first proceed with the *discretization* of space and time by introducing a *mesh* composed of *discrete points* in space and time as follows:

$$
x_i = x_L + i\Delta x, \quad i = 0, 1, ...
$$
(1.87)

$$
t^n = n\Delta t, \quad n = 0, 1, ...
$$
(1.88)

where $\Delta x$ is the interval length of the mesh or the *mesh spacing* and $\Delta t$ is the time step. The indices $i$ and $n$ identify a node of the mesh in space and time, as depicted in Fig. 5. Let us introduce the short notation for the numerical solution, defined by

$$
q_i^n = q(x_i, t^n).
$$
(1.89)

First, we consider the *finite difference* method to solve (1.86). A simple finite difference applied to the time derivative in (1.86) results to be

$$
q_t(x_i, t^n) \approx \frac{q_i^{n+1} - q_i^n}{\Delta t},
$$
(1.90)

which is obviously consistent with the *definition* of the derivative

$$
q_t(x_i, t^n) = \lim_{\Delta t \to 0} \frac{q(x_i, t^n + \Delta t) - q(x_i, t^n)}{\Delta t}.
$$
(1.91)

While the continuous derivative (1.91) uses *infinitesimal* differences, the discrete derivative (1.90) uses *finite* differences.

FIGURE 6. The discretization of the derivatives in space following the finite difference approach.

The particular finite difference (1.90) applied to a time derivative is also called *first order Euler* time discretization. To discretize the derivative in space there are at least the following three different possibilities, as sketched in Fig. 6:

- The so called *backward difference*

$$(1.92) \qquad q_x\left(x_i, t^n\right) \approx \frac{q_i^n - q_{i-1}^n}{\Delta x}.$$

- The *forward difference*

$$(1.93) \qquad q_x\left(x_i, t^n\right) \approx \frac{q_{i+1}^n - q_i^n}{\Delta x}.$$

- The *central difference*

$$(1.94) \qquad q_x\left(x_i, t^n\right) \approx \frac{q_{i+1}^n - q_{i-1}^n}{2\Delta x}.$$

**5.1. Some basic explicit finite difference schemes.** It follows that three different numerical methods are possible to solve the linear advection equation (1.86) using the explicit Euler time discretization:

- The first scheme is based on the *backward difference* discretization of the space derivative

$$(1.95) \qquad q_i^{n+1} = q_i^n - a\frac{\Delta t}{\Delta x}\left(q_i^n - q_{i-1}^n\right),$$

Note that according to the sign of the constant velocity $a$, we obtain an *upwind* scheme if $a > 0$ and a *downwind* method with $a < 0$.

- The second method applies the *forward difference* discretization in space

$$(1.96) \qquad q_i^{n+1} = q_i^n - a\frac{\Delta t}{\Delta x}\left(q_{i+1}^n - q_i^n\right).$$

On the contrary to the approach above, the scheme (1.96) leads to a *downwind* method if $a > 0$ and an *upwind* scheme if $a < 0$.

FIGURE 7. Visualization of the numerical methods derived from
the finite difference method applying the three different approaches
to discretize the space derivatives.

- The last scheme uses the *central difference* approach in space

$$(1.97) \qquad q_i^{n+1} = q_i^n - a\frac{\Delta t}{2\Delta x}\left(q_{i+1}^n - q_{i-1}^n\right).$$

A sketch of the upwind or downwind property of the above numerical methods
in function of the sign of the propagation speed $a$ is shown in Fig. 7.

As we will see later, not all of the above schemes are useful in practice since
not all of the above schemes are stable.

To assure the stability of an explicit method for solving a linear evolution-
ary partial differential equation a constraint on the time step $\Delta t$ is required. It
derives from the Courant-Friedrichs-Lewy (CFL) condition [18], that in the one

FIGURE 8. The plot of the CFL condition on the x-t plane.

dimensional case is defined by

$$(1.98) \qquad c \leq 1, \quad \text{with } c = a \frac{\Delta t}{\Delta x},$$

where $c$ is the so-called CFL number. The CFL condition leads to a restriction on the maximum admissible numerical time step at fixed grid size. The physical meaning of the CFL number is the measure for the progress of a disturbance over a time step $\Delta t$ related to the grid distance $\Delta x$. If the domain of dependence of the PDE due to the physical celerity is not contained in the domain of dependence of the discrete scheme given by the CFL number, then the numerical scheme will be unstable. This means that numerical information propagates slower than the physical wave, as shown in Fig. 8.

**5.2. The Lax-Friedrichs method.** However, the CFL condition is not sufficient for stability. In fact, via a von Neumann stability analysis shown later in this manuscript, it is possible to prove that *all* the central and downwind methods are unconditionally unstable for any CFL number $c$. Hence, Lax and Friedrichs [**52**] proposed a modification of the original central finite difference method replacing $q_i^n$ in (1.97) by the average of the backward and forward values, as follows

$$(1.99) \qquad q_i^{n+1} = \frac{1}{2} \left( q_{i+1}^n + q_{i-1}^n \right) - a \frac{1}{2} \frac{\Delta t}{\Delta x} \left( q_{i+1}^n - q_{i-1}^n \right).$$

It can be proven that the scheme (1.99) is stable.

**5.3. The Lax-Wendroff method.** In this section we present the scheme introduced by Lax and Wendroff [**53**]. Let us consider the Taylor series in time of the function $q_i^{n+1}$ computed at the point $(x_i, t^n)$

$$(1.100) \qquad q_i^{n+1} = q_i^n + \Delta t q_t + \frac{\Delta t^2}{2} q_{tt} + \dots$$

The derivatives in time are obtained from the original PDE, in this case the linear advection equation (1.86), through the so-called Lax–Wendroff or Cauchy–Kovalewski procedure, as follows

$$(1.101) \qquad q_t = -a q_x,$$

$$(1.102) \qquad q_{tt} = -a q_{xt}, \quad q_{xt} = -a q_{xx} \quad \Rightarrow q_{tt} = a^2 q_{xx}.$$

Afterwards, the derivatives in space are approximated using the central finite difference approach

$$(1.103) \qquad q_x \approx \frac{q_{i+1}^n - q_{i-1}^n}{2\Delta x}, \quad q_{xx} \approx \frac{q_{i+1}^n - 2q_i^n + q_{i-1}^n}{\Delta x^2}.$$

Finally, we substitute the spatial derivatives above into the Taylor expansion and obtain the second order Lax-Wendroff scheme

$$(1.104) \qquad q_i^{n+1} = q_i^n - \frac{a}{2}\frac{\Delta t}{\Delta x}\left(q_{i+1}^n - q_{i-1}^n\right) + \frac{a^2 \Delta t^2}{2\Delta x^2}\left(q_{i+1}^n - 2q_i^n + q_{i-1}^n\right).$$

**5.4. The method of lines (MOL).** This technique for solving partial differential equations dates back to the 1960ies [**71**] and consists of a semi-discrete numerical method. It proceeds by first discretizing the spatial derivatives only (e.g. through the backward finite difference method) and leaving the time variable still continuous

$$(1.105) \qquad \frac{\partial q}{\partial t} \approx -a\frac{q_i^n - q_{i-1}^n}{\Delta x}.$$

This leads to a system of ordinary differential equations ODE in which the unknown quantities are the $q_i$ at the grid points:

$$(1.106) \qquad \frac{d}{dt}\vec{q}(t) = \vec{f}(\vec{q}(t)) \quad \text{with } \vec{q}(t) = \begin{pmatrix} q_1(t) \\ q_2(t) \\ ... \\ q_N(t) \end{pmatrix}.$$

Here, $N$ denotes the number of the nodes on the mesh. To discretize the system (1.106) in time we can use any numerical method for solving ODE, such as the explicit or implicit first order Euler schemes, the second order accurate semi–implicit Crank-Nicholson method, the explicit second order scheme of Heun, any higher order Runge-Kutta scheme or any method from the family of Adams–Bashforth and Adams–Moulton–type multi–stage methods.

**5.5. Implicit Euler method together with the backward finite difference discretization.** Applying the implicit Euler scheme coupled with the backward finite difference method to (1.106) we obtain the following discrete equation

$$(1.107) \qquad q_i^{n+1} = q_i^n - a\frac{\Delta t}{\Delta x}\left(q_i^{n+1} - q_{i-1}^{n+1}\right).$$

This leads to a tridiagonal system of $N$ equations, as follows

$$(1.108) \qquad q_i^{n+1}(1+c) - cq_{i-1}^n = q_i^n, \quad \text{for } i = 1, .., N$$

that can be written in matrix form as

$$(1.109) \qquad \begin{pmatrix} c+1 & 0 & 0 & ... & 0 \\ -c & c+1 & 0 & ... & 0 \\ 0 & -c & c+1 & ... & 0 \\ ... & ... & ... & ... & ... \\ 0 & 0 & ... & -c & c+1 \end{pmatrix} \begin{pmatrix} q_1^{n+1} \\ q_2^{n+1} \\ q_3^{n+1} \\ ... \\ q_N^{n+1} \end{pmatrix} = \begin{pmatrix} q_1^n + cb_L \\ q_2^n \\ q_3^n \\ ... \\ q_N^n \end{pmatrix}.$$

Here, we supposed $a > 0$ to introduce the left boundary condition.

**5.6. The Thomas algorithm.** The Thomas algorithm consists of a simplified form of the Gauss algorithm and is used to solve tridiagonal system of equations

$$(1.110) \qquad \underline{\underline{A}} \cdot \vec{x} = \vec{b},$$

where $\vec{x}$ is the vector of the unknowns. The tridiagonal matrix $\underline{\underline{A}}$ and the right–hand–side vector $\vec{b}$ are given by

$$(1.111)$$
$$\underline{\underline{A}} = \begin{pmatrix} d_1 & u_1 & 0 & 0 & ... & 0 & 0 & 0 & 0 \\ l_2 & d_2 & u_2 & 0 & ... & 0 & 0 & 0 & 0 \\ ... & ... & ... & ... & ... & ... & ... & ... & ... \\ 0 & 0 & 0 & 0 & ... & 0 & l_{N-1} & d_{N-1} & u_{N-1} \\ 0 & 0 & 0 & 0 & ... & 0 & 0 & l_N & d_N \end{pmatrix}, \quad \vec{b} = \begin{pmatrix} b_1 \\ b_2 \\ ... \\ b_{N-1} \\ b_N \end{pmatrix},$$

where $N$ is the number of unknowns. The first step of the Thomas algorithm consists of a forward elimination:

$$(1.112) \qquad \begin{aligned} & u_1 = u_1/d_1 \\ & b_1 = b_1/d_1 \\ & \text{for } i = 2 : N \\ & \quad g = 1/(d_i - u_{i-1}l_i) \\ & \quad u_i = u_i \cdot g \\ & \quad b_i = g \cdot (b_i - l_i b_{i-1}) \\ & \text{end} \end{aligned}$$

Then, a backward substitution leads to the final solution of the system:

$$(1.113) \qquad \begin{aligned} & x_N = b_N \\ & \text{for } i = N - 1 : -1 : 1 \\ & \quad x_i = b_i - u_i x_{i+1} \\ & \text{end} \end{aligned}$$

**5.7. Implicit Euler method together with the central finite difference discretization.** To solve the advection equation (1.106) we could also use the implicit Euler scheme together with a central finite difference method to compute the derivative in space

$$(1.114) \qquad q_i^{n+1} = q_i^n - \frac{a\Delta t}{2\Delta x}\left(q_{i+1}^{n+1} - q_{i-1}^{n+1}\right),$$

which can also be written as

$$(1.115) \qquad \frac{c}{2}q_{i+1}^{n+1} + q_i^{n+1} - \frac{c}{2}q_{i-1}^n = q_i^n, \quad \text{for } i = 1, .., N$$

In matrix form this results as

$$(1.116)$$
$$\underline{\underline{A}} = \begin{pmatrix} 1 & \frac{c}{2} & 0 & 0 & ... & 0 & 0 & 0 & 0 \\ -\frac{c}{2} & 1 & \frac{c}{2} & 0 & ... & 0 & 0 & 0 & 0 \\ ... & ... & ... & ... & ... & ... & ... & ... & ... \\ 0 & 0 & 0 & 0 & ... & 0 & -\frac{c}{2} & 1 & \frac{c}{2} \\ 0 & 0 & 0 & 0 & ... & 0 & 0 & 1 & \frac{c}{2} \end{pmatrix} \begin{pmatrix} q_1^{n+1} \\ q_2^{n+1} \\ ... \\ q_{N-1}^{n+1} \\ q_N^{n+1} \end{pmatrix} = \begin{pmatrix} q_1^n + \frac{c}{2}b_L \\ q_2^n \\ ... \\ q_{N-1}^n \\ q_N^n - \frac{c}{2}b_R \end{pmatrix}.$$

The linear tridiagonal system of equations above can be solved again using the Thomas algorithm presented in the previous section.

EXAMPLE 6. *In the following we derive the Beam-Warming method. It is based on the same idea as the Lax-Wendroff scheme, but it uses an upwind stencil, i.e. it uses the points $x_i$, $x_{i-1}$ and $x_{i-2}$, assuming $a > 0$. We compute the finite difference approximations of the first and second space derivative based on these*

*three points using a Taylor series expansion in space computed at the point $(x_i, t^n)$, as follows:*

$$(1.117) \qquad q_{i-2}^n = q_i^n - 2\Delta x q_x + \frac{4\Delta x^2}{2} q_{xx} - ...,$$

$$(1.118) \qquad q_{i-1}^n = q_i^n - \Delta x q_x + \frac{\Delta x^2}{2} q_{xx} - ....$$

*From (1.117) and (1.118) one obtains the finite difference expressions for the first and second spatial derivative as*

$$(1.119) \qquad q_x = \frac{q_{i-2}^n - 4q_{i-1}^n + 3q_i^n}{2\Delta x},$$

$$(1.120) \qquad q_{xx} = \frac{q_{i-2}^n - 2q_{i-1}^n + q_i^n}{\Delta x^2}.$$

*Like in the Lax–Wendroff scheme the solution $q_i^{n+1}$ at the new time $t^{n+1}$ is given by a temporal Taylor series expansion as*

$$(1.121) \qquad q_i^{n+1} = q_i^n + \Delta t q_t + \frac{\Delta t^2}{2} q_{tt} + ...$$

*Applying the Cauchy-Kovalewski procedure to the linear advection equation we obtain again*

$$(1.122) \qquad q_t = -a q_x \qquad \text{and} \qquad q_{tt} = +a^2 q_{xx}.$$

*Replacing the time derivatives in the time Taylor series by space derivatives and inserting the finite difference approximations shown above leads to the Beam-Warming method*

$$(1.123) \quad q_i^{n+1} = q_i^n - \frac{a\Delta t}{2\Delta x} \left( q_{i-2}^n - 4q_{i-1}^n + 3q_i^n \right) + \frac{a^2\Delta t^2}{2\Delta x^2} \left( q_{i-2}^n - 2q_{i-1}^n + q_i^n \right).$$

EXAMPLE 7. *The method of Fromm [33] is based on an arithmetic average of the Lax-Wendroff method and the Beam-Warming scheme. In this approach, the derivatives in space are given by an arithmetic average of the backward finite differences based on the stencil $x_i$, $x_{i-1}$ and $x_{i-2}$ and the central finite difference. We obtain*

$$(1.124) \qquad q_{xx} = \frac{1}{2} \frac{q_{i-2}^n - 2q_{i-1}^n + q_i^n}{\Delta x^2} + \frac{1}{2} \frac{q_{i+1}^n - 2q_i^n + q_{i-1}^n}{\Delta x^2}.$$

*and*

$$(1.125) \qquad q_x = \frac{1}{2} \frac{q_{i-2}^n - 4q_{i-1}^n + 3q_i^n}{2\Delta x} + \frac{1}{2} \frac{q_{i+1}^n - q_{i-1}^n}{2\Delta x}.$$

*Based on the same ideas as the ones used for deriving the Lax–Wendroff and the Beam–Warming schemes we obtain Fromm's method:*

(1.126)
$$q_i^{n+1} = q_i^n - \frac{a\Delta t}{4\Delta x} \left( q_{i-2}^n - 5q_{i-1}^n + 3q_i^n + q_{i+1}^n \right) + \frac{a^2\Delta t^2}{4\Delta x^2} \left( q_{i-2}^n - q_{i-1}^n - q_i^n + q_{i+1}^n \right).$$

*It is very interesting to note that the Fromm scheme is the least dispersive and hence the least oscillatory linear second order scheme.*

# Analysis of Error, Stability and Monotonicity

## 1. Local discretization error

All the explicit finite difference (FD) methods presented in the previous section can be written in the compact form

$$(2.1) \qquad q_i^{n+1} = H\left(q_{i-l}^n, q_{i-l+1}^n, ..., q_{i-1}^n, q_i^n, q_{i+1}^n, ..., q_{i+r-1}^n, q_{i+r}^n\right),$$

where the operator $H$ is a real function of $(l + r + 1)$ variables and $l$ and $r$ are integer numbers that define the size of the *stencil* (support) of the numerical method.

EXAMPLE 8. *We consider the upwind finite difference method*

$$(2.2) \qquad \begin{array}{ll} q_i^{n+1} = q_i^n - \frac{a\Delta t}{\Delta x}\left(q_i^n - q_{i-1}^n\right), & \text{for } a > 0, \\ q_i^{n+1} = q_i^n - \frac{a\Delta t}{\Delta x}\left(q_{i+1}^n - q_i^n\right), & \text{for } a < 0. \end{array}$$

*The two different discretizations above can be written together in the following unified form*

$$(2.3)$$
$$q_i^{n+1} = q_i^n - \frac{\Delta t}{\Delta x}\left(\frac{a - |a|}{2}q_{i+1}^n - \frac{a - |a|}{2}q_i^n\right) - \frac{\Delta t}{\Delta x}\left(\frac{a + |a|}{2}q_i^n - \frac{a + |a|}{2}q_{i-1}^n\right),$$

*or shorter as*

$$(2.4) \qquad q_i^{n+1} = q_i^n - \frac{\Delta t}{\Delta x}a^-\left(q_{i+1}^n - q_i^n\right) - \frac{\Delta t}{\Delta x}a^+\left(q_i^n - q_{i-1}^n\right),$$

*where the $a^\pm$ are defined by $a^\pm = \frac{1}{2}\left(a \pm |a|\right)$. Hence, we have*

$$(2.5) \quad H = H\left(q_{i-1}^n, q_i^n, q_{i+1}^n\right) = \left(1 + \frac{\Delta t}{\Delta x}(a^- - a^+)\right)q_i^n + \frac{\Delta t}{\Delta x}a^+q_{i-1}^n - \frac{\Delta t}{\Delta x}a^-q_{i+1}^n.$$

**1.1. Linear methods.** A numerical method is called *linear* if it can be written in the following form

$$(2.6) \qquad q_i^{n+1} = \sum_{j=-l}^{r} b_j q_{i+j}^n,$$

where the coefficients $b_j$ do *not* depend on the solution $q_i^n$, but only on the parameters of the PDE and its discretization, i.e. on the advection speed $a$, on the mesh length $\Delta x$ and on the time step $\Delta t$.

**1.2. Operator form.** It is often useful to define a general numerical method as an operator $L_h$ acting on the discrete solution $q_h = q_h(x_h, t_h)$ in the form

$$(2.7) \qquad L_h\left(q_h\right) = 0.$$

Using the notation (2.1) of a general explicit scheme we have

$$(2.8) \qquad L_h\left(q_h\left(x_h, t_h\right)\right) = q_i^{n+1} - H\left(q_{i-l}^n, ..., q_{i-1}^n, q_i^n, q_{i+1}^n, ..., q_{i+r}^n\right) = 0,$$

where $x_h = (x_{i-l}, ..., x_{i-1}, x_i, x_{i+1}, ..., x_{i+r})$ and $t_h = \left(t^n, t^{n+1}\right)$ are the discretizations of space and time, respectively, and the numerical solution $q_h\left(x_h, t_h\right)$ is defined in the discrete points $(x, t) \in x_h \times t_h$.

**1.3. Local truncation error.** The local truncation error or the *local discretization error* is due to the truncation of the infinite Taylor series to form the discrete algorithm. It refers to one time step and is defined by

$$\tag{2.9} \tau = \frac{L_h\left(q\left(x_h, t_h\right)\right)}{\Delta t},$$

where $q\left(x_h, t_h\right)$ is the *exact solution* of the PDE computed at the discrete point $\left(x_h, t_h\right)$. Here, the exact solution $q\left(x, t\right)$ is expanded using a Taylor series in space and time with respect to the point $\left(x_i, t^n\right)$.

DEFINITION 1. *A method is called k-th order accurate in space and n-th order accurate in time if*

$$\tag{2.10} \tau = \mathcal{O}(\Delta x^k) + \mathcal{O}(\Delta t^n).$$

EXAMPLE 9. *We consider the upwind finite difference method applied to the linear scalar advection equation assuming $a > 0$. The related numerical operator is given by*

$$\tag{2.11} L_h\left(q_h\left(x_h, t_h\right)\right) = q_i^{n+1} - q_i^n + \frac{a\Delta t}{\Delta x}\left(q_i^n - q_{i-1}^n\right) = 0.$$

*The exact solution is given by the Taylor series expansion as follows:*

$$\tag{2.12} q_i^{n+1} = q_i^n + \Delta t q_t + \frac{\Delta t^2}{2}q_{tt} + \mathcal{O}(\Delta t^3),$$

$$\tag{2.13} q_{i-1}^n = q_i^n - \Delta x q_x + \frac{\Delta x^2}{2}q_{xx} - \mathcal{O}(\Delta x^3).$$

*Now, we compute the local truncation error*

$$\tag{2.14} \begin{aligned} \tau = \quad & \frac{1}{\Delta t}\left(q_i^n + \Delta t q_t + \frac{\Delta t^2}{2}q_{tt} + \mathcal{O}\left(\Delta t^3\right) - q_i^n\right) + \\ & \frac{1}{\Delta t}\frac{a\Delta t}{\Delta x}\left(q_i^n - q_i^n + \Delta x q_x - \frac{\Delta x^2}{2}q_{xx} - \mathcal{O}\left(\Delta x^3\right)\right), \end{aligned}$$

*which results in*

$$\tag{2.15} \tau = \overbrace{q_t + a q_x}^{=0} + \frac{\Delta t}{2}q_{tt} - \frac{a\Delta x}{2}q_{xx} + \mathcal{O}\left(\Delta x^2\right) + \mathcal{O}\left(\Delta t^2\right).$$

*Here, the first term vanishes because the exact solution $q$ obviously must satisfy the PDE exactly. We therefore obtain the final result*

$$\tag{2.16} \tau = \frac{\Delta t}{2}q_{tt} - \frac{a\Delta x}{2}q_{xx} + \mathcal{O}\left(\Delta x^2\right) + \mathcal{O}\left(\Delta t^2\right) = \mathcal{O}\left(\Delta t\right) + \mathcal{O}\left(\Delta x\right).$$

*From the leading error terms (lowest exponents of $\Delta x$ and $\Delta t$ in the above equation) we find that they are proportional to $\Delta x^1$ and $\Delta t^1$, hence the upwind finite difference scheme is of first order of accuracy in space and time.*

EXAMPLE 10. *In the following, we compute the local truncation error of the central finite difference scheme. The numerical operator is defined by*

$$\tag{2.17} L_h = q_i^{n+1} - q_i^n + \frac{a\Delta t}{2\Delta x}\left(q_{i+1}^n - q_{i-1}^n\right) = 0.$$

*We apply the Taylor series expansion in space to the term $q_{i+1}^n$*

$$\tag{2.18} q_{i+1}^n = q_i^n + \Delta x q_x + \frac{\Delta x^2}{2}q_{xx} + \frac{\Delta x^3}{6}q_{xxx} + ...$$

*The same procedure is carried out for $q_n^{n+1}$ and $q_{i-1}^n$, see Eqns. (2.12) and (2.13). Then, the local truncation error related to the central FD method results as*

$$
\begin{aligned}
(2.19) \quad \tau = \ & \frac{1}{\Delta t}\left(q_i^n + \Delta t q_t + \frac{\Delta t^2}{2}q_{tt} + \mathcal{O}\left(\Delta t^3\right) - q_i^n\right) + \\
& \frac{1}{\Delta t}\frac{a\Delta t}{2\Delta x}\left(+q_i^n + \Delta x q_x + \frac{\Delta x^2}{2}q_{xx} + \frac{\Delta x^3}{6}q_{xxx} + \mathcal{O}\left(\Delta x^4\right)\right) + \\
& \frac{1}{\Delta t}\frac{a\Delta t}{2\Delta x}\left(-q_i^n + \Delta x q_x - \frac{\Delta x^2}{2}q_{xx} + \frac{\Delta x^3}{6}q_{xxx} - \mathcal{O}\left(\Delta x^4\right)\right).
\end{aligned}
$$

*Simplifying the relation above we obtain*

$$
(2.20) \quad \tau = \overbrace{q_t + aq_x}^{=0} + \frac{\Delta t}{2}q_{tt} + \frac{a\Delta x^2}{6}q_{xxx} + \mathcal{O}\left(\Delta x^3\right) + \mathcal{O}\left(\Delta t^2\right) = \mathcal{O}\left(\Delta t\right) + \mathcal{O}\left(\Delta x^2\right).
$$

*Hence, the central finite difference scheme is of first order of accuracy in time and of second order in space. However, it is not stable.*

EXAMPLE 11. *In the case of the Lax-Friedrichs scheme (1.99) the numerical operator is given by*

$$
(2.21) \quad L_h = q_i^{n+1} - \frac{1}{2}\left(q_{i+1}^n + q_{i-1}^n\right) + \frac{a\Delta t}{2\Delta x}\left(q_{i+1}^n - q_{i-1}^n\right) = 0.
$$

*We compute the local truncation error using Taylor series expansions of the exact solution in space and time, as in the examples above:*

$$
\begin{aligned}
(2.22) \quad \tau = \ & \frac{1}{\Delta t}\left(q_i^n + \Delta t q_t + \frac{\Delta t^2}{2}q_{tt} + \mathcal{O}\left(\Delta t^3\right)\right) - \\
& \frac{1}{\Delta t}\frac{1}{2}\left(q_i^n + \Delta x q_x + \frac{\Delta x^2}{2}q_{xx} + \mathcal{O}\left(\Delta x^3\right) + q_i^n - \Delta x q_x + \frac{\Delta x^2}{2}q_{xx} + \mathcal{O}\left(\Delta x^3\right)\right) + \\
& \frac{1}{\Delta t}\frac{a\Delta t}{2\Delta x}\left(q_i^n + \Delta x q_x + \frac{\Delta x^2}{2}q_{xx} + \mathcal{O}\left(\Delta x^3\right) - q_i^n + \Delta x q_x - \frac{\Delta x^2}{2}q_{xx} + \mathcal{O}\left(\Delta x^3\right)\right).
\end{aligned}
$$

*After some manipulations, using the definition of the Courant number $c = a\Delta t/\Delta x$ the local truncation error becomes*

$$
(2.23) \quad \tau = \overbrace{q_t + aq_x}^{=0} + \frac{\Delta t}{2}q_{tt} - \frac{a\Delta x}{2c}q_{xx} + \mathcal{O}\left(\Delta x^2\right) + \mathcal{O}\left(\Delta t^2\right) = \mathcal{O}\left(\Delta t\right) + \mathcal{O}\left(\Delta x\right).
$$

*This proves that the Lax-Friedrichs scheme is of first order of accuracy in space and in time.*

EXERCISE 2. *Compute the local truncation error of the Lax–Wendroff scheme, the Beam–Warming method and the scheme of Fromm.*

## 2. The modified equation approach or the method of differential approximation

In general, the numerical solution $q_h$ does *not* satisfy the original PDE

$$
(2.24) \quad \frac{\partial q}{\partial t} + a\frac{\partial q}{\partial x} = 0
$$

exactly, but only approximately, i.e. in general we have

$$
(2.25) \quad \frac{\partial q_h}{\partial t} + a\frac{\partial q_h}{\partial x} \neq 0.
$$

To quantify the errors, we will consider the so–called *modified equation* or *equivalent differential equation*

$$
(2.26) \quad \frac{\partial q_h}{\partial t} + a\frac{\partial q_h}{\partial x} = \sum_{l=2}^{\infty} c_l \frac{\partial^l q_h}{\partial x_l},
$$

which is solved *exactly* by the numerical method. A detailed analysis of the modified equation (2.26) gives very important insight into the numerical scheme, in particular concerning its stability, the accuracy and its diffusive and dispersive behaviour. The

error coefficients $c_l$ in equation (2.26) are coefficients depending on the advection speed $a$, on the mesh size $\Delta x$ and on the time step.

**2.1. Computation of the modified equation.** The modified equation analysis or the so–called method of differential approximation goes back to the work of the Russian mathematician Y.I. Shokin [**73**] and the US mathematicians Warming and Hyett [**90**]. From the general form of the numerical scheme

$$(2.27) \qquad\qquad q_i^{n+1} = H(q_{i-l}^n, ..., q_i^n, ..., q_{i+r}^n),$$

the modified equation is computed by expanding the discrete solution $q_h$ at the discrete points $x_h$ and $t_h$ in a Taylor series in space and time with respect to point $x_i$ and $t^n$, similar to the local truncation error analysis, i.e. we obtain

$$(2.28)$$
$$q + \Delta t q_t + \frac{1}{2}\Delta t^2 q_{tt} + \frac{1}{6}\Delta t^3 q_{ttt} + ... =$$
$$H(q - l\Delta x q_x + \frac{1}{2}(l\Delta x)^2 q_{xx} - ..., ..., q, ..., q + r\Delta x q_x + \frac{1}{2}(r\Delta x)^2 q_{xx} + ...).$$

Eqn. (2.28) is the so–called $\Gamma$–form of the differential approximation of the scheme (2.27), where still higher order time derivatives appear. In order to obtain the final form of the differential approximation, we must then use the $\Gamma$–form (2.28) in order to replace time derivatives and mixed space–time derivatives by pure space derivatives using successive differentiation of (2.28) with respect to $x$ and $t$. This procedure, when applied to Eqn. (2.28) is the so–called Warming–Hyett procedure [**90**], which is similar to the Cauchy-Kovalewski procedure, where successive differentiation with respect to $x$ and $t$ is applied to the original PDE. Due to the increased complexity of the PDE (2.28) with respect to the original PDE, the Warming–Hyett procedure is algebraically *much* more complicated than the Cauchy–Kovalewski procedure. To get expressions which can be handled, one must determine a priori the maximum order of $\Delta x$ and $\Delta t$ in the modified equation and all higher order terms are immediately dropped as soon as they appear during the Warming–Hyett procedure. The final result of the method of differential approximation is then the so–called $\Pi$–form of the differential approximation (2.26), which contains the original PDE on the left hand side and the error terms involving only pure space derivatives on the right hand side:

$$\frac{\partial q_h}{\partial t} + a\frac{\partial q_h}{\partial x} = \sum_{l=2}^{\infty} c_l \frac{\partial^l q_h}{\partial x_l}.$$

**2.2. Physical interpretation of the error.** The terms on the right–hand side of (2.26) provide in particular a very interesting *physical* interpretation of the error using the harmonic analysis of Fourier applied to (2.26). We consider the evolution of one single Fourier mode (monochromatic wave) as follows:

$$(2.29) \qquad q\,(x,t) = q_0 e^{i(kx-\omega t)}, \quad \text{with} \quad k = \frac{2\pi}{\lambda}, \ \omega = \frac{2\pi}{T}, \ i^2 = -1,$$

where $k$ denotes the wave number and $\omega$ is the angular frequency. Its derivatives in space and in time are the following

$$(2.30) \qquad\qquad q_t = -i\omega q, \quad q_x = ikq, \quad \frac{\partial^l q}{\partial x^l} = (ik)^l \, q.$$

Substitution of these relations into (2.26) leads to the so-called *dispersion relation* of the PDE (2.26)

$$(2.31) \qquad\qquad -i\omega + aik = \sum_{l=2}^{\infty} c_l \, (ik)^l \, .$$

Then, the angular frequency $\omega$ is defined by

$$(2.32) \qquad \omega = ak + i \sum_{l=2}^{\infty} c_l \, (ik)^l, \quad \text{with } i^l = \begin{cases} (-1)^m, & \text{if } l = 2m, \\ i \, (-1)^m, & \text{if } l = 2m + 1. \end{cases}$$

Hence, we obtain the final expression of the dispersion relation:

$$(2.33) \qquad \omega = ak + i \sum_{m=1}^{\infty} (-1)^m \, c_{2m} k^{2m} - \sum_{m=1}^{\infty} (-1)^m \, c_{2m+1} k^{2m+1}.$$

The *exact* solution of (2.26) for one single Fourier mode is given by

$$(2.34) \quad q\,(x,t) = q_0 e^{ik \cdot \left( x - t \left( a - \overbrace{\sum_{m=1}^{\infty} (-1)^m \, c_{2m+1} k^{2m}}^{\text{dispersion error}} \right) \right)} \cdot e^{t \cdot \overbrace{\sum_{m=1}^{\infty} (-1)^m \, c_{2m} k^{2m}}^{\text{diffusion error}}},$$

where the dispersion and diffusion error terms can be clearly identified. Relation (2.34) can be written as

$$(2.35) \qquad q\,(x,t) = q_0 \cdot e^{-dt} \cdot e^{i(k(x - v(k)t))},$$

where $d$ represents the diffusion error and $v(k)$ is the wave speed, which in this case depends on the wavenumber $k$.

The physical phenomenon of *dispersion* means that the wave speed is a function of the wavenumber $k$. An example of a dispersive phenomenon is the formation of a rainbow, which is caused by the separation of white light into different components of different wavelengths (different colors) due to optical dispersion. In our case, the velocity of one single Fourier mode is given by

$$(2.36) \qquad v\,(k) = a - \sum_{m=1}^{\infty} (-1)^m \, c_{2m+1} k^{2m},$$

instead of being

$$(2.37) \qquad v\,(k) = a,$$

which would be the correct wave speed for the original linear scalar advection equation (3.1). We note in particular that dispersive errors are introduced only by the coefficients $c_l$ in front of the *odd* derivatives in the errors terms of the equivalent PDE, i.e. by the coefficients $c_3$, $c_5$, etc. In contrast to dispersion, which modifies the wave speed but not the amplitude, the diffusion error causes a reduction or amplification of the amplitude of the wave during its propagation, but does not have any influence on the wave speed. From (2.34) we clearly see that the two effects are separated. The diffusion error is quantified by the term

$$(2.38) \qquad d = - \sum_{m=1}^{\infty} (-1)^m \, c_{2m} k^{2m},$$

which must be positive or zero for a scheme to be stable, i.e. the amplitudes of all the Fourier modes must be *non–increasing*. As we can clearly see from this equation, the diffusion error is only governed by the error coefficients $c_l$ in front of the *even* spatial derivatives, i.e. by the coefficients $c_2$, $c_4$, etc. It is important to note the alternating sign in Eqn. (2.38), which means that for stability we must have $c_2 \geq 0$ but $c_4 \leq 0$ and so on.

EXAMPLE 12. *We consider the backward finite difference (FD) method (1.95) applied to the linear advection equation. Let us assume $a > 0$, hence, the*

*backward FD scheme is upwind. Using Taylor series expansions of $q_{i-1}^n$ and $q_i^{n+1}$ we obtain*

(2.39)

$$q_i^n + \Delta t q_t + \frac{\Delta t^2}{2} q_{tt} + \mathcal{O}\left(\Delta t^3\right) = q_i^n - \frac{a\Delta t}{\Delta x}\left(q_i^n - q_i^n + \Delta x q_x - \frac{\Delta x^2}{2} q_{xx} + \mathcal{O}\left(\Delta x^3\right)\right).$$

*Neglecting the terms of higher oder in (2.39), we obtain the $\Gamma$–form of the differential approximation*

(2.40)
$$q_t + a q_x = \frac{a\Delta x}{2} q_{xx} - \frac{\Delta t}{2} q_{tt}.$$

*To obtain the $\Pi$–form (2.26) from (2.40) we compute the second derivative in time $q_{tt}$ via the procedure of Warming-Hyett deriving (2.40) in time and in space:*

(2.41)
$$\begin{aligned} q_{tt} &= -a q_{xt} + \frac{a\Delta x}{2} q_{xxt} - \frac{\Delta t}{2} q_{ttt} \\ q_{tx} &= -a q_{xx} + \frac{a\Delta x}{2} q_{xxx} - \frac{\Delta t}{2} q_{ttx}. \end{aligned}$$

*Substituting $q_{xt}$ into the relation for the second derivative in time we can write $q_{tt}$ in terms of pure spatial derivatives*

(2.42)
$$q_{tt} = a^2 q_{xx} + \mathcal{O}\left(\Delta x, \Delta t\right).$$

*Here, we stop at zeroth order terms in the Warming–Hyett procedure since we are only interested in the leading error term of the scheme. If one wanted to obtain also the higher order error terms in the differential approximation of the scheme, then also higher order terms would have to be retained during the Warming–Hyett procedure. We obtain*

(2.43)
$$q_t + a q_{xx} = \frac{a\Delta x}{2} q_{xx} - \frac{a^2\Delta t}{2} q_{xx} = \frac{a\Delta x}{2}\left(1 - c\right) q_{xx},$$

*hence the leading error coefficient of the modified equation is*

$$c_2 = \frac{a\Delta x}{2}\left(1 - c\right), \quad \text{with the Courant number } c = \frac{a\Delta t}{\Delta x}.$$

*For the solution to be stable, the coefficient must be positive, i.e. we obtain the following result:*

(2.44)
$$c_2 \geq 0 \quad \Rightarrow \quad c \leq 1,$$

*which is the classical Courant–Friedrichs–Lewy stability condition [**18**]. We find that the scheme is first order accurate and that the diffusion error is the smaller the larger the Courant number (within the stability limit, of course). For $c = 1$ the error term vanishes, and one can show that also all the higher order error terms vanish. In this very particular case of $c = 1$, the upwind scheme reproduces the exact solution of the PDE.*

EXAMPLE 13. *Now, we study the forward finite difference (FD) method (1.96) for the linear scalar advection equation with $a > 0$, leading to a downwind scheme. Using appropriate Taylor series expansions with respect to the point $(x_i, t^n)$ we obtain*

(2.45)
$$q_i^n + \Delta t q_t + \frac{\Delta t^2}{2} q_{tt} = q_i^n - \frac{a\Delta t}{\Delta x}\left(q_i^n + \Delta x q_x - \frac{\Delta x^2}{2} - q_i^n\right),$$

*which simplifies to*

(2.46)
$$q_t + a q_x = -\frac{a\Delta x}{2} q_{xx} - \frac{\Delta t}{2} q_{tt}.$$

*Applying the procedure of Warming-Hyett, similar to the previous example, we obtain*

(2.47)
$$q_{tt} = a^2 q_{xx} + \mathcal{O}\left(\Delta x, \Delta t\right),$$

*which substituted into (2.46) provides the following $\Pi$–form of the differential approximation:*

$$(2.48) \qquad q_t + aq_{xt} = c_2 q_{xx}, \quad \text{with } c_2 = -\frac{a\Delta x}{2}\left(1 + c\right).$$

*Since for positive $a$ also the Courant number $c$ is always positive, we immediately obtain that*

$$(2.49) \qquad c_2 < 0, \quad \forall c > 0,$$

*hence the downwind scheme is unconditionally unstable.*

EXAMPLE 14. *The aim in the following is to analyze the implicit backward finite difference scheme assuming $a > 0$,*

$$(2.50) \qquad q_i^{n+1} = q_i^n - \frac{a\Delta t}{\Delta x}\left(q_i^{n+1} - q_{i-1}^{n+1}\right).$$

*The Taylor series expansion in space and time of $q_{i-1}^{n+1}$ yields*

$$(2.51) \qquad q_{i-1}^{n+1} = q_i^n - \Delta x q_x + \Delta t q_t - \Delta x \Delta t q_{xt} + \frac{\Delta x^2}{2}q_{xx} + \frac{\Delta t^2}{2}q_{tt}.$$

*Substituting the Taylor expansions of $q_{i-1}^{n+1}$ and $q_i^{n+1}$ into (2.50) we obtain*

$$(2.52) \qquad q_t + aq_x = -\frac{\Delta t}{2}q_{tt} - a\Delta t q_{xt} + \frac{a\Delta x}{2}q_{xx}.$$

*The Warming-Hyett procedure applied to the relation above provides*

$$(2.53) \qquad q_{tx} = -aq_{xx} + \mathcal{O}\left(\Delta x, \Delta t\right), \qquad q_{tt} = -aq_{xt} + \mathcal{O}\left(\Delta x, \Delta t\right).$$

*Finally, the modified equation of the PDE up to first order error terms is given by*

$$(2.54) \qquad q_t + aq_x = c_2 q_{xx}, \quad \text{with } c_2 = \frac{a\Delta x}{2}\left(1 + c\right).$$

*The positivity condition of the coefficient $c_2$ is always satisfied*

$$(2.55) \qquad c_2 > 0, \quad \forall c > 0,$$

*hence the scheme is unconditionally stable.*

EXAMPLE 15. *The explicit central finite difference scheme (1.97) is analyzed. Using appropriate Taylor series expansions the $\Gamma$–form of the differential approximation results as*

$$(2.56) \qquad q_t + aq_x = -\frac{a\Delta x^2}{6}q_{xxx} - \frac{\Delta t}{2}q_{tt} - \frac{\Delta t^2}{6}q_{ttt}.$$

*From the Warming-Hyett procedure we obtain the higher order time derivatives as follows:*

$$(2.57) \qquad \begin{aligned} q_{ttt} &= -a^3 q_{xxx} + \mathcal{O}\left(\Delta x, \Delta t\right), \\ q_{tt} &= a^2 q_{xx} + \mathcal{O}\left(\Delta x, \Delta t\right). \end{aligned}$$

*The modified equation of the PDE results as*

$$(2.58) \quad q_t + aq_x = c_2 q_{xx} + c_3 q_{xxx}, \quad \text{with } c_2 = -c\frac{a\Delta x}{2}, \quad c_3 = \frac{a\Delta x^2}{6}\left(c^2 - 1\right).$$

*For the coefficient $c_2$ we obtain*

$$(2.59) \qquad c_2 < 0, \quad \forall c > 0,$$

*hence the explicit central finite difference scheme is unconditionally unstable.*

EXERCISE 3. *Prove that the equivalent differential equation for the Lax–Wendroff method applied to the linear scalar advection equation is up to third order error terms given by*

$$(2.60) \qquad q_t + aq_x = \frac{1}{6}a\Delta x^2(c^2 - 1)q_{xxx} + \frac{1}{8}a\Delta x^3 c(c^2 - 1)q_{xxxx}.$$

## 3. The von Neumann stability analysis

The stability of numerical schemes can also be investigated by carrying out the so–called *von Neumann* stability analysis. The basic idea of this procedure is to assume that the discrete solution can be developed in a Fourier series, instead of the Taylor expansions used in the method of differential approximation. The numerical solution is defined by a component of the Fourier series, as follows

$$(2.61) \qquad q_j^n = \hat{Q}^n e^{i\theta j}, \quad \text{with } i^2 = -1, \quad \theta = k\Delta x = 2\pi\frac{\Delta x}{\lambda_k},$$

where the integer $j$ denotes the spatial index (to avoid confusion with the imaginary unit $i$), $\hat{Q}^n$ is the amplitude at time $t^n$, $0 \leq \theta \leq \pi$ is an angle, $k$ is the wave number and $\lambda_k \geq 2\Delta x$ is the associated wave length. According to the Shannon theorem the smallest wavelength of a periodic function that can still be represented on a discrete mesh of size $\Delta x$ is $\lambda_{\min} = 2\Delta x$. Substituting the definition (2.61) in the numerical scheme we obtain the amplification factor $G$ given by

$$(2.62) \qquad G = \frac{\hat{Q}^{n+1}}{\hat{Q}^n}, \quad \text{with } G \in \mathbb{C}.$$

It defines the ratio of the wave amplitudes at the new and the old time, respectively. In order to obtain a *stable* scheme, the wave amplitude must be *non–increasing*, i.e. the following condition must hold for stability:

$$(2.63) \qquad |G| \leq 1,$$

which is equivalent to

$$(2.64) \qquad |G|^2 = \Re(G)^2 + \Im(G)^2 \leq 1.$$

EXAMPLE 16. *We propose an application of the von Neumann stability analysis to the explicit backward finite difference method (1.95), that is upwind in the case of $a > 0$. First, we replace the function $q_j^n$ with the Fourier mode (2.61)*

$$(2.65) \qquad \hat{Q}^{n+1}e^{i\theta j} = \hat{Q}^n e^{i\theta j} - c\hat{Q}^n e^{i\theta j} + c\hat{Q}^n e^{i\theta(j-1)}.$$

*Subsequently, we compute the amplification factor $G$*

$$(2.66) \qquad G = \frac{\hat{Q}^{n+1}}{\hat{Q}^n} = 1 - c + ce^{-i\theta} = 1 - c + c\left(\cos\theta - i\sin\theta\right),$$

*and the square of its absolute value is given by the sum of the square of the real and the imaginary part of $G$, i.e. for stability we must have*

$$(2.67) \qquad |G|^2 = \left(1 - c + c\cos\theta\right)^2 + c^2\sin^2\theta = 1 - 2c\left(1 - c\right)\left(1 - \cos\theta\right) \leq 1.$$

*In order to verify the stability criterion (2.63) it is equivalent to write*

$$(2.68) \qquad \overbrace{2c}^{>0}\left(1 - c\right)\overbrace{\left(1 - \cos\theta\right)}^{\geq 0} \geq 0,$$

*which becomes a condition on the CFL number*

$$(2.69) \qquad c \leq 1.$$

*This result is again the famous CFL condition for the Courant number* [**18**] *and is identical with the result obtained by the method of differential approximation shown previously. This condition guarantees that small perturbations do not grow in time.*

EXAMPLE 17. *In this example we carry out the von Neumann stability analysis of the explicit central finite difference scheme (1.97) proving that it is unconditionally unstable. Substitution of the Fourier ansatz (2.61) into the numerical scheme (1.97) yields*

$$(2.70) \qquad \hat{Q}^{n+1} e^{i\theta j} = \hat{Q}^n e^{i\theta j} - \hat{Q}^n \frac{c}{2} \left( e^{i\theta(j+1)} - e^{i\theta(j-1)} \right),$$

*hence the amplification factor G becomes*

$$(2.71) \qquad G = \frac{\hat{Q}^{n+1}}{\hat{Q}^n} = 1 - \frac{c}{2} \left( e^{i\theta} - e^{-i\theta} \right) = (1 - c\, i \sin\theta).$$

*We obtain the square of the absolute value of G as*

$$(2.72) \qquad |G|^2 = \Re(G)^2 + \Im(G)^2 = 1 + c^2 \sin^2\theta \geq 1, \quad \forall c > 0.$$

*From (2.72) we can conclude that the stability condition (2.63) is never satisfied, hence the explicit central finite difference method is unconditionally unstable.*

EXAMPLE 18. *We consider the explicit forward finite difference scheme (1.96) and assume $a > 0$. Then, the numerical method is downwind. Using the von Neumann method we obtain*

$$(2.73) \qquad \hat{Q}^{n+1} e^{i\theta j} = \hat{Q}^n e^{i\theta j} - c\hat{Q}^n \left( e^{i\theta(j+1)} - e^{i\theta j} \right).$$

*The amplification factor G is then given by*

$$(2.74) \qquad G = \frac{\hat{Q}^{n+1}}{\hat{Q}^n} = 1 - c(e^{i\theta} - 1) = 1 + c - c(\cos\theta + i\sin\theta).$$

*The stability criterion for the square of the absolute value of the amplification factor is then*

$$(2.75) \qquad |G|^2 = ((1+c) - c\cos\theta)^2 + c^2 \sin^2\theta = 1 + 2c\,(1+c)\,(1 - \cos\theta) \leq 1.$$

*It is obvious that the stability criterion (2.63) is not satisfied*

$$(2.76) \qquad \overbrace{2c}^{>0} \overbrace{(1+c)}^{>0} \overbrace{(1-\cos\theta)}^{\geq 0} \leq 0.$$

*The explicit downwind finite difference scheme is therefore unconditionally unstable.*

EXAMPLE 19. *We analyze the stability of the Lax-Wendroff scheme (1.104). Substituting the Fourier ansatz (2.61) into (1.104) we obtain*

$$(2.77)$$
$$\hat{Q}^{n+1} e^{i\theta j} = \hat{Q}^n e^{i\theta j} - \frac{c}{2} \hat{Q}^n \left( e^{i\theta(j+1)} - e^{i\theta(j-1)} \right) + \frac{c^2}{2} \hat{Q}^n \left( e^{i\theta(j+1)} - 2e^{i\theta j} + e^{i\theta(j-1)} \right).$$

*To verify the stability of the method we compute again the amplification factor, which reduces to*

$$(2.78) \qquad G = \left(1 - c^2\right) + c^2 \cos\theta - i\, c \sin\theta,$$

*and the square of its absolute value must satisfy*

$$(2.79) \qquad |G|^2 = 1 + c^2 \left(1 - \cos^2\theta\right) - 2c^2 \left(1 - \cos\theta\right) + c^4 \left(1 - \cos\theta\right)^2 \leq 1,$$

*that is equivalent to*

$$(2.80) \qquad 1 + \overbrace{c^2 \left(1 - \cos\theta\right)^2}^{\geq 0} \left(c^2 - 1\right) \leq 1.$$

*The stability condition for the Lax-Wendroff scheme* (1.104) *is verified when*

(2.81)
$$c \leq 1,$$

*i.e. the scheme is stable under the classical CFL condition.*

EXERCISE 4. *Prove that the Beam–Warming method for the linear scalar advection equation with $a > 0$ is stable up to Courant number two. Hint: the square of the absolute value of the amplification factor should result as $|G|^2 = 1 + \alpha(1 - \cos\theta)^2$ with $\alpha = c^4 - 4c^3 + 5c^2 - 2c = c(c-2)(c-1)^2$.*

## 4. Monotonicity and the Godunov theorem

We consider a numerical method written in the compact form

(2.82)
$$q_i^{n+1} = H\left(q_{i-l}^n, q_{i-l+1}^n, ..., q_{i-1}^n, q_i^n, q_{i+1}^n, ..., q_{i+r-1}^n, q_{i+r}^n\right),$$

where $H$ can also be a *nonlinear* operator. The scheme is called monotone if $H$ is a non-decreasing function of all its arguments, i.e. when

(2.83)
$$\frac{\partial H}{\partial q_k^n} \geq 0, \quad i - l \leq k \leq i + r.$$

From the definition above it follows that a linear method given by (2.6) is monotone if and only if all its coefficients $b_k$ in (2.6) are non-negative. We therefore must have

(2.84)
$$b_k \geq 0, \quad \forall -l \leq k \leq +r.$$

EXAMPLE 20. *We consider the backward finite difference scheme* (1.96), *which is upwind for $a > 0$:*

(2.85)
$$q_i^{n+1} = H\left(q_{i-1}^n, q_i^n\right) = q_i^n(1-c) + q_{i-1}^n c.$$

*According to definition* (2.83) *we differentiate the function $H$ with respect to its two arguments $q_{i-1}^n$ and $q_i^n$ and obtain the following criteria for monotonicity:*

(2.86)
$$\frac{\partial H}{\partial q_{i-1}^n} = c \geq 0, \quad \frac{\partial H}{\partial q_i^n} = 1 - c \geq 0.$$

*While $c \geq 0$ is always verified, the condition $1 - c \geq 0$ requires $c \leq 1$ for monotonicity, which is the same condition required for the stability of the scheme, see* (2.69). *In other words: if the CFL condition is satisfied, the function $H$ is always non– decreasing and thus the explicit upwind backward FD method is monotone.*

LEMMA 1. *(Roe) A linear method written in form* (2.6) *is of order $p > 0$, in space and time in the sense of the local discretization error*

(2.87)
$$\tau = \frac{L_h}{\Delta t} = \mathcal{O}\left(\Delta x^p, \Delta t^p\right),$$

*if and only if*

(2.88)
$$\sum_{k=-l}^{r} b_k k^m = (-c)^m, \quad \forall m = 0, 1, ..., p \quad \text{with } c = \frac{a \Delta t}{\Delta x}, \ 0^0 := 1.$$

PROOF. The operator of a general linear scheme is given by

(2.89)
$$L_h\left(q_h\left(x_h, t_h\right)\right) = q_i^{n+1} - \sum_{k=-l}^{r} b_k q_{i+k}^n = 0.$$

We compute the local truncation error of the method and replace the expressions $q_i^{n+1}$ and $q_{i+k}^n$ in (2.89) with appropriate Taylor series expansions in time and space, computed with respect to the point $(x_i, t^n)$:

(2.90) $\tau = \dfrac{L_h}{\Delta t} = \dfrac{1}{\Delta t}\left[q_i^n + \displaystyle\sum_{m=1}^{\infty} \frac{\Delta t^m}{m!}\frac{\partial^m q}{\partial t^m} - \sum_{k=-l}^{r} b_k\left(q_i^n + \sum_{m=1}^{\infty} \frac{(k\Delta x)^m}{m!}\frac{\partial^m q}{\partial x^m}\right)\right],$

which is equivalent to

$$(2.91) \quad \tau = \frac{1}{\Delta t} \left[ q_i^n - \sum_{k=-l}^{r} b_k q_i^n + \sum_{m=1}^{\infty} \frac{\Delta t^m}{m!} \frac{\partial^m q}{\partial t^m} - \sum_{k=-l}^{r} b_k \sum_{m=1}^{\infty} \frac{(k\Delta x)^m}{m!} \frac{\partial^m q}{\partial x^m} \right].$$

Since by hypothesis the local truncation error must be of order $p$ in space and in time, i.e.

$$(2.92) \qquad\qquad \tau = \mathcal{O}\left( \Delta x^p, \Delta t^p \right),$$

all terms of order $p$ or lower must vanish in (2.91). For the leading error term we obtain the condition

$$(2.93) \qquad\qquad q_i^n \left( 1 - \sum_{k=-l}^{r} b_k \right) = 0.$$

This is the case $m = 0$ in Eqn. (2.88), which thus has been proven. For the higher order terms we have the conditions

$$(2.94) \qquad \sum_{m=1}^{p} \left( \frac{\Delta t^m}{m!} \frac{\partial^m q}{\partial t^m} - \sum_{k=-l}^{r} b_k \frac{(k\Delta x)^m}{m!} \frac{\partial^m q}{\partial x^m} \right) = 0, \quad \text{for } 0 < m \le p.$$

To simplify this expression we apply the Cauchy-Kovalewski procedure to the linear advection equation obtaining the following general result (the proof is very simple using complete induction):

$$(2.95) \qquad\qquad \frac{\partial^m q}{\partial t^m} = (-a)^m \frac{\partial^m q}{\partial x^m}.$$

Substituting the time derivatives (2.95) into (2.94) we obtain the conditions

$$(2.96) \qquad \sum_{m=1}^{p} \left( \frac{\Delta t^m}{m!} (-a)^m \frac{\partial^m q}{\partial x^m} - \sum_{k=-l}^{r} b_k \frac{(k\Delta x)^m}{m!} \frac{\partial^m q}{\partial x^m} \right) = 0, \quad \text{for } 0 < m \le p,$$

which after factorization yields

$$(2.97) \qquad \sum_{m=1}^{p} \frac{\Delta x^m}{m!} \frac{\partial^m q}{\partial x^m} \left[ \frac{\Delta t^m (-a)^m}{\Delta x^m} - \sum_{k=-l}^{r} b_k k^m \right] = 0$$

The first term in the square brackets can be simplified using the definition of the Courant number $c$. Since this relation must hold for any value of the spatial derivatives, the term in the square brackets must vanish, i.e. we have the condition

$$(2.98) \qquad (-c)^m - \sum_{k=-l}^{r} b_k k^m = 0, \quad \text{for } 0 < m \le p,$$

which completes the proof of the lemma.

$\square$

THEOREM 1. *(Godunov, 1959). For the linear scalar advection equation* $q_t + aq_x = 0$ *there exists* **no** *monotone linear method of the form* (2.6) *of order of accuracy greater or equal than two.*

PROOF. To prove the theorem, we restrict ourselves to the second order case. Let us define coefficients $s_m$ as follows:

$$(2.99) \qquad\qquad s_m = \sum_{k=-l}^{r} k^m b_k.$$

Using a linear method of second order we need to determine the coefficients $s_0$, $s_1$, $s_2$. According to the Lemma of Roe (2.88), they result as $s_0 = 1$, $s_1 = -c$, $s_2 = c^2$. We therefore have

$$(2.100) \qquad s_2 = \sum_{k=-l}^{r} k^2 b_k = \sum_{k=-l}^{r} (k+c)^2 b_k - 2c \overbrace{\sum_{k=-l}^{r} k b_k}^{s_1 = -c} - c^2 \overbrace{\sum_{k=-l}^{r} b_k}^{s_0 = 1} = c^2.$$

Using the expressions for $s_0$ and $s_1$ according to the lemma of Roe, this is equivalent to

$$(2.101) \qquad s_2 = \sum_{k=-l}^{r} (k+c)^2 b_k + 2c^2 - c^2 = c^2,$$

and hence

$$(2.102) \qquad \sum_{k=-l}^{r} \overbrace{(k+c)^2}^{>0} b_k = 0.$$

Following the definition of a monotone linear method (2.84) the coefficients $b_k$ have to be non–decreasing, i.e. $b_k \geq 0$. Hence, the only possibility to satisfy the condition (2.102) preserving the monotonicity implies to impose

$$(2.103) \qquad b_k = 0, \quad \forall k = -l, ..., r,$$

but this leads to a contradiction because according to the Roe lemma we must have for zeroth order consistency

$$(2.104) \qquad s_0 = \sum_{k=-l}^{r} b_k = 1,$$

hence the $b_k$ can not be all zero!                                          $\square$

# Nonlinear Hyperbolic Equations

## 1. Nonlinear scalar hyperbolic equations

We consider the following first order nonlinear scalar equation

$$(3.1) \qquad \frac{\partial q}{\partial t} + \frac{\partial}{\partial x} f(q) = 0,$$

where $q$ denotes the conservative variable and $f(q)$ is the nonlinear flux. The equation (3.1) is called *conservation law*. Classical examples from physics are the conservation of mass, momentum and energy. If the conservative variable $q$ and the flux $f$ are sufficiently regular, the derivatives in (3.1) are defined and the PDE can be written also in the so–called *quasi-linear* form using the chain rule:

$$(3.2) \qquad \frac{\partial q}{\partial t} + a(q) \frac{\partial q}{\partial x} = 0, \quad \text{with } a(q) = \frac{df}{dq} = f'(q).$$

In (3.2) the term $a(q)$ is the characteristic velocity, generally function of $q$. An important property of $a(q)$ is its monotonicity, for which there are three cases:

- Convex flux, when $a(q)$ is monotone and increasing

$$(3.3) \qquad \frac{d}{dq} a(q) = a'(q) = f''(q) > 0, \quad \forall q.$$

- Concave flux, when $a(q)$ is monotone and decreasing

$$(3.4) \qquad \frac{d}{dq} a(q) = a'(q) = f''(q) < 0, \quad \forall q.$$

- Non-convex and non-concave flux if

$$(3.5) \qquad \exists q : \frac{d}{dq} a(q) = a'(q) = f''(q) = 0.$$

EXAMPLE 21. *We consider the nonlinear Burgers equation, defined by*

$$(3.6) \qquad \frac{\partial q}{\partial t} + \frac{\partial}{\partial x} \left( \frac{1}{2} q^2 \right) = 0, \quad f(q) = \frac{1}{2} q^2.$$

*We compute the characteristic velocity $a(q)$ and its derivative*

$$(3.7) \qquad a(q) = \frac{df}{dq} = q, \quad a' = f'' = 1 > 0, \forall q.$$

*We find that the flux of the Burgers equation is convex. This means that larger values of q propagate with higher velocity than the small ones.*

EXAMPLE 22. *The following PDE represents a simple model for traffic flow*

$$(3.8) \qquad \frac{\partial q}{\partial t} + \frac{\partial}{\partial x} f(q) = 0, \quad f(q) = u_{\max} \left( 1 - \frac{q}{q_{\max}} \right) q,$$

*where q denotes the vehicle density per length unit, $u_{\max}$ is the maximum velocity (speed limit) and $q_{\max}$ is the maximum vehicle density. The term $a(q)$ and its*

*derivative $a'$ are given by*

$$(3.9) \qquad a(q) = f'(q) = u_{\max}\left(1 - \frac{2q}{q_{\max}}\right), \quad a' = f'' = -\frac{2u_{\max}}{q_{\max}} < 0, \quad \forall q.$$

*The flux is concave. Here, the larger values of $q$ propagate at slower velocities than the smaller values of $q$.*

EXAMPLE 23. *In this case we consider a PDE with cubic flux*

$$(3.10) \qquad \frac{\partial q}{\partial t} + \frac{\partial}{\partial x} f(q) = 0, \quad f(q) = q^3.$$

*We obtain*

$$(3.11) \qquad a(q) = f'(q) = 3q^2, \quad a' = f'' = 6q \begin{cases} < 0, & \text{if } q < 0, \\ = 0, & \text{if } q = 0, \\ > 0, & \text{if } q > 0. \end{cases}$$

*Hence, the cubic flux is neither convex nor concave.*

**1.1. Solution along the characteristics.** We consider the initial value problem

$$(3.12) \qquad \frac{\partial q}{\partial t} + \frac{\partial f}{\partial x} = 0, \qquad q(x,0) = h(x),$$

which is written in quasi-linear form as

$$(3.13) \qquad \frac{\partial q}{\partial t} + a(q)\frac{\partial q}{\partial x} = 0, \qquad q(x,0) = h(x).$$

In the nonlinear case, the characteristic curves are defined by the ODE

$$(3.14) \qquad \frac{dx}{dt} = a(q(x(t),t)), \qquad x(0) = x_0.$$

Now we compute the material derivative of the variable $q$ along the characteristic

$$(3.15) \qquad \frac{dq}{dt} = \frac{\partial q}{\partial t}\frac{dt}{dt} + \frac{\partial q}{\partial x}\frac{dx}{dt} = \frac{\partial q}{\partial t} + a(q)\frac{\partial q}{\partial x} = 0.$$

We find that even in the nonlinear case, the solution remains *constant* along the characteristic curves, hence also the characteristic velocity $a(q)$ is constant along the characteristic. It follows that the characteristic curves are also *straight lines* even in the *nonlinear* case. Starting from the initial condition in (3.13) we obtain the equation of the characteristic curves as

$$(3.16) \qquad x = x_0 + a(h(x_0))t,$$

which is a nonlinear scalar algebraic equation for the foot of the characteristic $x_0$ in terms of $x$ and $t$. The solution of the nonlinear PDE (3.13) is then given by the initial condition evaluated at the foot of the characteristic $x_0$ as

$$(3.17) \qquad q(x,t) = h(x_0) = h(x - a(h(x_0))t).$$

To verify this solution we compute the derivatives of $q(x,t)$ in time and in space

$$(3.18) \qquad \frac{\partial q}{\partial t} = h'(x_0)\frac{\partial x_0}{\partial t}, \qquad \frac{\partial q}{\partial x} = h'(x_0)\frac{\partial x_0}{\partial x}.$$

Similarly, we differentiate (3.16)

$$(3.19) \qquad \begin{array}{l} \frac{\partial x}{\partial t} = \frac{\partial x_0}{\partial t} + a + a'h'\frac{\partial x_0}{\partial t}t = 0 \Rightarrow \frac{\partial x_0}{\partial t} = \frac{-a}{1+a'h't}, \\ \frac{\partial x}{\partial x} = \frac{\partial x_0}{\partial x} + a + a'h'\frac{\partial x_0}{\partial x}t = 1 \Rightarrow \frac{\partial x_0}{\partial x} = \frac{1}{1+a'h't}. \end{array}$$

Substituting into (3.15) we verify that (3.17) is the solution of the PDE

$$(3.20) \qquad \frac{\partial q}{\partial t} + a(q)\frac{\partial q}{\partial x} = -\frac{ah'}{1+a'h't} + \frac{ah'}{1+a'h't} = 0.$$
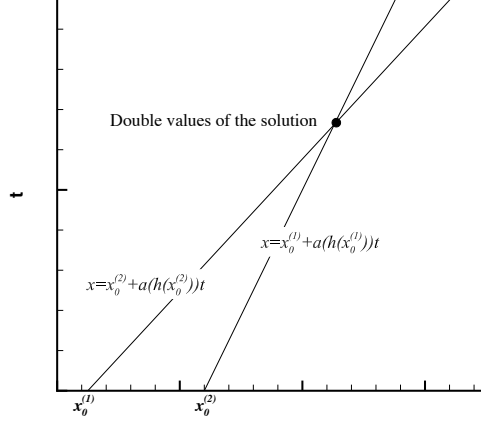
FIGURE 1. The intersection of two characteristic curves leads to a multi–valued solution in the intersection point.

Note that the solution of (3.12) is formally derived as in the linear case using the characteristic curves, but it is *implicit*. Then, it necessary to implement an iterative scheme, such the Newton or bisection method, to compute $x_0$ from $x$ and $t$ solving the following equation

$$(3.21) \qquad g(x_0) = x_0 - x + a(h(x_0))t = 0.$$

Moreover, the method of characteristics breaks down when two characteristic curves intersect, which we will see later in the section of shock waves. Since each characteristic curve transports a constant solution, at the intersection point there would be two different values of $q$ defined, according to the different values transported by each characteristic. In this case we would have a so–called multi–valued solution at the intersection point $(x, t)$, which is sketched in Fig. 1.

**1.2. Rarefaction waves.** As an example, we consider the following initial value problem

$$(3.22) \quad \begin{aligned} &\text{PDE:} \quad q_t + f_x = 0, \qquad f'' > 0, \\ &\text{IC:} \quad q(x,0) = h(x) = \begin{cases} q_L, & \text{if } x < x_L, \\ q_L + \frac{q_R - q_L}{x_R - x_L}(x - x_L), & \text{if } x_L \le x \le x_R, \\ q_R & \text{if } x > x_R, \end{cases} \end{aligned}$$

with $a(q_L) \le a(q_R)$. The initial condition of the problem (3.22) consists of two piecewise constant states and of a piecewise linear distribution of the variable $q$ between the points $x_L$ and $x_R$. The solution of (3.22) is given by a rarefaction wave, plotted in Fig. 2. Its borders are defined by the two particular characteristic curves

$$(3.23) \qquad \begin{aligned} x &= x_L + a(q_L)t, \\ x &= x_R + a(q_R)t, \end{aligned}$$

which are called the *head* and the *tail* of the rarefaction wave with the foot points located at $x_L$ and $x_R$, respectively. The characteristics associated with the tail and the head of the rarefaction divide the half-plane $x - t$ into three regions $R_0$, $R_1$ and $R_2$.
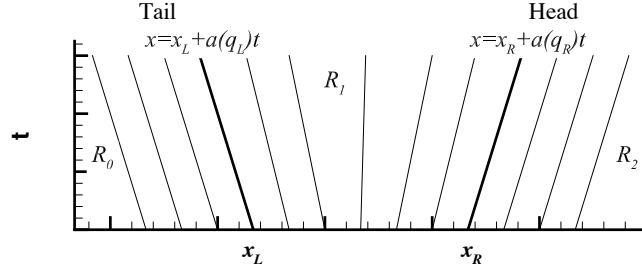
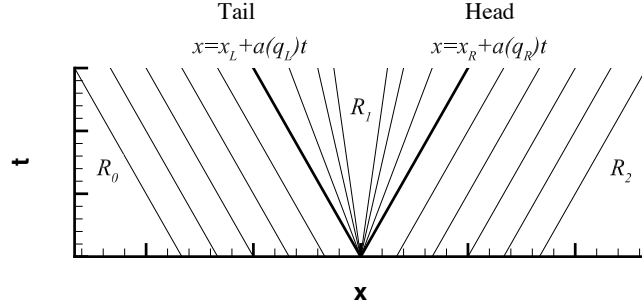FIGURE 2. The solution of the initial value problem (3.22).



FIGURE 3. The solution of the Riemann problem (3.24).

Now, we consider the Riemann problem

$$
(3.24) \qquad
\begin{aligned}
q_t + f_x &= 0, \quad f'' > 0, \\
q(x,0) = h(x) &= \begin{cases} q_L, & \text{if } x < x_0, \\ q_R, & \text{if } x \geq x_0, \end{cases}
\end{aligned}
$$

with $a(q_L) \leq a(q_R)$. This problem can be considered as the limit case of the previous one (3.22) with $(x_L, x_R) \to x_0$. Thus, we find self-similar solutions $q(x,t)$ depending only on one variable $\xi = \frac{x-x_0}{t}$ so that

$$
(3.25) \qquad q(x,t) = q(\xi) = q\left(\frac{x-x_0}{t}\right).
$$

We compute the derivatives in space and in time of the unknown $q(x,t)$

$$
(3.26) \qquad \frac{\partial q}{\partial t} = -\frac{x-x_0}{t^2} q', \quad \frac{\partial q}{\partial x} = \frac{1}{t} q',
$$

where $q'$ means the derivative with respect to $\xi$. Substitution into (3.22) yields

$$
(3.27) \qquad \left(-\frac{x-x_0}{t^2} + a(q(\xi))\frac{1}{t}\right) Q' = 0, \quad \Rightarrow \quad a(q(\xi)) = \frac{x-x_0}{t} = \xi.
$$

For the Riemann problem (3.24), the foot points of the head and the tail of the rarefaction wave coincide with the initial location of the discontinuity $x_0$, as depicted in Fig. 3. The complete solution of the Riemann problem is given by

$$
(3.28) \quad q(x,t) = \begin{cases} q_L, & \text{if } \frac{x-x_0}{t} < a(q_L), \\ \text{root of } g(q) = a(q) - \frac{x-x_0}{t} = 0, & \text{if } a(q_L) \leq \frac{x-x_0}{t} \leq a(q_R), \\ q_R, & \text{if } \frac{x-x_0}{t} > a(q_R), \end{cases}
$$

where the intermediate states between the head and the tail of the rarefaction wave are computed by solving the nonlinear scalar algebraic equation

$$(3.29) \qquad a(q) = \xi,$$

which can be solved numerically by using for example a Newton–Raphson method or the bisection algorithm.

**1.3. Shock waves.** As a motivation for the introduction of shock waves, first introduced into the theory of nonlinear partial differential equations by Bernhard Riemann in [**67, 68**], we solve following initial value problem for the Burgers equation (3.6) written in quasi–linear form

$$(3.30) \qquad \begin{aligned} &\text{PDE:} \quad \frac{\partial q}{\partial t} + q \frac{\partial q}{\partial x} = 0, \\ &\text{IC:} \quad q(x,0) = h(x) = \begin{cases} 1, & \text{if } x < 0, \\ 1 - x, & \text{if } 0 \le x \ge 1, \\ 0, & \text{if } x \ge 1. \end{cases} \end{aligned}$$

The characteristic curves divide the domain in three regions, i.e. $R_0$, $R_1$ and $R_2$, as shown in Fig. 4. The lines $c_0$ and $c_1$ are defined by

$$(3.31) \qquad \begin{aligned} c_0 : \quad & x = 0 + 1 \cdot t = t, \\ c_1 \quad & x = 1 + 0 \cdot t = 1. \end{aligned}$$

In the regions $R_0$ and $R_2$ all the characteristic curves are parallel to the line $c_0$ and $c_1$, respectively. In the middle (region $R_1$) the characteristics converge on the intersection point $P$ and the foot of the characteristic is computed from (3.16) as

$$(3.32) \qquad x = x_0 + a(h(x_0))t = x_0 + (1 - x_0)t = x_0(1 - t) + t.$$

Hence, the foot of the characteristic in terms of $x$ and $t$ is given by

$$(3.33) \qquad x_0 = \frac{x - t}{1 - t},$$

and thus the solution along the characteristic curve in the region $R_1$ is

$$(3.34) \qquad q(x,t) = h(x_0) = 1 - x_0 = 1 - \frac{x - t}{1 - t}.$$

To conclude, we summarize the complete solution until the intersection time $t = 1$:

$$(3.35) \qquad q(x,t) = \begin{cases} 1, & \text{if } x < t, \\ \frac{1-x}{1-t}, & \text{if } t \le x \le 1, \\ 0, & \text{if } x > 1. \end{cases}$$

Fig. 5 shows the evolution in time of the solution of (3.30). Note that at time $t \ge 1$ the solution $q$ has multiple values for $x \ge 1$. The problem of multi–valued solutions has been solved by Bernhard Riemann already in 1860 [**67, 68**] by the introduction of the concept of *shock waves*, i.e. a *discontinuity* in the solution that guarantees the conservation of the variable $q$ through the condition $A_1 = A_2$. The terms $A_1$ and $A_2$ are plotted in Fig. 6 and denote the area given by the intersection between the multi–valued solution (solid line) and the shock profile assumed by Riemann (dashed line). The discontinuity is localized so that the solution is *conservative*.

**1.4. Conservation laws in integral form.** We analyze a general conservative PDE (3.1). Let us define a control volume in space $I = [x_L, x_R]$ and integrate the PDE over this spatial control volume:

$$(3.36) \qquad \int_{x_L}^{x_R} \frac{\partial q}{\partial t} dx + \int_{x_L}^{x_R} \frac{\partial}{\partial x} f(q) \, dx = 0.$$

FIGURE 4. The characteristic curve corresponding to the problem (3.30) that define three regions $R_0$, $R_1$ and $R_2$.



FIGURE 5. The solution of (3.30) at times $t = \frac{1}{2}, 1, \frac{3}{2}, 2$.



FIGURE 6. The solution at times $t = 2$ that presents a shock wave with multiple values of $q$. Assuming the equality of $A_1$ and $A_2$ the shock wave can be localized so that the solution is conservative.

Integrating the second term by parts we obtain

$$(3.37) \qquad \frac{\partial}{\partial t} \int_{x_L}^{x_R} q(x,t)\, dx + f(q(x_R,t)) - f(q(x_L,t)) = 0.$$

This means that the temporal variation of the conservative variable $q(x,t)$ inside the spatial interval $I = [x_L; x_R]$ is only due to the difference of the fluxes on the boundaries of the control volume $I$.

When the control volume is defined in space *and* in time $V = [x_L; x_R] \times [t_1; t_2]$, we obtain

$$(3.38) \qquad \int_{t_1}^{t_2} \int_{x_L}^{x_R} \frac{\partial q}{\partial t} \, dx dt + \int_{t_1}^{t_2} \int_{x_L}^{x_R} \frac{\partial}{\partial x} f(q) \, dx dt = 0,$$

that is equivalent to

$$(3.39) \qquad \int_{x_L}^{x_R} q(x, t_2) \, dx = \int_{x_L}^{x_R} q(x, t_1) \, dx - \int_{t_1}^{t_2} \left( f(q(x_R, t)) - f(q(x_L, t)) \right) dt.$$

Thus, the conservative quantity at time $t_2$ is equal to the conservative quantity at time $t_1$ subtracted by the difference of the integrals of the fluxes in time computed on the spatial boundaries of the control volume.

**1.5. Rankine-Hugoniot conditions.** The Rankine–Hugoniot conditions relate the jump of the conserved quantities across a shock wave with the propagation velocity $s$ of the discontinuity. They are derived from integral conservation over a control volume that includes the shock.

Let us assume that the discontinuity travels with a velocity $s$ and that it is an isolated wave that divides two piecewise constant states on the left $q_L$ and on the right $q_R$ of the shock. Then, the PDE in integral form over a space–time control volume $V = [x_L; x_R] \times [t_1; t_2]$ with $x_L \le x_d \le x_R$ and $x_L \le x_d + s(t_2 - t_1)$ is given by

$$(3.40) \qquad \int_{x_L}^{x_R} \left( q(x, t_2) - q(x, t_1) \right) dx + \left( \int_{t_1}^{t_2} \left( f(q(x_R, t)) - f(q(x_L, t)) \right) dt \right) = 0,$$

where $x_d$ is the location of the shock wave at time $t_1$. Evaluation of the integrals yields

$$(3.41) \qquad \begin{aligned} & q_L \left( x_d + s(t_2 - t_1) \right) + q_R \left( x_R - x_d - s(t_2 - t_1) \right) - \\ & -q_L(x_d - x_L) - q_R(x_R - x_d) + \left( f(q_R) - f(q_L) \right)(t_2 - t_1) = 0. \end{aligned}$$

This leads to the Rankine-Hugoniot relations

$$(3.42) \qquad s(q_R - q_L) = f(q_R) - f(q_L),$$

which relate the jump of the conservative quantity $q$ with the speed of the shock wave $s$ and the jump of the flux over the discontinuity.

EXAMPLE 24. *We consider the Burgers equation (3.6) in conservative form given by the flux $f(q) = \frac{1}{2}q^2$. The relations of Rankine-Hugoniot provide the shock speed $s$, as follows*

$$(3.43) \qquad \left( \frac{1}{2}q_R^2 - \frac{1}{2}q_L^2 \right) = s(q_R - q_L).$$

*Simplifying we obtain*

$$(3.44) \qquad s = \frac{1}{2}(q_R + q_L).$$

**1.6. Riemann problem with shock wave.** We consider the initial value problem

$$(3.45) \qquad \begin{aligned} &\text{PDE:} \quad \frac{\partial q}{\partial t} + \frac{\partial}{\partial x} f(q) = 0, \\ &\text{IC:} \quad q(x, 0) = h(x) = \begin{cases} q_L, & \text{if } x < x_0, \\ q_R, & \text{if } x \ge x_0, \end{cases} \end{aligned}$$

where we suppose $a(q_L) > a(q_R)$. The solution of this initial value problem is given by

$$(3.46) \qquad q(x, t) = \begin{cases} q_L, & \text{if } \frac{x - x_0}{t} < s, \\ q_R, & \text{if } \frac{x - x_0}{t} \ge s, \end{cases}$$
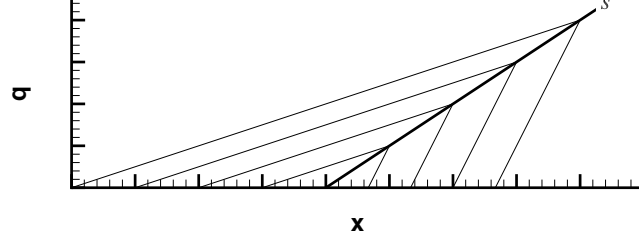
FIGURE 7. The solution of the initial value problem that presents a shock that propagates in space with a celerity $s$.

where the shock speed $s$ is determined by the use of the Rankine–Hugoniot conditions (3.42). The characteristic curves intersect and lead to a discontinuity (shock wave), traveling with speed $s$, as shown in Fig. 7.

**1.7. Non-uniqueness of solutions and the Lax entropy condition.**
Solving nonlinear hyperbolic conservation laws using the weak (integral) formulation allows us to admit discontinuities in the solution, thus extending the set of possible solutions even if they are not sufficiently smooth. However, this does not guarantee uniqueness. It is necessary to introduce further conditions that lead to a unique solution starting from an initial value problem.

EXAMPLE 25. *We solve the following initial value problem for the Burgers equation* (3.6)*:*

$$
(3.47) \qquad
\begin{aligned}
&\text{PDE:} \quad \frac{\partial q}{\partial t} + q \frac{\partial q}{\partial x} = 0, \\
&\text{IC:} \qquad q(x,0) = h(x) = \begin{cases} 0, & \text{if } x < 0, \\ 1, & \text{if } x \geq 0, \end{cases}
\end{aligned}
$$

*We have two possible solution alternatives, namely a rarefaction wave or a shock wave.*

*In the first case, we assume a rarefaction wave solution, plotted in Fig. 8 and given by*

$$
(3.48) \qquad q(x,t) = \begin{cases} 0, & \text{if } \frac{x}{t} < 0, \\ \frac{x}{t}, & \text{if } 0 \leq \frac{x}{t} \leq 1, \\ 1, & \text{if } \frac{x}{t} > 1. \end{cases}
$$

*This first solution is a classical solution, based on the theory of characteristics and therefore it is also a weak solution since all classical solutions that satisfy the PDE in its strong differential form also satisfy the PDE in its weak integral form.*

*In the second case, the solution contains a discontinuity (rarefaction shock), as plotted in Fig. 9. It is given by*

$$
(3.49) \qquad q(x,t) = \begin{cases} 0, & \text{if } \frac{x}{t} < s, \\ 1, & \text{if } \frac{x}{t} \geq s, \end{cases} \quad \text{with } s = \frac{1}{2}(q_L + q_R) = \frac{1}{2}.
$$

*The second solution, obtained by using the Rankine–Hugoniot conditions, satisfies the integral form of the conservation law and therefore is also a weak solution of our conservation law. This can be proven by using the integral form of the conservation law* (3.39) *over the space–time control volume* $[x_L; x_R] \times [t_1; t_2]$*, see Fig. 10. With* $s = 1/2$ *we obtain*

$$
(3.50) \qquad x_R - st_2 - (x_R - st_1) + \frac{1}{2}(t_2 - t_1) = 0.
$$

FIGURE 8. The solution of the Burgers equation including the rarefaction fan.



FIGURE 9. The solution of the Burgers equation including the rarefaction shock.



FIGURE 10. The space-time control volume used for the integral form (3.39) of the conservation law.

*This proves indeed that also the rarefaction shock solution satisfies the integral form of the conservation law and hence is a possible weak solution of the problem. Nevertheless, according to the entropy criterion of Lax it is not an admissible solution from the physical point of view, since shock waves are only produced if the characteristic curves enter the shock wave. In contrast to this, as plotted in Fig. 9, for the rarefaction shock solution, the characteristic curves leave the shock wave.*

Since hyperbolic conservation laws allow more than one solution if we also admit weak solutions, we need an additional criterion to select the physically correct one. A criterion introduced by Lax is the so–called Lax entropy condition

$$(3.51) \qquad\qquad\qquad a\left(q_L\right) > s > a\left(q_R\right),$$

where $s$ is the velocity of the shock obtained from the Rankine–Hugoniot condition. The Lax entropy condition requires that the characteristic curves enter the shock

wave from both sides, or in other words the shock is *compressed* from both sides (left and right). When the criterion (3.51) is not satisfied, then the unique solution that is physically admissible is the rarefaction wave.

## 2. Nonlinear hyperbolic systems

We consider the following system of hyperbolic conservation laws

$$(3.52) \qquad \frac{\partial \vec{Q}}{\partial t} + \frac{\partial}{\partial x} \vec{f}\left(\vec{Q}\right) = 0, \quad \text{with } \vec{Q} \in \Omega_Q \subset \mathbb{R}^m, \vec{f} \in \mathbb{R}^m,$$

in which the PDE are nonlinear and the terms $\vec{Q}$ and $\vec{f}(\vec{Q})$ are the vector of the conserved quantities and the flux vector, respectively.

$$(3.53) \qquad \vec{Q} = (q_1, q_2, ..., q_m)^T, \quad \vec{f}\left(\vec{Q}\right) = (f_1, f_2, ..., f_m)^T.$$

In quasi-linear form the expression above results in

$$(3.54) \qquad \frac{\partial \vec{Q}}{\partial t} + \underline{\underline{A}}\left(\vec{Q}\right) \frac{\partial \vec{Q}}{\partial x} = 0, \quad \text{with} \quad \underline{\underline{A}}\left(\vec{Q}\right) = \frac{\partial \vec{f}}{\partial \vec{Q}},$$

where the matrix $\underline{\underline{A}}(\vec{Q})$ is the Jacobian matrix of the flux vector. The system (3.52) is called *hyperbolic* if the Jacobian matrix of the flux has $m$ real eigenvalues and a complete set of linearly independent eigenvectors for all states $\vec{Q} \in \Omega_Q$.

Similarly to (3.39) we define a space-time control volume $V = [x_L; x_R] \times [t_1, t_2]$ and write the nonlinear system (3.52) in integral form

$$(3.55)$$
$$\int_{x_L}^{x_R} \vec{Q}\left(x, t_2\right) dx = \int_{x_L}^{x_R} \vec{Q}\left(x, t_1\right) dx + \int_{t_1}^{t_2} \left( \vec{f}\left(\vec{Q}\left(x_R, t\right)\right) - \vec{f}\left(\vec{Q}\left(x_L, t\right)\right) \right) dt.$$

The corresponding Rankine-Hugoniot conditions at a shock wave propagating with speed $s$ are given by

$$(3.56) \qquad \vec{f}\left(\vec{Q}\left(x_R, t\right)\right) - \vec{f}\left(\vec{Q}\left(x_L, t\right)\right) = s\left(\vec{Q}\left(x_R, t\right) - \vec{Q}\left(x_L, t\right)\right).$$

To understand the nature of the characteristic curves we analyze the $i^{th}$ eigenvalue $\lambda_i$ of the matrix $\underline{\underline{A}}$ and the associated eigenvector $\vec{r}_i$. The characteristic field is defined to be *linearly degenerate* if

$$(3.57) \qquad \nabla \lambda_i\left(\vec{Q}\right) \cdot \vec{r}_i\left(\vec{Q}\right) = 0, \quad \forall \vec{Q} \in \Omega_Q.$$

Otherwise, it is *genuinely nonlinear* when

$$(3.58) \qquad \nabla \lambda_i\left(\vec{Q}\right) \cdot \vec{r}_i\left(\vec{Q}\right) \neq 0, \quad \forall \vec{Q} \in \Omega_Q.$$

Here, the gradient $\nabla \lambda$ means the gradient with respect to the vector of conserved quantities, i.e. we use the notation

$$(3.59) \qquad \nabla \lambda = \frac{\partial \lambda}{\partial \vec{Q}}.$$

The phase-space or state-space $\Omega_Q$ is the space of the vectors

$$(3.60) \qquad \vec{Q} = (q_1, q_2, ..., q_m)^T \in \Omega_{\vec{Q}} \subseteq \mathbb{R}^m.$$

EXAMPLE 26. *We consider the equations of isentropic gasdynamics given by*

$$(3.61) \qquad \frac{\partial}{\partial t}\left(\begin{array}{c} \rho \\ \rho u \end{array}\right) + \frac{\partial}{\partial x}\left(\begin{array}{c} \rho u \\ \rho u^2 + p \end{array}\right) = 0,$$

*where the variables $\rho$, $u$ and $p$ denote the density, the velocity in $x$–direction and the pressure, respectively. The system (3.61) is not closed because the number of*

*equations is not sufficient to determine the flow field. Hence, a closure relation is needed, called the equation of state (EOS). A very simple choice consists of the isentropic EOS, given by*

$$(3.62) \qquad p = k\rho^\gamma, \quad \text{with } \gamma = \frac{c_P}{c_V}, \quad \gamma \in \mathbb{R}, k \in \mathbb{R},$$

*where $c_V$ and $c_P$ are the heat capacity at constant volume and the heat capacity at constant pressure, respectively. Replacing (3.62) into (3.61) and writing the system in terms of the conservative variables $\vec{Q} = (q_1, q_2)^T = (\rho, \rho u)^T$ we obtain*

$$(3.63) \qquad \frac{\partial}{\partial t} \begin{pmatrix} q_1 \\ q_2 \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} q_2 \\ \frac{q_2^2}{q_1} + kq_1^\gamma \end{pmatrix} = 0.$$

*It follows that the Jacobian matrix of the flux vector of isentropic gasdynamics results as*

$$(3.64) \qquad \underline{\underline{A}}\left(\vec{Q}\right) = \frac{\partial \vec{f}}{\partial \vec{Q}} = \begin{pmatrix} 0 & 1 \\ C^2 - \frac{q_2^2}{q_1^2} & 2\frac{q_2}{q_1} \end{pmatrix},$$

*where c is the sound speed, given by*

$$(3.65) \qquad C^2 = \frac{dp}{d\rho} = \gamma k q_1^{\gamma-1} = \frac{\gamma p}{\rho}.$$

*The matrix $\underline{\underline{A}}$ can be written in physical variables, i.e. primitive variables, as follows*

$$(3.66) \qquad \underline{\underline{A}} = \begin{pmatrix} 0 & 1 \\ C^2 - u^2 & 2u \end{pmatrix}.$$

*Now, we compute the eigenvalues from the characteristic polynomial*

$$(3.67) \qquad \left\| \underline{\underline{A}} - \lambda \underline{\underline{I}} \right\| = \lambda^2 - 2u\lambda - C^2 + u^2 = 0,$$

*the roots of which are $\lambda_{1,2} = u \pm C$. Using the first eigenvalue we obtain the following right eigenvector*

$$(3.68) \quad \begin{pmatrix} C - u & 1 \\ C^2 - u^2 & 2u + C - u \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \Rightarrow \vec{r}_1 = \begin{pmatrix} 1 \\ u - C \end{pmatrix}.$$

*Similarly, the right eigenvector associated with the second eigenvalue is*

$$(3.69) \quad \begin{pmatrix} -C - u & 1 \\ C^2 - u^2 & u - C \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \Rightarrow \vec{r}_1 = \begin{pmatrix} 1 \\ u + C \end{pmatrix}.$$

*In matrix form we have*
$$(3.70)$$
$$\underline{\underline{R}} = \begin{pmatrix} 1 & 1 \\ u - C & u + C \end{pmatrix}, \quad \underline{\underline{\Lambda}} = \begin{pmatrix} u - C & 0 \\ 0 & u + C \end{pmatrix}, \quad \underline{\underline{R}}^{-1} = \frac{1}{2C} \begin{pmatrix} C + u & -1 \\ C - u & 1 \end{pmatrix}.$$

EXAMPLE 27. *We consider the shallow water model, that consists of a set of hyperbolic PDE*

$$(3.71) \qquad \frac{\partial}{\partial t} \begin{pmatrix} h \\ hu \end{pmatrix} + \frac{\partial}{\partial x} \begin{pmatrix} hu \\ hu^2 + \bar{p} \end{pmatrix} = 0,$$

*where the variables $h$ and $u$ are the water depth and the velocity in $x$–direction, respectively. The term $\bar{p}$ denotes the averaged pressure force along the water depth. Since the shallow water equations are derived from depth–integrating the incompressible Navier-Stokes equations assuming that the horizontal length scale is much larger than the vertical one, the vertical pressure is assumed to be hydrostatic. Then, the average pressure force (normalized with the fluid density $\rho$) is given by*

$$(3.72) \qquad \bar{p} = \frac{1}{2} gh^2,$$

*where $g$ is the gravity acceleration along the vertical direction. We write the system (3.71) in terms of the conservative variables $\vec{Q} = (q_1, q_2)^T = (h, hu)^T$, as follows*

$$(3.73) \qquad \frac{\partial}{\partial t}\begin{pmatrix} q_1 \\ q_2 \end{pmatrix} + \frac{\partial}{\partial x}\begin{pmatrix} q_2 \\ \frac{q_2^2}{q_1} + \frac{1}{2}gq_1^2 \end{pmatrix} = 0.$$

*The Jacobian matrix $\underline{\underline{A}}$ is given by*

$$(3.74) \qquad \underline{\underline{A}}\left(\vec{Q}\right) = \frac{\partial \vec{f}}{\partial \vec{Q}} = \begin{pmatrix} 0 & 1 \\ -\frac{q_2^2}{q_1^2} + C^2 & 2\frac{q_2}{q_1} \end{pmatrix},$$

*where $C$ is the celerity of the gravity waves that propagate on the water surface and is equal to $C = \sqrt{gh}$. Similarly, in primitive variables the Jacobian matrix is written as*

$$(3.75) \qquad \underline{\underline{A}} = \begin{pmatrix} 0 & 1 \\ C^2 - u^2 & 2u \end{pmatrix}.$$

*Note that there is a very strong similarity between the system of isentropic compressible gasdynamics (3.61) and the shallow water model. Hence, we can write by analogy as follows*

(3.76)

$$\underline{\underline{R}} = \begin{pmatrix} 1 & 1 \\ u - C & u + C \end{pmatrix}, \quad \underline{\underline{\Lambda}} = \begin{pmatrix} u - C & 0 \\ 0 & u + C \end{pmatrix}, \quad \underline{\underline{R}}^{-1} = \frac{1}{2C}\begin{pmatrix} C + u & -1 \\ C - u & 1 \end{pmatrix}.$$

## 3. The exact solution of the Riemann problem for isothermal gasdynamics

For the Godunov flux [37] described later in this manuscript we need the *exact solution* of the local Riemann problem at element interfaces. We already know how to solve the Riemann problem in the case of general linear hyperbolic systems and for general scalar nonlinear hyperbolic equations. In this section and the following one, we will give a strategy for the solution of the Riemann problem for two particular nonlinear hyperbolic systems, namely the governing equations of isothermal compressible gasdynamics augmented by the transport of a passive scalar and the shallow water equations. We restrict ourselves to these two particular systems since the general nonlinear case is very difficult to handle. For many important but rather complicated hyperbolic systems the exact solution of the Riemann problem is still the topic of current research.

In this section, we consider the isothermal compressible gasdynamics, augmented by the transport equation of a passive scalar, such as small dust particles. The governing PDE system consists of the conservation equations of mass and momentum for the fluid as well as the mass conservation equation for the passive scalar given by

$$(3.77) \qquad \begin{aligned} \frac{\partial \rho}{\partial t} + \frac{\partial}{\partial x}(\rho u) &= 0, \\ \frac{\partial \rho u}{\partial t} + \frac{\partial}{\partial x}(\rho u^2 + p) &= 0, \\ \frac{\partial \rho \psi}{\partial t} + \frac{\partial}{\partial x}(\rho u \psi) &= 0. \end{aligned}$$

where $\rho$ is the mass density of the fluid, $u$ is the velocity, $\psi$ is the concentration of the passive scalar and $p = a^2\rho$ is the pressure according to the isothermal hypothesis. Recall that the ideal gas law is $p/\rho = RT$ and for the isothermal case the square of

the sound speed is hence $a^2 = RT$, where $R$ is the specific gas constant and $T$ is the temperature. The nonlinear PDE system can be written as

$$\frac{\partial Q}{\partial t} + \frac{\partial}{\partial x} f(Q) = 0, \qquad Q, f \in \mathbb{R}^3, \quad t \in \mathbb{R}_0^+, \quad x \in \mathbb{R},$$

with the vector of conserved variables $Q$ and the flux vector $f$ being defined as

$$(3.78) \qquad Q = \begin{pmatrix} \rho \\ \rho u \\ \rho \psi \end{pmatrix} = \begin{pmatrix} q_1 \\ q_2 \\ q_3 \end{pmatrix}, f = \begin{pmatrix} \rho u \\ \rho u^2 + p \\ \rho u \psi \end{pmatrix} = \begin{pmatrix} q_2 \\ q_2^2/q_1 + a^2 q_1 \\ q_2 q_3/q_1 \end{pmatrix}.$$

The Jacobian matrix $A$ of the flux $f$ with respect to the vector of conserved variables $Q$ is

$$(3.79) \qquad A = \begin{pmatrix} 0 & 1 & 0 \\ -q_2^2/q_1^2 + a^2 & 2q_2/q_1 & 0 \\ -q_2 q_3/q_1^2 & q_3/q_1 & q_2/q_1 \end{pmatrix} = \begin{pmatrix} 0 & 1 & 0 \\ a^2 - u^2 & 2u & 0 \\ -u\psi & \psi & u \end{pmatrix}.$$

For (3.78) the diagonal matrix of eigenvalues is $\Lambda = \mathrm{diag}(u - a, u, u + a)$ and the corresponding matrix of right eigenvectors $R$ and its inverse $R^{-1}$ are given by

$$(3.80) \qquad R = \begin{pmatrix} 1 & 0 & 1 \\ u - a & 0 & u + a \\ \psi & 1 & \psi \end{pmatrix}, \quad R^{-1} = \frac{1}{2a} \begin{pmatrix} a + u & -1 & 0 \\ -2a\psi & 0 & 2a \\ a - u & 1 & 0 \end{pmatrix}.$$

For the first and the third eigenvalue we find that the associated characteristic fields are *genuinely nonlinear* since

$$(3.81) \quad \nabla\lambda_{1,3} = \nabla(u \mp a) = \frac{\partial}{\partial Q}(q_2/q_1 \mp a) = (-q_2/q_1^2, 1/q_1, 0)^T = (-u/\rho, 1/\rho, 0)^T.$$

and hence

$$(3.82) \qquad \nabla\lambda_{1,3} \cdot \vec{r}_{1,3} = (-u/\rho, 1/\rho, 0)^T \cdot (1, u \mp a, 0)^T = \mp a/\rho \neq 0, \quad \forall Q \in \Omega_Q.$$

The second characteristic field, associated to $\lambda_2 = u$ is linearly degenerate, which can be seen as follows:

$$(3.83) \qquad \nabla\lambda_2 = \nabla u = \frac{\partial}{\partial Q}(q_2/q_1) = (-q_2/q_1^2, 1/q_1, 0)^T = (-u/\rho, 1/\rho, 0)^T.$$

$$(3.84) \qquad \nabla\lambda_2 \cdot \vec{r}_2 = (-u/\rho, 1/\rho, 0)^T \cdot (0, 0, 1)^T = 0, \quad \forall Q \in \Omega_Q.$$

The general solution of the Riemann problem for the nonlinear PDE system (3.78) consists of three waves associated with the three eigenvalues of the PDE. The intermediate wave, associated to $\lambda_2$ is always a so–called *contact discontinuity*, due to the linear degeneracy of the second characteristic field. The nonlinear waves associated with the first and the third eigenvalue can be either a shock wave, in the case of compression of the characteristics, or a centered rarefaction wave, according to the Lax entropy condition. A sketch for a typical configuration can be found in Fig. 11.

**3.1. Shock waves and Rankine–Hugoniot relations.** The general Rankine–Hugoniot conditions for a system of conservation laws read

$$(3.85) \qquad\qquad s\,(Q_R - Q_L) = (f(Q_R) - f(Q_L)).$$

Since the algebraic manipulations can be performed much easier for a steady shock wave, we use the principle of Galilean invariance of Newtonian mechanics which states that the governing PDE remain unchanged under a linear transformation of the velocity of the form

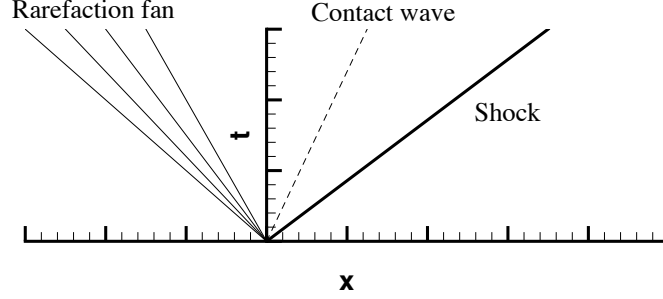$$(3.86) \qquad\qquad \hat{u} = u - s.$$

FIGURE 11. The plot for a typical solution of the Riemann problem for the nonlinear PDE system (3.78).

Here, $u$ is the fluid velocity in the fixed laboratory frame, $s$ is the (constant) velocity of an observer that moves with the speed of the shock wave $s$ and $\hat{u}$ is the transformed velocity in the co–moving frame. Note that in particular we have $\hat{s} = 0$, since the shock wave is steady in the co–moving frame. With the use of the principle of Galilean invariance, the Rankine–Hugoniot relations written in the co–moving frame reduce to

$$(3.87) \qquad 0 \cdot (\hat{Q}_R - \hat{Q}_L) = f(\hat{Q}_R) - f(\hat{Q}_L), \qquad \Rightarrow \qquad f(\hat{Q}_R) = f(\hat{Q}_L).$$

In isothermal gasdynamics this leads to

$$(3.88) \qquad \begin{aligned} \rho_R \hat{u}_R &= \rho_L \hat{u}_L := M, \\ \rho_R \hat{u}_R^2 + a^2 \rho_R &= \rho_L \hat{u}_L^2 + a^2 \rho_L, \\ \rho_R \hat{u}_R \psi_R &= \rho_L \hat{u}_L \psi_L. \end{aligned}$$

Performing algebraic manipulations on the equations above we obtain

$$(3.89) \qquad \begin{aligned} \psi_R &= \psi_L, \\ M &= -a^2 \frac{\rho_R - \rho_L}{\hat{u}_R - \hat{u}_L}, \\ \hat{u}_R &= \frac{M}{\rho_R}, \quad \hat{u}_L = \frac{M}{\rho_L}. \end{aligned}$$

Finally, the transformed velocity referred to the states on the left and on the right of the shock discontinuity is given by

$$(3.90) \qquad \hat{u}_R = -a\sqrt{\frac{\rho_L}{\rho_R}}, \quad \hat{u}_L = -a\sqrt{\frac{\rho_R}{\rho_L}}.$$

Then, the velocity jump across the shock in the fixed laboratory frame results

$$(3.91) \qquad \hat{u}_R - \hat{u}_L = u_R - u_L = a\frac{\rho_R - \rho_L}{\sqrt{\rho_R \rho_L}},$$

and the speed of the shock wave is equal to

$$(3.92) \qquad s = u_R - \hat{u}_R = u_R + a\sqrt{\frac{\rho_L}{\rho_R}}.$$

The Lax entropy condition (3.51) applied to isothermal gasdynamics requires that the eigenvalue $\lambda_3 (Q_{L,R}) = u_{L,R} + a$ satisfies the following relation

$$(3.93) \qquad \lambda_3 (Q_L) > s > \lambda_3 (Q_R),$$

that is equivalent to

$$(3.94) \qquad \overbrace{u_R - a\frac{\rho_R - \rho_L}{\sqrt{\rho_R \rho_L}}}^{=u_L} + a > u_R + a\sqrt{\frac{\rho_L}{\rho_R}} > u_R + a, \quad \Rightarrow \quad \rho_L > \rho_R.$$

The Lax entropy condition therefore requires that the shock wave produces a compression in the fluid, i.e. an *increase* of density across the shock wave.

**3.2. Rarefaction waves and Riemann invariants.** For a general quasilinear hyperbolic system

$$(3.95) \qquad \frac{\partial Q}{\partial t} + A(Q)\frac{\partial Q}{\partial x} = 0,$$

we seek self-similar solutions so that

$$(3.96) \qquad Q(x,t) = v(\xi), \quad \text{with } \xi = \frac{x}{t}.$$

In terms of the self-similar variable $v$ the PDE results

$$(3.97) \qquad \frac{\partial v}{\partial t} + A(v)\frac{\partial v}{\partial x} = 0.$$

Using the chain–rule, we obtain

$$(3.98) \qquad \frac{dv}{d\xi}\frac{\partial \xi}{\partial t} + A(v)\frac{dv}{d\xi}\frac{\partial \xi}{\partial x} = -v'\frac{x}{t^2} + A(v)\frac{v'}{t} = (A(v) - \xi I)v' = 0.$$

The trivial solution of (3.98) is $v' = 0$. Non–trivial solutions are obtained for the case where $\xi$ is an eigenvalue of $A$ and $v'$ is the associated eigenvector of $A$, i.e. when

$$(3.99) \qquad \xi = \lambda_k(v(\xi)), \qquad \text{and} \qquad v' = r_k(v(\xi)) \cdot \alpha(\xi).$$

Here, $\alpha(\xi)$ is an appropriate scaling factor, since the eigenvectors of $A$ are determined only up to a multiplicative constant.

Following the method of characteristics, a characteristic curve $C_k$ is defined by the ODE

$$(3.100) \qquad \frac{dx}{dt} = \lambda_k(Q(x,t)) = \lambda_k(v(\xi)),$$

where the index $k$ denotes the number of the characteristic field under consideration. If we fix the value of $\xi$, the variable $v(\xi)$ and the eigenvalue $\lambda_k$ are constants. In the case of a simple centered rarefaction wave, the solution is thus constant along the characteristic $C_k$, which therefore is a straight line, as in the linear and the nonlinear scalar case and as plotted in Fig. 12.

We consider the first characteristic field ($\lambda_1 = u - a$) for isothermal gasdynamics. Let us introduce a change of coordinates $\xi \to \zeta$ so that $\alpha := 1$. Then, the system is given by

$$(3.101) \qquad \frac{dv(\zeta)}{d\zeta} = r_1(v(\zeta)), \quad \Rightarrow \quad \frac{d}{d\zeta}\begin{pmatrix} \rho \\ \rho u \\ \rho \psi \end{pmatrix} = \begin{pmatrix} 1 \\ u - a \\ \psi \end{pmatrix}.$$

This leads to the following relations

$$(3.102) \quad \frac{d\rho}{d\zeta} = 1, \quad \frac{d}{d\zeta}(\rho u) = u\frac{d\rho}{d\zeta} + \rho\frac{du}{d\zeta} = u - a, \quad \frac{d}{d\zeta}(\rho\psi) = \psi\frac{d\rho}{d\zeta} + \rho\frac{d\psi}{d\zeta} = \psi.$$

Hence, we integrate across the first wave and obtain

$$(3.103) \qquad \begin{array}{ll} \int_{\rho_L}^{\rho} d\rho = \int_0^1 d\zeta, & \Rightarrow \quad \rho(\zeta) = \rho_L + \zeta, \\ \int_{u_L}^{u} du = \int_0^1 -\frac{a}{\rho(\zeta)}d\zeta = \int_0^1 -\frac{a}{\rho_L + \zeta}d\zeta, & \Rightarrow \quad u(\zeta) = u_L + a\ln\frac{\rho_L}{\rho(\zeta)}, \\ \int_{\psi_L}^{\psi} d\psi = 0, & \Rightarrow \quad \psi = \psi_L = const, \end{array}$$
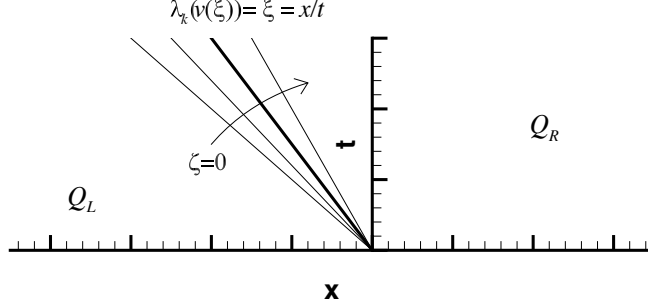
FIGURE 12. Sketch of an isolated centered rarefaction fan, for which the characteristic curves consist of straight lines.

where $L$ indicates the left state. From (3.103) we can derive the so–called *Riemann invariants* for the system of isothermal gasdynamics. They consist of the following relations that hold true across the wave structure

$$(3.104) \qquad u + a \ln \rho = u_L + a \ln \rho_L := w_1^1 = const., \quad \psi = \psi_L := w_1^2 = const.$$

Here, the notation $w_k^j$ denotes the $j$–th Riemann invariant associated with the $k$–th characteristic field. To compute the solution $Q = (\rho, \rho u, \rho \psi)$ in a particular point $\xi$ along a characteristic curve we use the equation $\xi = \lambda_1 = u - a$, see Fig. 12. Using (3.104) one derives the primitive variables

$$(3.105) \qquad\qquad u = a + \xi, \quad \rho = \rho_L e^{\frac{u_L - a - \xi}{a}}, \quad \psi = \psi_L.$$

We now consider the contact discontinuity with $\lambda_2 = u$, obtaining the following system of ODE across the wave structure:

$$(3.106) \qquad\qquad \frac{d}{d\zeta} \begin{pmatrix} \rho \\ \rho u \\ \rho \psi \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix},$$

which correspond to

$$(3.107) \quad \frac{d\rho}{d\zeta} = 0, \quad \frac{d}{d\zeta}(\rho u) = u \frac{d\rho}{d\zeta} + \rho \frac{du}{d\zeta} = 0, \quad \frac{d}{d\zeta}(\rho\psi) = \psi \frac{d\rho}{d\zeta} + \rho \frac{d\psi}{d\zeta} = 1.$$

With easy algebraic manipulations the Riemann invariants are defined by

$$(3.108) \qquad\qquad \rho = \rho_L := w_2^1 = const., \quad u = u_L := w_2^2 = const.$$

This result is also consistent with the Rankine–Hugoniot relations. Let us consider a Riemann problem with the initial states, such that $u_L = u_R$ and $\rho_L = \rho_R$. We verify that the Rankine–Hugoniot conditions (3.56) are satisfied:

$$(3.109) \qquad u \begin{pmatrix} \rho_R - \rho_L \\ \rho_R u_R - \rho_L u_L \\ \rho_R \psi_R - \rho_L \psi_L \end{pmatrix} = \begin{pmatrix} \rho_R u_R \\ \rho_R u_R^2 + a\rho_R^2 \\ \rho_R u_R \psi_R \end{pmatrix} - \begin{pmatrix} \rho_L u_L \\ \rho_L u_L + a\rho_L^2 \\ \rho_L u_L \psi_L \end{pmatrix}.$$

We found that the velocity $u$ is a Riemann invariant for the second characteristic field, which is identical with the eigenvalue $\lambda_2 = u$. We will analyze this feature in
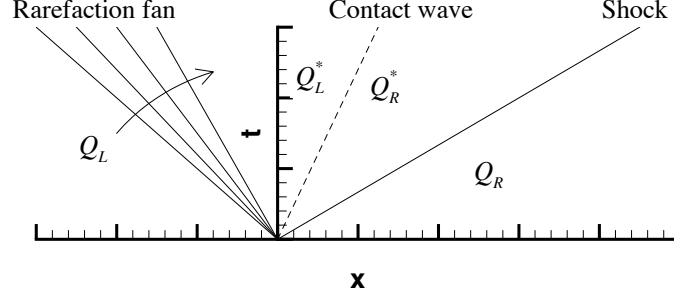
FIGURE 13. The plot for the Riemann solution.

a more general setting in the following. In general, a Riemann invariant $w_k(Q) = w_k(v(\xi))$ associated with the characteristic field $k$ satisfies

$$(3.110) \qquad \frac{d}{d\xi} w_k = 0,$$

hence, with the chain rule we obtain

$$(3.111) \qquad \frac{\partial w_k}{\partial Q} \cdot \frac{dQ}{d\xi} = 0, \qquad \Rightarrow \qquad \nabla w_k \cdot v' = 0,$$

since $dQ/d\xi = v'$. Since $v' = r_k(v(\xi)) \cdot \alpha(\xi)$, it follows that

$$(3.112) \qquad \nabla w_k \cdot r_k = 0.$$

Comparing with the definition of linear degenerate fields (3.57), it results that the eigenvalue $\lambda_k$ is a Riemann invariant, i.e. $w_k = \lambda_k$, if the field $k$ is linearly degenerate.

**3.3. Solution of the Riemann problem.** We seek the exact solution of the Riemann problem for isothermal gasdynamics. A typical wave structure is sketched in Fig. 13.

Across the rarefaction fan the Riemann invariants are constant. Therefore, we use the relations (3.104) to compute the state $Q_L^*$ in the star region on the left of the contact discontinuity as

$$(3.113) \qquad u_L^* + a \ln \rho_L^* = u_L + a \ln \rho_L, \quad \psi_L^* = \psi_L,$$

that leads to

$$(3.114) \qquad \rho_L^* = \rho_L e^{\frac{u_L - u_L^*}{a}}, \quad \psi_L^* = \psi_L.$$

Similarly, across the contact wave the density and the velocity are defined by

$$(3.115) \qquad \rho_R^* = \rho_L = \rho^*, \quad u_R^* = u_L = u^*.$$

The Rankine–Hugoniot conditions (3.56) across the shock give

$$(3.116) \qquad u_R^* = u_R + a \frac{\rho_R^* - \rho_R}{\sqrt{\rho_R^* \rho_R}}, \quad s = u_R + a \sqrt{\frac{\rho_R^*}{\rho_R}}, \quad \psi_R^* = \psi_R.$$

Through algebraic manipulations we obtain the following relations

$$(3.117) \qquad u_L + a \ln \frac{\rho_L}{\rho^*} = u_R + a \ln \frac{\rho^* - \rho_R}{\sqrt{\rho^* \rho_R}},$$

$$(3.118) \qquad g\left(\rho^*\right) = a \ln \frac{\rho^*}{\rho_L} + a \ln \frac{\rho^* - \rho_R}{\sqrt{\rho^* \rho_R}} + u_R - u_L = 0.$$

Note that (3.117) and (3.118) hold true only when the exact solution follows the configuration plotted in Fig. 13. For a general wave structure we have to find the zero of the following nonlinear function in terms of the density variable $\rho^*$

$$(3.119) \qquad g\left(\rho^*\right) = \phi_L + \phi_R + u_R - u_L = 0,$$

where the value of $\phi_{L,R}$ depends on the nature of the nonlinear wave (shock or rarefaction):

$$(3.120) \qquad \phi_L = \begin{cases} a \ln \frac{\rho^*}{\rho_L}, & \text{if } \rho^* \le \rho_L, \quad \text{(rarefaction)} \\ a \frac{\rho^* - \rho_L}{\sqrt{\rho^* \rho_L}}, & \text{if } \rho^* > \rho_L, \qquad \text{(shock)} \end{cases}$$

$$(3.121) \qquad \phi_R = \begin{cases} a \ln \frac{\rho^*}{\rho_R}, & \text{if } \rho^* \le \rho_R, \quad \text{(rarefaction)} \\ a \frac{\rho^* - \rho_R}{\sqrt{\rho^* \rho_R}}, & \text{if } \rho^* > \rho_R. \qquad \text{(shock)} \end{cases}$$

Solving the equation (3.119) requires an iterative procedure. A wide used and proper choice consists of the Newton method, that assures quick convergence to the sought solution $\rho^*$.

## 4. The exact solution of the Riemann problem for the shallow water equations

We consider the shallow water equations that describe the evolution of a homogeneous, incompressible free surface flow under the condition that the horizontal length scale is much greater in comparison with the vertical one. This leads to neglecting the vertical accelerations of the fluid and to a hydrostatic pressure distribution along the water depth. Then, the mathematical model is strongly simplified and results as

$$\begin{aligned} \frac{\partial h}{\partial t} + \frac{\partial}{\partial x}(hu) &= 0, \\ (3.122) \qquad \frac{\partial hu}{\partial t} + \frac{\partial}{\partial x}(hu^2 + \bar{p}) &= 0, \end{aligned}$$

where $h$ is the depth of the fluid, $u$ is the velocity and $\bar{p}$ is the pressure according to the assumption of hydrostatic distribution, $\bar{p} = \frac{1}{2}gh^2$. The system of equations can be written in matrix form

$$\frac{\partial Q}{\partial t} + \frac{\partial}{\partial x}f(Q) = 0, \qquad Q \in \Omega_Q \subset \mathbb{R}^2, \quad t \in \mathbb{R}_0^+, \quad x \in \mathbb{R}.$$

where the vector of conserved variables $Q$ and the flux vector $f$ are defined as

$$(3.123) \qquad Q = \begin{pmatrix} h \\ hu \end{pmatrix} = \begin{pmatrix} q_1 \\ q_2 \end{pmatrix}, f = \begin{pmatrix} hu \\ hu^2 + \bar{p} \end{pmatrix} = \begin{pmatrix} q_2 \\ q_2^2/q_1 + \frac{1}{2}gq_1^2 \end{pmatrix}.$$

The Jacobian matrix $A$ of the flux $f$ with respect to the vector of conserved variables $Q$ is given by

$$(3.124) \qquad A = \begin{pmatrix} 0 & 1 \\ -q_2^2/q_1^2 + a^2 & 2q_2/q_1 \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ a^2 - u^2 & 2u \end{pmatrix},$$

where $a = \sqrt{gh}$ denotes celerity of the gravity waves that propagate on the water surface. As computed in example 27, the eigenvalues are given by $\lambda_{1,2} = u \mp a$ and the corresponding matrix of right eigenvectors $R$ and its inverse $R^{-1}$ are the following

$$(3.125) \qquad R = \begin{pmatrix} 1 & 1 \\ u-a & u+a \end{pmatrix}, \qquad R^{-1} = \frac{1}{2a}\begin{pmatrix} a+u & -1 \\ a-u & 1 \end{pmatrix}.$$

We first compute the nature of the characteristic fields and find for both eigenvalues

$$(3.126) \quad \begin{aligned} \nabla\lambda_1 &= \frac{\partial}{\partial Q}\left(\frac{q_2}{q_1} - \sqrt{gq_1}\right) = \left(-\frac{q_2}{q_1^2} - \frac{1}{2}\sqrt{\frac{g}{q_1}}, \frac{1}{q_1}\right)^T = \left(-\frac{u}{h} - \frac{1}{2}\sqrt{\frac{g}{h}}, \frac{1}{h}\right)^T, \\ \nabla\lambda_2 &= \frac{\partial}{\partial Q}\left(\frac{q_2}{q_1} + \sqrt{gq_1}\right) = \left(-\frac{q_2}{q_1^2} + \frac{1}{2}\sqrt{\frac{g}{q_1}}, \frac{1}{q_1}\right)^T = \left(-\frac{u}{h} + \frac{1}{2}\sqrt{\frac{g}{h}}, \frac{1}{h}\right)^T, \end{aligned}$$

that the fields are *genuinely nonlinear*

$$(3.127) \quad \begin{aligned} \nabla\lambda_1 \cdot \vec{r}_1 &= \left(-\frac{u}{h} - \frac{1}{2}\sqrt{\frac{g}{h}}, \frac{1}{h}\right)^T \cdot \left(1, u - \sqrt{gh}\right)^T = -\frac{3}{2}\sqrt{\frac{g}{h}} \neq 0, \quad \forall Q, \\ \nabla\lambda_2 \cdot \vec{r}_2 &= \left(-\frac{u}{h} + \frac{1}{2}\sqrt{\frac{g}{h}}, \frac{1}{h}\right)^T \cdot \left(1, u + \sqrt{gh}\right)^T = +\frac{3}{2}\sqrt{\frac{g}{h}} \neq 0, \quad \forall Q. \end{aligned}$$

The solution of the Riemann problem using the shallow water model produces two waves, namely the rarefaction fan and the shock. Since only one vertical level $h$ is computed, the shallow water equations can not directly encompass any physical quantity that varies with height. Hence, they are especially suitable to model phenomena which have very large length scales compared to the depth $h$, such as rivers and tides.

**4.1. Shock waves and Rankine–Hugoniot conditions.** To achieve the relation associated with the shock we use the Rankine–Hugoniot conditions. We apply the principle of Galilean invariance of Newtonian mechanics and observe the shock wave in a moving frame, such that the linear transformation of the velocity is given by the relation (3.86). Then, the Rankine–Hugoniot conditions (3.56) for the shallow water equations become

$$(3.128) \quad \begin{aligned} h_R\hat{u}_R &= h_L\hat{u}_L := \mu, \\ h_R\hat{u}_R^2 + \frac{1}{2}gh_R^2 &= h_L\hat{u}_L^2 + \frac{1}{2}gh_L^2. \end{aligned}$$

With algebraic manipulations we obtain

$$(3.129) \quad \begin{aligned} \hat{u}_R &= \frac{\mu}{h_R}, \quad \hat{u}_L = \frac{\mu}{h_L}, \\ \mu\left(\hat{u}_R - \hat{u}_L\right) &= \frac{1}{2}g\left(h_L^2 - h_R^2\right), \quad \Rightarrow \quad \mu = -\frac{1}{2}g\frac{h_R^2 - h_L^2}{\hat{u}_R - \hat{u}_L}, \end{aligned}$$

that provide the velocity jump across the shock wave expressed in terms of the initial water depths $h_L$, $h_R$

$$(3.130) \quad \hat{u}_R - \hat{u}_L = u_R - u_L = \sqrt{\frac{1}{2}g\frac{h_L + h_R}{h_Rh_L}}\left(h_R - h_L\right),$$

and the celerity

$$(3.131) \quad s = u_R - \hat{u}_R = u_R + \sqrt{\frac{1}{2}g\frac{h_L}{h_R}\left(h_L + h_R\right)}.$$

To satisfy the Lax entropy condition (3.51) for the shallow water equations we need that

$$(3.132) \quad u_L + \sqrt{gh_L} > s > u_R + \sqrt{gh_R}, \quad \Rightarrow \quad h_L > h_R.$$

Then, the Lax entropy condition requires that the water depth increases across the shock wave.

**4.2. Rarefaction waves and Riemann invariants.** We consider the self–similar solutions associated with (3.122). In the case of the shallow water model, the relations (3.98) lead to the following system

$$(3.133) \quad \frac{d}{d\zeta}\begin{pmatrix} h \\ hu \end{pmatrix} = \begin{pmatrix} 1 \\ \lambda_k \end{pmatrix},$$

with $\alpha := 1$ and $\xi \to \zeta$. Let us assume the solution structure plotted in Fig. 14 with a rarefaction wave on the left and a shock on the right. For the first eigenvalue $\lambda_1 = \zeta = u - a$, with $a = \sqrt{gh}$ we have

$$(3.134) \qquad \frac{dh}{d\zeta} = 1, \quad \frac{d}{d\zeta}(hu) = u\frac{dh}{d\zeta} + h\frac{du}{d\zeta} = u + h\frac{du}{d\zeta} = u - \sqrt{gh}.$$

Hence, we integrate across the first wave structure

$$(3.135)$$
$$\int_{h_L}^{h} dh = \int_0^\zeta d\zeta \quad \Rightarrow \quad h(\zeta) = h_L + \zeta,$$
$$\int_{u_L}^{u} du = -\int_0^\zeta d\zeta \sqrt{\tfrac{g}{h}} = -\int_0^\zeta d\zeta \sqrt{\tfrac{g}{h_L + \zeta}} \quad \Rightarrow \quad u = u_L - \left(2\sqrt{gh(\zeta)} - 2\sqrt{gh_L}\right),$$

where the index $L$ denotes the left state $Q_L = (h_L, h_L u_L)$ of the Riemann problem. The relations above lead to the Riemann invariant that holds true across the wave structure given by a rarefaction fan

$$(3.136) \qquad u^* + 2\sqrt{gh^*} = u_L + 2\sqrt{gh_L}.$$

**4.3. Solution of the Riemann problem.** The exact Riemann solution for the shallow water equations is computed using the Riemann invariants and the Rankine–Hugoniot conditions. We assume a typical solution for the PDE as plotted in 14. It consists of a left rarefaction wave and a right shock. The region inside the wave structures, the so-called star region, is defined by a state $Q^* = (h^*, u^*)$ that is unknown.

The Riemann invariants for the shallow water equations across the left rarefaction fan give

$$(3.137) \qquad u^* + 2\sqrt{gh^*} = u_L + 2\sqrt{gh_L}.$$

The Rankine-Hugoniot conditions provide the jump in the velocity field across the shock wave, as follows

$$(3.138) \qquad u_R - u^* = -\sqrt{\frac{1}{2}g\frac{h^* + h_R}{h^* h_R}}\,(h^* - h_R).$$

With algebraic manipulations the relations (3.137) and (3.138) reduce to the following function $g$ expressed in terms of $h^*$

$$(3.139) \qquad g(h^*) = 2\left(\sqrt{gh^*} - \sqrt{gh_L}\right) + \sqrt{\frac{1}{2}g\frac{h^* + h_R}{h^* h_R}}\,(h^* - h_R) + u_R - u = 0,$$

the root of which provides the state in the star region. Nevertheless, this solution is true only for the wave structure initially assumed. In the general case the state $Q^*$ is the root of

$$(3.140) \qquad g(h^*) = \Phi_L + \Phi_R + u_R - u_L = 0,$$

where the function $\Phi_{L,R}$ is defined by

$$(3.141) \qquad \Phi_{L,R} = \begin{cases} 2\left(\sqrt{gh^*} - \sqrt{gh_{L,R}}\right), & \text{if } h^* \leq h_{L,R}, \\ \sqrt{\frac{1}{2}g\frac{h^* + h_{L,R}}{h^* h_{L,R}}}\,(h^* - h_{L,R}), & \text{if } h^* > h_{L,R}. \end{cases}$$
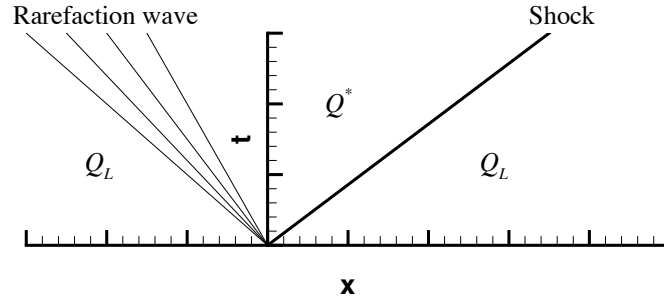
FIGURE 14. The plot for a typical solution of the Riemann problem for the swallow water equations.

# The finite volume method

We consider the following initial–boundary value problem given by

(4.1)
$$
\begin{aligned}
&\text{PDE:} && \tfrac{\partial}{\partial t} Q + \tfrac{\partial}{\partial x} f(Q) = 0, \quad x \in [x_L, x_R], \, Q \in \Omega_Q \subset \mathbb{R}^m, \\
&\text{IC:} && Q(x, 0) = h(x), \\
&\text{BC:} && Q(x_L, t) = B_L(t), \quad Q(x_R, t) = B_R(t),
\end{aligned}
$$

where $Q$ is the vector of the conservative variables, $f(Q)$ is the flux vector, $h$ is the initial condition and $B_{L,R}$ are the boundary conditions on the left and the right. To define the finite volume method we use the integral formulation of the conservation law. In Fig. 1 we show the discretization of the computational domain in space and time. The spatial control volumes are centered in $x_i$ and are defined by $I_i = [x_{i-\frac{1}{2}}; x_{i+\frac{1}{2}}]$, where the boundaries of the interval are equal to $x_{i\mp\frac{1}{2}} = x_i \mp \frac{\Delta x}{2}$, with $\Delta x$ the mesh size. For an equidistant grid, $\Delta x$ is given by $\Delta x = \frac{x_R - x_L}{i_{\max}}$. The temporal control volume is $T_n = [t^n; t^{n+1}]$, with the time step $\Delta t = t^{n+1} - t^n$. Integrating the PDE (4.1) over the space-time control volume $I_i \times T_n$, using the integral form 3.39, we obtain:

(4.2)
$$
\int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} Q(x, t^{n+1}) \, dx = \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} Q(x, t^n) \, dx - \int_{t^n}^{t^{n+1}} \left( f\left(Q(x_{i+\frac{1}{2}}, t)\right) - f\left(Q(x_{i-\frac{1}{2}}, t)\right) \right) dt.
$$

Since the spatial domain is divided into finite control volumes or cells, we can define the cell-average value computed at time $t^n$ as

(4.3)
$$
Q_i^n = \frac{1}{\Delta x} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} Q(x, t^n) \, dx.
$$

Similarly, the time-averaged flux calculated at the cell interface $x_{i+\frac{1}{2}}$ is given by

(4.4)
$$
f_{i+\frac{1}{2}} = \frac{1}{\Delta t} \int_{t^n}^{t^{n+1}} f\left(Q(x_{i+\frac{1}{2}}, t^n)\right) dt.
$$

Substituting the definitions (4.3) and (4.4) into (4.2) we can therefore derive the following *exact* relation, which is nothing else than a reformulation of the principle of integral conservation:

(4.5)
$$
Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x} \left( f_{i+\frac{1}{2}} - f_{i-\frac{1}{2}} \right).
$$

So far, no numerical scheme has been introduced yet and the exact solution of (4.1) satisfies also (4.5) exactly. However, we will now use the formulation (4.5) in order to construct a discrete method for the solution of (4.1). Since the scheme is based on the integral conservation over finite control volumes, it is called a *finite volume* scheme. Remember that the finite volume (FV) approach is philosophically different from a finite difference (FD) scheme. In the FD framework we use the

*strong differential* form of the PDE and approximate derivatives at discrete points $x_i$, $t^n$. In the FV framework, we discretize the *weak integral* form of the PDE. The main ingredient here is to define the numerical flux at the cell interfaces $f_{i\pm\frac{1}{2}}$ as function of the cell averages $Q_i^n$ and $Q_{i\pm1}^n$, since in the finite volume framework only the cell–averaged solution is known. This means that the finite volume scheme does not provide the discrete value of $Q$ at any spatial point $x_i$, but its integral average over the spatial control volume. However, in order to compute the flux at the element interfaces, we need to define the values of $Q$ at the interfaces $x_{i\pm\frac{1}{2}}$. This procedure is called *reconstruction step*. In the simplest of all cases, the solution $Q_i^n$ is supposed to be *piecewise constant* with value $Q_i^n$ within each cell. This piece–wise constant solution inside each cell leads to discontinuities at the cell interfaces $x_{i\pm1}$, see Fig. 2, since the numerical solution has two values at the interface:

$$(4.6) \qquad Q_{i+\frac{1}{2}}^- = Q_i^n, \qquad \text{and} \qquad Q_{i+\frac{1}{2}}^+ = Q_{i+1}^n,$$

where $Q_{i+\frac{1}{2}}^-$ denotes the solution at the left of the interface and $Q_{i+\frac{1}{2}}^+$ denotes the solution at the right of the interface $x_{i+\frac{1}{2}}$. To resolve these discontinuities, we need a so–called *numerical flux* that has to be a function of both states $Q_{i+\frac{1}{2}}^-$ and $Q_{i+\frac{1}{2}}^+$ on the left and on the right of the interface:

$$(4.7) \qquad f_{i+\frac{1}{2}} = f_{i+\frac{1}{2}}\left(Q_{i+\frac{1}{2}}^-, Q_{i+\frac{1}{2}}^+\right) = f_{i+\frac{1}{2}}\left(Q_i^n, Q_{i+1}^n\right).$$

## 1. Basic numerical fluxes

The choice of the numerical flux can be for example:

- The central flux that is given by

$$(4.8) \qquad f_{i+\frac{1}{2}}^c = \frac{1}{2}\left(f\left(Q_{i+1}^n\right) + f\left(Q_i^n\right)\right).$$

- The Lax-Friedrichs flux, defined as follows

$$(4.9) \qquad f_{i+\frac{1}{2}}^{LF} = \frac{1}{2}\left(f\left(Q_{i+1}^n\right) + f\left(Q_i^n\right)\right) - \frac{1}{2}\frac{\Delta x}{\Delta t}\left(Q_{i+1}^n - Q_i^n\right).$$

- The flux of Lax–Wendroff in two possible versions. The first one consists of

$$(4.10) \qquad f_{i+\frac{1}{2}}^{LW1} = \frac{1}{2}\left(f\left(Q_{i+1}^n\right) + f\left(Q_i^n\right)\right) - \frac{1}{2}\frac{\Delta t}{\Delta x}A_{i+\frac{1}{2}}^2\left(Q_{i+1}^n - Q_i^n\right),$$

where the term $A_{i+\frac{1}{2}}$ is an averaged Jacobian matrix defined at $x_{i+\frac{1}{2}}$, computed for example by the arithmetic average of the two Jacobians computed with the left and the right state, respectively. The second Lax–Wendroff–type flux is given by the following two–stage procedure

$$(4.11) \qquad f_{i+\frac{1}{2}}^{LW2} = f\left(Q_{i+\frac{1}{2}}^{LW}\right),$$

where the state $Q_{i+\frac{1}{2}}^{LW}$ is computed from

$$(4.12) \qquad Q_{i+\frac{1}{2}}^{LW} = \frac{1}{2}\left(Q_{i+1}^n + Q_i^n\right) - \frac{1}{2}\frac{\Delta t}{\Delta x}\left(f\left(Q_{i+1}^n\right) - f\left(Q_i^n\right)\right).$$

- Another possible choice consists of the FORCE flux, introduced by Toro and Billet [80]

$$(4.13) \qquad f_{i+\frac{1}{2}}^{FO} = \frac{1}{2}\left(f_{i+\frac{1}{2}}^{LF} + f_{i+\frac{1}{2}}^{LW}\right).$$
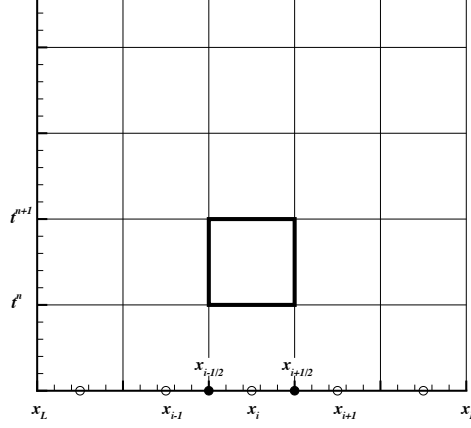
FIGURE 1. The space-time discretization and the space-time control volume $I_i \times T_n$ (bold line) referred to the finite volume method.



FIGURE 2. The assumption of a piece–wise constant solution using first order accurate finite volume schemes. Note the discontinuities at the element interfaces.

- The most accurate flux, but also the most complex one is the Godunov flux (1959), [**37**].

$$(4.14) \qquad f_{i+\frac{1}{2}}^{GO} = f\left(Q_{i+\frac{1}{2}}^{GO}\right),$$

where the so–called Godunov state $Q_{i+\frac{1}{2}}^{GO} = Q^{RP}(0,t)$ is the *exact* solution of a local Riemann problem defined at the cell interface. The local Riemann problem is given by the initial value problem

$$(4.15) \qquad \begin{array}{ll} \text{PDE:} & Q_t + f_x = 0, \\ \text{IC:} & Q(x,0) = \left\{ \begin{array}{ll} Q_i^n, & \text{if } x < 0, \\ Q_{i+1}^n, & \text{if } x \geq 0, \end{array} \right. \end{array}$$

and the exact solution of (4.15) is denoted by $Q^{RP}(x,t)$.

## 2. Properties of finite volume schemes

**2.1. Exact conservation.** A very important feature of finite volume methods is that they are *exactly* conservative, since the variation in time of a quantity, whose evolution obeys a conservation law (4.1), is only given by the exchange through the element interfaces. Hence, a method is conservative when it can be written as follows

$$(4.16) \qquad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x}\left(f_{i+\frac{1}{2}} - f_{i-\frac{1}{2}}\right),$$

FIGURE 3. Convergence to the weak solution according to the Lax-Wendroff theorem (1960), which requires consistency and stability of the conservative numerical scheme.

with $f_{i+\frac{1}{2}} = f_{i+\frac{1}{2}}\left(Q_{i-l}^n, Q_{i-l+1}^n, ..., Q_i^n, Q_{i+1}^n, ..., Q_{i+r}^n\right)$ and $r, l \geq 0$. The finite volume scheme is exactly conservative by construction, which is easy to prove at the aid of the telescopic property of the resulting sum:

$$(4.17) \quad \begin{aligned} \sum_{i=1}^{i_{\max}} \Delta x Q_i^{n+1} &= \left(\sum_{i=1}^{i_{\max}} \Delta x Q_i^n\right) - \Delta t \left(\sum_{i=1}^{i_{\max}} f_{i+\frac{1}{2}} - f_{i-\frac{1}{2}}\right) = \\ &\left(\sum_{i=1}^{i_{\max}} \Delta x Q_i^n\right) - \Delta t \left(\overbrace{f_{i_{\max}+\frac{1}{2}}}^{\text{right BC}} - \overbrace{f_{1-\frac{1}{2}}}^{\text{left BC}}\right). \end{aligned}$$

This exact conservation property is crucial to solve nonlinear hyperbolic conservation laws with discontinuous solutions. According to the theorem of Hou and Le Floch (1994) [**45**] only a conservative method (4.16) is able to capture shock waves properly, with the right velocity $s$ and the right states $Q_L$, $Q_R$. As proven by Hou and Le Floch, a non-conservative scheme converges to the *wrong solution* with wrong shock speed and the wrong states at the shock wave.

**2.2. Consistency with the PDE.** A conservative method (4.16) is consistent if

$$(4.18) \quad f_{i+\frac{1}{2}}\left(Q, Q, ..., Q\right) = f\left(Q\right).$$

**2.3. Convergence to weak solutions.** The Lax–Wendroff theorem (1960) [**53**] states that if the numerical solution of a conservative method (4.16) converges, then it converges to a weak solution of the conservation law. For convergence, the method needs to be stable and consistent, see the sketch of Fig. (3).

**2.4. Monotonicity.** Regarding the property of monotonicity, a necessary condition so that the finite volume scheme is monotone for a scalar PDE $q_t + f_x = 0$ is given by

$$(4.19) \quad \frac{\partial}{\partial q_i^n} f_{i+\frac{1}{2}}\left(q_i^n, q_{i+1}^n\right) \geq 0, \quad \text{and} \quad \frac{\partial}{\partial q_{i+1}^n} f_{i+\frac{1}{2}}\left(q_i^n, q_{i+1}^n\right) \leq 0.$$

Hence, the numerical flux $f_{i+\frac{1}{2}}$ has to be a non-decreasing function of the first argument $q_i^n$ and a non-increasing function of the second one $q_{i+1}^n$. This is proved using the operator $H$ given by

$$(4.20) \quad q_i^{n+1} = H\left(q_{i-1}^n, q_i^n, q_{i+1}^n\right) = q_i^n - \frac{\Delta t}{\Delta x}\left(f_{i+\frac{1}{2}} - f_{i-\frac{1}{2}}\right).$$

The general monotonicity condition (2.83) requires that the numerical flux satisfies the following relations:

$$(4.21) \quad \frac{\partial H}{\partial q_{i-1}^n} = \frac{\Delta t}{\Delta x}\frac{\partial f_{i-\frac{1}{2}}}{\partial q_{i-1}^n} \geq 0, \quad \Rightarrow \quad \frac{\partial f_{i-\frac{1}{2}}}{\partial q_{i-1}^n} \geq 0,$$

and

$$(4.22) \qquad \frac{\partial H}{\partial q_{i+1}^n} = -\frac{\Delta t}{\Delta x}\frac{\partial f_{i+\frac{1}{2}}}{\partial q_{i+1}^n} \geq 0, \qquad \Rightarrow \qquad \frac{\partial f_{i+\frac{1}{2}}}{\partial q_{i+1}^n} \leq 0.$$

EXAMPLE 28. *The aim of this example is to show that the finite volume scheme with the central flux and the central finite difference method coincide in the one-dimensional case. We consider the linear scalar advection equation (3.1), with a constant positive velocity $a \in \mathbb{R}^+$. Using the central flux (4.8) in the finite volume scheme we obtain*

$$(4.23)\quad q_i^{n+1} = q_i^n - \frac{\Delta t}{2\Delta x}\left(aq_{i+1}^n + aq_i^n - aq_{i-1}^n - aq_i^n\right) = q_i^n - \frac{a\Delta t}{2\Delta x}\left(q_{i+1}^n - q_{i-1}^n\right).$$

*The numerical scheme derived above is the same obtained from the finite difference approach with the central flux (1.97).*

*Moreover, according to the definition (4.20) we have*

$$(4.24) \qquad\qquad q_i^{n+1} := H\left(q_{i-1}^n, q_i^n, q_{i+1}^n\right),$$

*and the necessary condition for the monotonicity is **not** satisfied, as can be easily seen below:*

$$(4.25) \qquad \frac{\partial H}{\partial q_i^n} = 1 > 0, \quad \frac{\partial H}{\partial q_{i-1}^n} = \frac{a\Delta t}{2\Delta x} > 0, \quad \frac{\partial H}{\partial q_{i+1}^n} = -\frac{a\Delta t}{2\Delta x} < 0.$$

*Since the partial derivative of $H$ with respect to $q_{i+1}^n$ is always negative, the central scheme is **not** monotone.*

EXAMPLE 29. *We consider the Lax-Friedrichs flux (4.9) in the finite volume case and prove the monotonicity of the resulting scheme, applied to the linear scalar advection equation. The scheme is*

$$
\begin{aligned}
(4.26) \qquad q_i^{n+1} = q_i^n - \frac{\Delta t}{\Delta x}&\left[\left(\frac{1}{2}\left(aq_i^n + aq_{i+1}^n\right) - \frac{1}{2}\frac{\Delta x}{\Delta t}\left(q_{i+1}^n - q_i^n\right)\right) - \right.\\
&\left.\left(\frac{1}{2}\left(aq_i^n + aq_{i-1}^n\right) - \frac{1}{2}\frac{\Delta x}{\Delta t}\left(q_i^n - q_{i-1}^n\right)\right)\right].
\end{aligned}
$$

*Simplifying we obtain*

$$(4.27) \qquad q_i^{n+1} = \frac{1}{2}\left(q_{i+1}^n + q_{i-1}^n\right) - \frac{a\Delta t}{2\Delta x}\left(q_{i+1}^n - q_{i-1}^n\right) := H\left(q_{i-1}^n, q_{i+1}^n\right),$$

*which coincides with the Lax–Friedrichs finite difference scheme. For monotonicity we must have*

$$(4.28) \qquad \frac{\partial H}{\partial q_{i-1}^n} = \frac{1}{2}(1+c) > 0, \quad \frac{\partial H}{\partial q_{i+1}^n} = \frac{1}{2}(1-c) > 0, \quad \forall 0 < c < 1,$$

*with the Courant number $c > 0$. We can conclude that the necessary condition for the monotonicity is satisfied under CFL condition $0 < c \leq 1$.*

EXAMPLE 30. *For the Godunov flux [37] we first have to determine the exact solution of the local Riemann problem*

$$
(4.29) \qquad \left.
\begin{array}{ll}
q_t + aq_x = 0, & (PDE)\\[4pt]
q(x,0) = \left\{\begin{array}{ll} q_i^n, & x \leq 0,\\ q_{i+1}^n, & x > 0. \end{array}\right. & (IC)
\end{array}
\right\} (RP)
$$

*at the element interface $x_{i+\frac{1}{2}}$, which means that we must find the self–similar solution of (4.29) at position $\xi = x/t = 0$. The exact solution of (4.29) in the entire $x - t$ half–plane is given by*

$$(4.30) \qquad\qquad q(x,t) = \left\{\begin{array}{ll} q_i^n, & x/t \leq a,\\ q_{i+1}^n, & x/t > a, \end{array}\right.$$

*and the solution at $\xi = 0$ is*

$$(4.31) \qquad q_{i+\frac{1}{2}}^{GO} = q(x/t = 0) = \begin{cases} q_i^n, & 0 \le a, \\ q_{i+1}^n, & 0 > a. \end{cases}$$

*For $a > 0$ we obtain*

$$(4.32) \qquad q_i^{n+1} = q_i^n - \frac{\Delta t}{\Delta x}(a q_i^n - a q_{i-1}^n),$$

*which is identical with the upwind finite difference scheme. It is easy to prove that the upwind scheme is monotone under CFL condition since for $a > 0$ we get*

$$(4.33) \qquad \frac{\partial H}{\partial q_i^n} = 1 - c > 0, \quad \forall 0 < c < 1, \qquad \frac{\partial H}{\partial q_{i+1}^n} = c > 0.$$

*For negative speed $a$ we obtain the corresponding upwind forward difference and monotonicity can be proven under the same condition.*

## 3. Approximate Riemann solvers

In general, computing the exact solutions of all the local Riemann problems arising at the element interfaces for the Godunov flux is very complex and time–consuming. Hence, a lot of effort has been put into the research on *approximate* Riemann solvers, which still take into account part of the physical wave structure of the Riemann problem at the element interface, but without aiming at its exact solution. Some popular approximate solvers are presented very briefly in the following. For details, the reader is referred to the excellent overview given in the textbook [**79**].

**3.1. Approximate-state Riemann solvers.** The approximate-state Riemann solvers provide an *approximate Godunov state* $Q_{i+\frac{1}{2}}^{AG}$ at the cell interface in order to compute the numerical flux $f_{i+\frac{1}{2}}$ following the original Godunov approach (4.14)

$$(4.34) \qquad f_{i+\frac{1}{2}} = f_{i+\frac{1}{2}}\left(Q_i^n, Q_{i+1}^n\right) = f\left(Q_{i+\frac{1}{2}}^{AG}\right).$$

The approximation carried out to compute the state $Q_{i+\frac{1}{2}}^{AG}$ is not based on the exact resolution of the Riemann problem as proposed by Godunov, but one assumes *a priori* a *fixed wave pattern*, even if this is physically not the correct one. This simplification leads often (but not always) to a direct analytical calculation of the state $Q^*$

3.1.1. *The two-rarefaction approximate–state Riemann solver (TRRS).* The so-called two-rarefaction approximate–state Riemann solver supposes *a priori* that the solution of the Riemann problem consists of two rarefaction fans, as shown in Fig. 4. Here, we consider the particular case of isothermal gasdynamics. From the Riemann invariants we found that the velocity is constant across a contact discontinuity so that $u_L^* = u_R^* = u^*$. In particular, for isothermal gasdynamics the contact wave is in density equilibrium $\rho_L^* = \rho_R^* = \rho^*$. Using the Riemann invariants associated with the two rarefaction waves we obtain the following relation in $\rho^*$

$$(4.35) \qquad a \ln \frac{\rho^*}{\rho_L} + a \ln \frac{\rho^*}{\rho_R} + u_R - u_L = 0,$$

the root of which provides the state $Q^* = (\rho^*, u^*)$ in the star region

$$(4.36) \qquad \rho^* = \sqrt{\rho_L \rho_R \exp\left(\frac{u_L - u_R}{a}\right)}, \quad u^* = u_L + a \ln \frac{\rho_L}{\rho^*}.$$
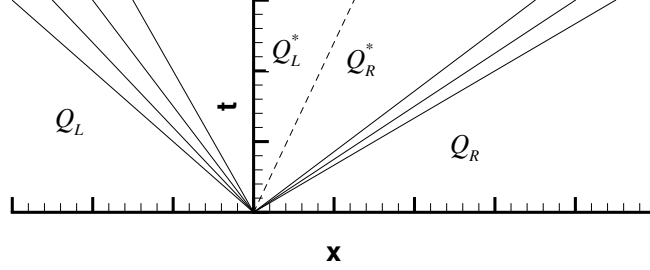
FIGURE 4. A solution of the Riemann problem that consists of two rarefaction waves and a contact discontinuity.



FIGURE 5. A solution of the Riemann problem that consists of two shock waves and a contact discontinuity.

Once the state $Q^*$ is known, one has to compute the approximate state $Q^{AG}_{i+\frac{1}{2}}$ as using the exact Riemann solver, but assuming the fixed structure with two rarefaction waves. Note that the two-rarefaction approximate state Riemann solver is exact when the exact solution actually consists only of two rarefaction waves separated by a contact discontinuity. The two–rarefaction approximate–state Riemann solver is consistent with the entropy condition, as the exact Riemann solver.

   3.1.2. *The two-shock approximate–state Riemann solver (TSRS).* In this section we consider the two-shock approximate–state Riemann solver for isothermal gasdynamics. It assumes a simplified solution of the Riemann problem that consists of two shock waves. Fig. 5 shows the fixed solution that we presuppose. At the contact discontinuity we have $u^*_L = u^*_R = u^*$ and $\rho^*_L = \rho^*_R = \rho^*$. These relations have to be combined with the Rankine-Hugoniot conditions which are valid across the shock waves and we obtain

(4.37)
$$a\frac{\rho^* - \rho_L}{\sqrt{\rho^*\rho_L}} + a\frac{\rho^* - \rho_R}{\sqrt{\rho^*\rho_R}} + \overbrace{u_R - u_L}^{=\Delta u} = 0,$$

FIGURE 6. The solution assumed by the HLL method.

that is equivalent to

$$(4.38) \qquad \rho^* + \frac{\Delta u}{a} \frac{\sqrt{\rho_R \rho_L}}{\sqrt{\rho_R} + \sqrt{\rho_L}} \sqrt{\rho^*} - \sqrt{\rho_L \rho_R} = 0,$$

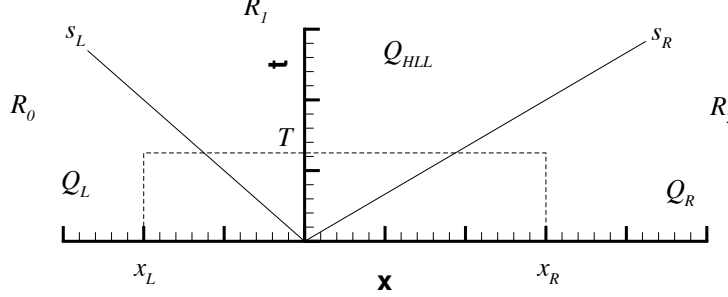Then, the state $Q^*$ is given by

$$(4.39) \qquad \sqrt{\rho^*} = -\frac{1}{2a} \Delta u \frac{\sqrt{\rho_L \rho_R}}{\sqrt{\rho_L} + \sqrt{\rho_R}} + \frac{1}{2} \sqrt{\Delta u^2 \frac{\rho_L \rho_R}{a^2 \left(\sqrt{\rho_L} + \sqrt{\rho_R}\right)^2} + 4\sqrt{\rho_L \rho_R}},$$

and

$$(4.40) \qquad u^* = u_L - a\frac{\rho^* - \rho_L}{\sqrt{\rho^* \rho_L}}.$$

The computation of the approximate state $Q^{AG}_{i+\frac{1}{2}}$ occurs as with the exact Riemann solver, but assuming a fixed structure of the wave pattern of the solution consisting of two shock waves. If the exact solution actually contains two shock waves and a contact discontinuity, this approximate state Riemann solver is exact. Unfortunately, on the contrary of the exact Riemann solver and the two–rarefaction approximate–state Riemann solver, the two–shock approximate–state Riemann solver is *not* consistent with the entropy condition. This means that unphysical solutions such as rarefaction shocks are admitted by the method.

**3.2. The method of Harten-Lax-Leer (HLL).** Harten, Lax and van Leer (1983) [**41**] proposed an approximate Riemann solver assuming a wave structure that is still simpler than the solution of the Riemann problem and only includes the fastest waves, as plotted in Fig. 6. In the HLL approach the solution consists of three regions $R_0$, $R_1$, $R_2$ with piece-wise constant states. In $R_0$ we have the initial state on the left $Q_L$, and in $R_2$ we have the initial right state $Q_R$. In the region $R_1$ we assume the intermediate state $Q_{\text{HLL}}$, which is an approximate one and is still to be defined. It is computed starting from the PDE

$$(4.41) \qquad \frac{\partial Q}{\partial t} + \frac{\partial f}{\partial x} = 0,$$

and integrating over the space-time control volume $[x_L; x_R] \times [0; T]$, sketched in dashed line in Fig. 6

$$(4.42) \qquad \int_{x_L}^{x_R} (Q(x, T) - Q(x, 0))\, dx + \int_0^T (f(Q(x_R, t)) - f(Q(x_L, t)))\, dt = 0,$$

which results in

$$(4.43) \qquad Q_L s_L + Q_{\mathrm{HLL}} (s_R - s_L) - Q_R s_R + f_R - f_L = 0,$$

with $f_L = f(Q_L)$ and $f_R = f(Q_R)$. Finally, the approximate state $Q_{\mathrm{HLL}}$ is given by

$$(4.44) \qquad Q_{\mathrm{HLL}} = \frac{s_R Q_R - s_L Q_L + f_L - f_R}{s_R - s_L}.$$

Note that the HLL method is *not* an approximate state Riemann solver,

$$(4.45) \qquad f_{i+\frac{1}{2}}^{\mathrm{HLL}} \neq f(Q_{\mathrm{HLL}})!$$

Indeed, it does not provide a state in order to compute the numerical flux as physical flux of that state. To compute the numerical flux using the HLL solver we have to apply *again* the PDE in integral form over the smaller space-time control volume $[0; x_R] \times [0; T]$ as follows

$$(4.46) \qquad \int_0^{x_R} (Q(x, T) - Q(x, 0))\, dx + \int_0^T (f(Q(x_R, t)) - f(Q(0, t)))\, dt = 0,$$

that leads to

$$(4.47) \qquad Q_{\mathrm{HLL}}\, s_R T + Q_R (x_R - s_R T) - Q_R x_R + T f_R - T f_{i+\frac{1}{2}}^{\mathrm{HLL}} = 0.$$

We finally obtain the numerical flux at the cell interface as

$$(4.48) \qquad f_{i+\frac{1}{2}}^{\mathrm{HLL}} = f_R + s_R (Q_{\mathrm{HLL}} - Q_R).$$

Substituting (4.44) into the relation above we obtain the numerical flux $f_{i+\frac{1}{2}}^{\mathrm{HLL}}$

$$(4.49) \qquad f_{i+\frac{1}{2}}^{\mathrm{HLL}} = \frac{s_R f_L - s_L f_R + s_L s_R (Q_R - Q_L)}{s_R - s_L}.$$

However, in the HLL method we still need two estimates for the two wave speeds $s_L$ and $s_R$. A simple choice has been proposed by Davis, as follows

$$(4.50) \qquad s_L = \min(0, u_L - a, u_R - a), \quad s_R = \max(0, u_L + a, u_R + a).$$

Another option is based on the assumption of a single wave speed estimate

$$(4.51) \qquad s_L = -s_{\max}, \quad s_R = +s_{\max}, \quad s_{\max} = \max(|u_L \pm a|, |u_R \pm a|),$$

that leads to the so-called Rusanov flux or local Lax–Friedrichs flux

$$(4.52)$$
$$f_{i+\frac{1}{2}}^{R} = \frac{s_{\max} f_L + s_{\max} f_R - s_{\max}^2 (Q_R - Q_L)}{2 s_{\max}} = \frac{1}{2}(f_L + f_R) - \frac{1}{2} s_{\max} (Q_R - Q_L).$$

Note that using the following assumption

$$(4.53) \qquad s_L = -\frac{\Delta x}{\Delta t}, \quad s_R = +\frac{\Delta x}{\Delta t},$$

the HLL method becomes the classical Lax-Friedrichs flux (4.9).

Since the intermediate waves are ignored, the solution of the HLL Riemann solver produces excessive smearing of contact waves (and of shear–waves in multi–dimensional computations). Some modifications have been proposed in order to improve the HLL scheme. One is the HLLE scheme that uses a different estimate of the velocities $s_L$ and $s_R$ following the Einfeldt approach, for details see [**29**]. The HLLEM method assumes no longer a constant intermediate state, but a *piecewise linear* distribution. This approach allows to resolve contact discontinuities according to Einfeld et al. [**30**]. Among all the possible HLL schemes, perhaps the most commonly used approximate Riemann solver of the HLL type is the HLLC method proposed by Toro et al. [**82**], which captures steady contact waves exactly.
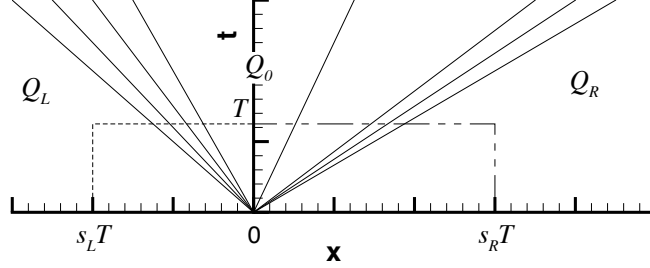
FIGURE 7. The space-time control volumes used to construct the Roe flux.

**3.3. The Roe method.** This approximate Riemann solver [**69**] is based on Roe's ingenious idea to apply a proper linearisation of the nonlinear Riemann problem (1.66) - (1.67). Using the quasi–linear form (3.95) of the PDE, Roe proposed to use some approximate Jacobian matrix in terms of the initial data $Q_L$ and $Q_R$, such that

$$(4.54) \qquad \tilde{A} = \tilde{A}(Q_L, Q_R).$$

The matrix $\tilde{A}$ is called the *Roe matrix*. In order to obtain the Roe flux, a *linearized* Riemann problem with initial data $Q_L$ and $Q_R$ and fixed system matrix $\tilde{A}$ is considered:

$$(4.55) \qquad \begin{aligned} &\frac{\partial \bar{Q}}{\partial t} + \tilde{A}\frac{\partial \bar{Q}}{\partial x} = 0, \\ &\bar{Q}(x,0) = \begin{cases} Q_L, & \text{if } x < 0 \\ Q_R, & \text{if } x \geq 0. \end{cases} \end{aligned}$$

For the Roe matrix $\tilde{A}$ the following properties are enforced:

- Hyperbolicity. It is required that the Roe matrix $\tilde{A}$ has $m$ real eigenvalues and $m$ linear independent eigenvectors.
- Consistency of $\tilde{A}$ with the Jacobian matrix $A(Q)$, that is

$$(4.56) \qquad \tilde{A}(Q,Q) = A(Q).$$

- Conservation.

$$(4.57) \qquad \tilde{A}(Q_L, Q_R) \cdot (Q_R - Q_L) = f(Q_R) - f(Q_L).$$

The intercell flux is given by a similar procedure used in the HLL approach. We consider the original PDE (4.41) in conservative form and integrate twice, first over the space-time control volume on the left $[s_L T; 0] \times [0; T]$, then over the space-time control volume on the right $[0; s_R T] \times [0; T]$, see Fig. 7. We obtain

$$(4.58) \qquad \int_{s_L T}^{0} Q(x,T)\, dx - (0 - s_L T) Q_L + T f_{i+\frac{1}{2}} - T f_L = 0,$$

$$(4.59) \qquad \int_{0}^{s_R T} Q(x,T)\, dx - (s_R T - 0) Q_R + T f_R - T f_{i+\frac{1}{2}} = 0.$$

In order to determine the integrals in the above expressions, the exact solution $\bar{Q}$ of the linearized Riemann problem (4.55) is used. The integration over the left and the right control volumes gives

$$(4.60) \qquad \int_{s_L T}^{0} \bar{Q}(x, T) \, dx - (0 - s_L T) Q_L + T\bar{F}(\bar{Q}_0) - T\bar{F}_L = 0,$$

$$(4.61) \qquad \int_{0}^{s_R T} \bar{Q}(x, T) \, dx - (s_R T - 0) Q_R + T\bar{F}_R - T\bar{F}(\bar{Q}_0) = 0.$$

with $\bar{F}_{L,R} = \tilde{A} \cdot \bar{Q}_{L,R}$ and $\bar{Q}_0 = \bar{Q}(0, T)$. Note that $\bar{Q}_0$ is the exact solution of the linear Riemann problem (4.55) for $\frac{x}{t} = 0$. We now introduce the following hypothesis:

$$(4.62)$$
$$\int_{s_L T}^{0} Q(x, T) \, dx = \int_{s_L T}^{0} \bar{Q}(x, T) \, dx, \quad \int_{0}^{s_R T} Q(x, T) \, dx = \int_{0}^{s_R T} \bar{Q}(x, T) \, dx,$$

Then, substituting into (4.60) and (4.61) we have

$$(4.63) \quad Tf_R - Tf_{i+\frac{1}{2}} - T\bar{F}_R + T\bar{F}(\bar{Q}_0) = 0, \quad \Rightarrow \quad f_{i+\frac{1}{2}} = f_R + \tilde{A}(\bar{Q}_0 - Q_R),$$

$$(4.64) \quad Tf_{i+\frac{1}{2}} - Tf_L - T\bar{F}(\bar{Q}_0 + T\bar{F}_L) = 0, \quad \Rightarrow \quad f_{i+\frac{1}{2}} = f_L + \tilde{A}(\bar{Q}_0 - Q_L).$$

Combining (4.63) with (4.64) we obtain the Roe flux as

$$(4.65) \qquad f_{i+\frac{1}{2}} = \frac{1}{2}(f_R + f_L) + \tilde{A}\bar{Q}_0 - \frac{1}{2}\tilde{A}(Q_R + Q_L).$$

The exact solution of the linearized Riemann problem (4.55) at $x/t = 0$ denoted by $\bar{Q}_0$ is given according to (1.72) by

$$(4.66) \qquad \bar{Q}_0 = \frac{1}{2}\tilde{R}\left(I + \operatorname{sign}\left(\tilde{\Lambda}\right)\right)\tilde{R}^{-1}Q_L + \frac{1}{2}\tilde{R}\left(1 - \operatorname{sign}(\tilde{\Lambda})\right)\tilde{R}^{-1}Q_R,$$

For the Roe flux we need the product $\tilde{A}\bar{Q}_0$, which is

$$(4.67)$$
$$\tilde{A}\bar{Q}_0 = \tilde{R}\tilde{\Lambda}\tilde{R}^{-1}\bar{Q}_0 = \frac{1}{2}\tilde{R}\left(\tilde{\Lambda} + \tilde{\Lambda}\operatorname{sign}(\tilde{\Lambda})\right)\tilde{R}^{-1}Q_L + \frac{1}{2}\tilde{R}\left(\tilde{\Lambda} - \tilde{\Lambda}\operatorname{sign}(\tilde{\Lambda})\right)\tilde{R}^{-1}Q_R.$$

With the definition of the absolute value $|\tilde{\Lambda}| = \tilde{\Lambda}\operatorname{sign}(\tilde{\Lambda})$ we have

$$(4.68) \qquad \tilde{A}\bar{Q}_0 = \frac{1}{2}\tilde{R}\left(\tilde{\Lambda} + \left|\tilde{\Lambda}\right|\right)\tilde{R}^{-1}Q_L + \frac{1}{2}\tilde{R}\left(\tilde{\Lambda} - \left|\tilde{\Lambda}\right|\right)\tilde{R}^{-1}Q_R,$$

which can be written more compactly using the notation $|A| = R|\Lambda|R^{-1}$ as

$$(4.69) \qquad \tilde{A}\bar{Q}_0 = \frac{1}{2}\left(\tilde{A} + |A|\right)Q_L + \frac{1}{2}\left(\tilde{A} - |A|\right)Q_R,$$

or

$$(4.70) \qquad \tilde{A}\bar{Q}_0 = \frac{1}{2}\tilde{A}(Q_R + Q_L) - \frac{1}{2}\left|\tilde{A}\right|(Q_R - Q_L),$$

Then, the Roe flux finally results as

$$(4.71) \qquad f_{i+\frac{1}{2}}^{\text{Roe}} = \frac{1}{2}(f(Q_R) + f(Q_L)) - \frac{1}{2}\left|\tilde{A}\right|(Q_R - Q_L).$$

This flux still requires the calculation of the Roe matrix $\tilde{A}$, which can be very complex. One usually applies the property (4.57) of matrix $\tilde{A}$ directly and uses a proper substitution of variables (parameter vector) in order to compute the so–called Roe state $\tilde{Q}$. Then, the Roe matrix $\tilde{A}$ is computed as the Jacobian of the PDE, evaluated at the Roe state

$$(4.72) \qquad \tilde{A}(Q_L, Q_R) = A\left(\tilde{Q}(Q_L, Q_R)\right),$$

with $\tilde{Q}(Q, Q) = Q$. This assures that also the properties of system hyperbolicity and consistency are satisfied. Unfortunately, the Roe method is *not* entropy satisfying. Harten introduced a technique, the so–called Harten entropy–fix, see [**79**] for details, to avoid this problem. An important feature of Roe's scheme is that it resolves exactly all stationary discontinuities.

EXAMPLE 31. *We use the property* (4.57) *to compute the Roe state* $\tilde{Q}$ *for isothermal gasdynamics assuming that* $\tilde{A} = A\left(\tilde{Q}(Q_L, Q_R)\right)$

(4.73)
$$\begin{pmatrix} 0 & 1 & 0 \\ a^2 - \tilde{u}^2 & 2\tilde{u} & 0 \\ -\tilde{u}\tilde{\psi} & \tilde{\psi} & \tilde{u} \end{pmatrix} \begin{pmatrix} \rho_R - \rho_L \\ \rho_R u_R - \rho_L u_L \\ \rho_R \psi_R - \rho_L \psi_L \end{pmatrix} = \begin{pmatrix} \rho_R u_R - \rho_L u_L \\ \rho_R u_R^2 + a^2 \rho_R - \rho_L u_L^2 - a^2 \rho_L \\ \rho_R u_R \psi_R - \rho_L u_L \psi_L \end{pmatrix}.$$

*Solving the system above we obtain the Roe state for isothermal gasdynamics*

(4.74)    $(\tilde{\rho} = \sqrt{\rho_L \rho_R}), \quad \tilde{u} = \dfrac{u_L \sqrt{\rho_L} + u_R \sqrt{\rho_R}}{\sqrt{\rho_L} + \sqrt{\rho_R}}, \quad \tilde{\psi} = \dfrac{\psi_L \sqrt{\rho_L} + \psi_R \sqrt{\rho_R}}{\sqrt{\rho_L} + \sqrt{\rho_R}}.$

As an alternative, Toumi (1992) [**86**] proposed to construct the Roe matrix $\tilde{A}$ by using an integral (weak) formulation of the conservation property (4.57), as follows

(4.75)    $\tilde{A}(Q_R - Q_L) = \displaystyle\int_{Q_L}^{Q_R} A(Q)\, dQ = \int_{Q_L}^{Q_R} \frac{\partial f}{\partial Q}\, dQ = \int_{f_L}^{f_R} df = f_R - f_L.$

Following a particular integration path, namely the simple segment–path

(4.76)    $Q(s) = Q_L + s(Q_R - Q_L), \quad 0 \le s \le 1,$

we obtain

(4.77)    $\tilde{A}(Q_R - Q_L) = \displaystyle\int_{Q_L}^{Q_R} A(Q)\, dQ = \left(\int_0^1 A(Q(s))\, ds\right)(Q_R - Q_L),$

that is equivalent to

(4.78)    $\tilde{A} = \displaystyle\int_0^1 A(Q(s))\, ds.$

With this weak formulation based on the segment path, the resulting Roe–type flux becomes

(4.79)    $f_{i+\frac{1}{2}}^{\text{Roe}'} = \dfrac{1}{2}(f(Q_R) + f(Q_L)) - \dfrac{1}{2}\left|\int_0^1 A(Q(s))\, ds\right|(Q_R - Q_L).$

The integral in (4.79) can be numerically approximated using any classical quadrature formula, such as the trapezoidal rule, or standard Newton-Cotes or Gauss formulae. The purely numerical computation of the Roe matrix has first been proposed in [**23, 28, 27, 8, 7**]. Note that the Roe matrix (4.78) in integral form satisfies the consistency (4.56) and the conservation (4.57) properties by construction, but the hyperbolicity of the system (the first condition for the Roe matrix) is *not* generally guaranteed. Moreover, also this version of the Roe method is *not* consistent with the entropy condition.

**3.4. The Osher method.** Osher and Solomon (1982) [**62**] proposed the following formulation of the numerical flux

(4.80)    $f_{i+\frac{1}{2}}^{OS} = \dfrac{1}{2}(f(Q_R) + f(Q_L)) - \dfrac{1}{2}\displaystyle\int_{Q_L}^{Q_R} |A(Q)|\, dQ.$

Note that the scheme above is differentiable due to its integral formulation. The original method of Osher and Solomon for nonlinear hyperbolic systems is very
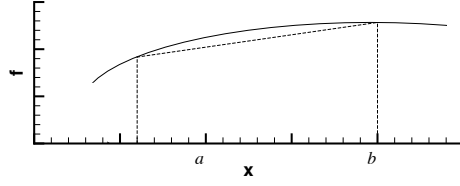
FIGURE 8. The approximation (dashed line) for the integral of a function (solid line) through the trapezoidal rule.

complex, since the exact computation of the integral for the dissipation term is not a trivial task. For more details and test problems see [**79**]. In [**27**], we first proposed a simple segment–path defined by (4.76) to calculate the integral in (4.80), as follows

$$(4.81) \qquad f^{OS}_{i+\frac{1}{2}} = \frac{1}{2} \left( f\left( Q_R \right) + f\left( Q_L \right) \right) - \frac{1}{2} \left( \int_0^1 \left| A\left( Q\left( s \right) \right) \right| ds \right) \left( Q_R - Q_L \right).$$

Compare the new formulation of the Osher–Solomon flux (4.81) with the Roe flux given by (4.79).

It is to remark that the method of Osher and Solomon is consistent with the entropy principle and solves the contact waves exactly. The hyperbolicity of the Jacobian matrix $A\left( Q\left( s \right) \right)$ is guaranteed. Moreover, the implementation is straightforward when a numerical quadrature formula for the integral approximation is used. A choice consists of the trapezoidal rule, given by

$$(4.82) \qquad \int_a^b f\left( x \right) dx \approx \frac{1}{2} \left( b - a \right) \left( f\left( a \right) + f\left( b \right) \right).$$

A sketch of the basic idea for the integral approximation is plotted in Fig. 8. With the trapezoidal rule the Osher–Solomon–type flux results in

$$(4.83) \qquad f^{OS}_{i+\frac{1}{2}} \approx \frac{1}{2} \left( f\left( Q_R \right) + f\left( Q_L \right) \right) - \frac{1}{2} \frac{\left| A\left( Q_L \right) \right| + \left| A\left( Q_R \right) \right|}{2} \left( Q_R - Q_L \right).$$

Similarly, for the Roe flux (4.79) one would have

$$(4.84) \qquad f^{R}_{i+\frac{1}{2}} \approx \frac{1}{2} \left( f\left( Q_R \right) + f\left( Q_L \right) \right) - \frac{1}{2} \left| \frac{A\left( Q_L \right) + A\left( Q_R \right)}{2} \right| \left( Q_R - Q_L \right).$$

## 4. Total Variation Diminishing (TVD) Schemes

Due to the Godunov theorem treated in Section 4, there exists no *linear* method of order of accuracy greater or equal than two that also assures the *monotonicity* of the scheme. To overcome the constraint expressed in this theorem we use the fact that the theorem only applies to linear methods. To construct monotone methods of second order of accuracy, we therefore have to develop *nonlinear* methods.

**4.1. Definitions and the TVD criterion of Harten.** The total variation (TV) of a discrete function $q_h^n(x_h)$ is defined as

$$(4.85) \qquad TV(q_h^n) = \sum_{i=-\infty}^{\infty} |q_{i+1}^n - q_i^n|,$$

with $q_i^n \to const.$ for $i \to \pm\infty$. A numerical method

$$(4.86) \qquad q_i^{n+1} = H(q_{i-l}^n, ...q_i^n, ...q_{i+r}^n)$$

is called TVD if

$$(4.87) \qquad \mathrm{TV}(q_h^{n+1}) \leq \mathrm{TV}(q_h^n).$$

THEOREM 2. *(Harten 1983). The set of monotone methods $S_{\mathrm{mon}}$ is a subset of TVD methods $S_{\mathrm{TVD}}$.*

$$(4.88) \qquad\qquad\qquad S_{\mathrm{mon}} \subset S_{\mathrm{TVD}}.$$

PROOF. For the proof see [**39**].                                              $\square$

THEOREM 3. *(Harten 1983). A numerical method of the form*

$$(4.89) \qquad\qquad q_i^{n+1} = q_i^n - C_{i-\frac{1}{2}}\Delta q_{i-\frac{1}{2}}^n + D_{i+\frac{1}{2}}\Delta q_{i+\frac{1}{2}}^n,$$

*with $\Delta q_{i-\frac{1}{2}}^n = q_i^n - q_{i-1}^n$ and $\Delta q_{i+\frac{1}{2}}^n = q_{i+1}^n - q_i^n$ is TVD if the following sufficient condition is satisfied:*

$$(4.90) \qquad\qquad C_{i+\frac{1}{2}} \geq 0, \quad D_{i+\frac{1}{2}} \geq 0, \quad 0 \leq C_{i+\frac{1}{2}} + D_{i+\frac{1}{2}} \leq 1.$$

PROOF. The numerical solution at grid points $x_i$ and $x_{i+1}$ at the new time $t^{n+1}$ is given by

$$
\begin{aligned}
q_i^{n+1} &= q_i^n - C_{i-\frac{1}{2}}(q_i^n - q_{i-1}^n) + D_{i+\frac{1}{2}}(q_{i+1}^n - q_i^n), \\
q_{i+1}^{n+1} &= q_{i+1}^n - C_{i+\frac{1}{2}}(q_{i+1}^n - q_i^n) + D_{i+\frac{3}{2}}(q_{i+2}^n - q_{i+1}^n),
\end{aligned}
$$

hence
$$(4.91)$$
$$q_{i+1}^{n+1} - q_i^{n+1} = (q_{i+1}^n - q_i^n)(1 - C_{i+\frac{1}{2}} - D_{i+\frac{1}{2}}) + D_{i+\frac{3}{2}}(q_{i+2}^n - q_{i+1}^n) + C_{i-\frac{1}{2}}(q_i^n - q_{i-1}^n).$$

Using the triangle inequality $|a + b| \leq |a| + |b|$ and $|ab| = |a||b|$ we obtain
$$(4.92)$$
$$|q_{i+1}^{n+1} - q_i^{n+1}| \leq |q_{i+1}^n - q_i^n||1 - C_{i+\frac{1}{2}} - D_{i+\frac{1}{2}}| + |D_{i+\frac{3}{2}}||q_{i+2}^n - q_{i+1}^n| + |C_{i-\frac{1}{2}}||q_i^n - q_{i-1}^n|.$$

Using the hypothesis of the theorem (4.90) we have

$$(4.93) \qquad
\begin{aligned}
|q_{i+1}^{n+1} - q_i^{n+1}| \leq |q_{i+1}^n - q_i^n| + C_{i-\frac{1}{2}}|q_i^n - q_{i-1}^n| - C_{i+\frac{1}{2}}|q_{i+1}^n - q_i^n| + \\
D_{i+\frac{3}{2}}|q_{i+2}^n - q_{i+1}^n| - D_{i+\frac{1}{2}}|q_{i+1}^n - q_i^n|.
\end{aligned}
$$

With this intermediate result we have for the total variation at time $t^{n+1}$

$$(4.94) \qquad
\begin{aligned}
\mathrm{TV}(q_h^{n+1}) = \sum_i |q_{i+1}^{n+1} - q_i^{n+1}| \leq \sum_i |q_{i+1}^n - q_i^n| + \\
\overbrace{\sum_i \left( C_{i-\frac{1}{2}}|q_i^n - q_{i-1}^n| - C_{i+\frac{1}{2}}|q_{i+1}^n - q_i^n| \right)}^{=0} + \\
\overbrace{\sum_i \left( D_{i+\frac{3}{2}}|q_{i+2}^n - q_{i+1}^n| - D_{i+\frac{1}{2}}|q_{i+1}^n - q_i^n| \right)}^{=0}.
\end{aligned}
$$

The last two sums vanish due to their telescopic property, hence we obtain

$$(4.95) \qquad\qquad \mathrm{TV}(q_h^{n+1}) \leq \mathrm{TV}(q_h^n),$$

which concludes the proof.

$\square$

**4.2. Higher order spatial reconstruction and the method of Kolgan.** The basic idea consists of changing of the data approximation inside each cell from a piece–wise constant distribution to piece–wise linear data. Since only cell averaged quantities are available, the slopes of the piecewise linear approximation have to be computed appropriately from the cell average of the cell and its neighbors. A sketch of piecewise constant and piecewise linear data is depicted in Fig. 9. Inside each cell we suppose a linear polynomial $P_i^n$ for the state vector $Q$ at time
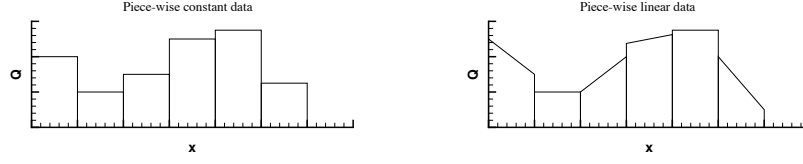
FIGURE 9. Piece–wise constant data (left) and piece-wise linear data (right). Note the discontinuities at the element interfaces.

$t^n$, such that

$$(4.96) \qquad P_i^n(x) = a + b(x - x_i),$$

where the index $i$ denotes the mesh element and $n$ is the time step number. The coefficients $a$ and $b$ are computed by using the principle of conservation. We first use a right–sided reconstruction stencil composed of the element itself and its right neighbor to determine a right reconstruction polynomial $P_i^R$, as follows

$$(4.97) \qquad \frac{1}{\Delta x} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} P_i^R(x)\,dx = Q_i^n, \quad \Rightarrow \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} a + b(x - x_i)\,dx = Q_i^n \Delta x,$$

$$(4.98) \qquad \frac{1}{\Delta x} \int_{x_{i+\frac{1}{2}}}^{x_{i+\frac{3}{2}}} P_i^R(x)\,dx = Q_{i+1}^n, \quad \Rightarrow \int_{x_{i+\frac{1}{2}}}^{x_{i+\frac{3}{2}}} a + b(x - x_i)\,dx = Q_{i+1}^n \Delta x.$$

Integration yields

$$(4.99) \qquad \left[ ax + \frac{b}{2}(x - x_i)^2 \right]_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} = Q_i^n \Delta x, \quad \Rightarrow a = Q_i^n,$$

$$(4.100) \qquad \left[ ax + \frac{b}{2}(x - x_i)^2 \right]_{x_{i+\frac{1}{2}}}^{x_{i+\frac{3}{2}}} = Q_{i+1}^n \Delta x, \quad \Rightarrow b = \frac{Q_{i+1}^n - Q_i^n}{\Delta x} = \frac{\Delta Q_{i+\frac{1}{2}}^n}{\Delta x}.$$

Similarly, imposing the conservation at the left side of the cell interface we compute the coefficients $a$ and $b$

$$(4.101) \qquad \frac{1}{\Delta x} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} P_i^L(x)\,dx = Q_i^n, \quad \Rightarrow a = Q_i^n,$$

$$(4.102) \qquad \frac{1}{\Delta x} \int_{x_{i-\frac{3}{2}}}^{x_{i-\frac{1}{2}}} P_i^L(x)\,dx = Q_{i-1}^n, \quad \Rightarrow b = \frac{Q_i^n - Q_{i-1}^n}{\Delta x} = \frac{\Delta Q_{i-\frac{1}{2}}^n}{\Delta x}.$$

Since we have two spatial reconstruction polynomials available, one from the left and one from the right, we now can make the scheme *nonlinear* by letting the actual choice of which reconstruction polynomial we finally retain depend on the data. One option that satisfies the maximum principle (no generation of new extrema) consists of

$$(4.103) \qquad \Delta Q_i^n = \mathrm{minmod}\left( \Delta Q_{i-\frac{1}{2}}^n, \Delta Q_{i+\frac{1}{2}}^n \right),$$

with $\Delta Q_{i-\frac{1}{2}}^n = Q_i^n - Q_{i-1}^n$ and $\Delta Q_{i+\frac{1}{2}}^n = Q_{i+1}^n - Q_i^n$, where the so–called minmod function is defined by

$$(4.104) \qquad \mathrm{minmod}(a, b) = \begin{cases} 0, & \text{if } ab \leq 0, \\ a, & \text{if } |a| < |b| \\ b, & \text{if } |a| \geq |b|. \end{cases}$$

Then, the linear polynomial for the cell $i$ becomes

$$(4.105) \qquad P_i^n(x) = Q_i^n + \frac{\Delta Q_i^n}{\Delta x}(x - x_i),$$

Once the polynomial is defined, we can compute the states at the cell interfaces $x_{i-\frac{1}{2}}$ and $x_{i+\frac{1}{2}}$

$$(4.106) \qquad Q_{i-\frac{1}{2}}^{n,+} = Q_i^n + \frac{\Delta Q_i^n}{\Delta x}\left(x_i - \frac{\Delta x}{2} - x_i\right) = Q_i^n - \frac{1}{2}\Delta Q_i^n,$$

$$(4.107) \qquad Q_{i+\frac{1}{2}}^{n,-} = Q_i^n + \frac{\Delta Q_i^n}{\Delta x}\left(x_i + \frac{\Delta x}{2} - x_i\right) = Q_i^n + \frac{1}{2}\Delta Q_i^n.$$

The numerical flux is computed including the values at the interfaces from the polynomials $P_i^n(x)$ and $P_{i+1}^n(x)$

$$(4.108) \qquad f_{i+\frac{1}{2}} = f_{i+\frac{1}{2}}\left(Q_{i+\frac{1}{2}}^{n,-}, Q_{i+\frac{1}{2}}^{n,+}\right),$$

with $Q_{i+\frac{1}{2}}^{n,+} = Q_{i+1}^n - \frac{1}{2}\Delta Q_{i+1}^n$ and $Q_{i+\frac{1}{2}}^{n,-} = Q_{i+1}^n + \frac{1}{2}\Delta Q_i^n$.

Kolgan (1972) [**50**] proposed the following scheme

$$(4.109) \qquad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x}\left(f_{i+\frac{1}{2}}\left(Q_{i+\frac{1}{2}}^{n,-}, Q_{i+\frac{1}{2}}^{n,+}\right) - f_{i-\frac{1}{2}}\left(Q_{i-\frac{1}{2}}^{n,-}, Q_{i-\frac{1}{2}}^{n,+}\right)\right).$$

To compute the numerical flux any Riemann solver can be used. This method is of second order of accuracy in space and of first order in time. It is linearly unconditionally unstable, but it is nonlinearly stable for a Courant number $c < \frac{1}{2}$.

## 4.3. Second order in space and in time: the MUSCL method of van Leer (1979).

An significant improvement in accuracy and stability is given by the MUSCL method of van Leer [**87**], who proposed to use a linear polynomial in space *and in time*, as follows

$$(4.110) \qquad P_i^n(x,t) = Q_i^n + \frac{\Delta Q_i^n}{\Delta x}(x - x_i) + \partial_t Q_i(t - t^n),$$

where the term $Q_t$ indicates the derivative in time of the variable $Q$. It can be computed using the PDE and the resulting spatial derivative of the flux is approximated by a central finite difference with respect to the cell center $x_i$ by using the boundary–extrapolated values from *within* the cell:
$$(4.111)$$

$$\partial_t Q_i = \frac{\partial Q}{\partial t} = -\frac{\partial f}{\partial x} \approx -\frac{f\left(Q_{i+\frac{1}{2}}^{n,-}\right) - f\left(Q_{i-\frac{1}{2}}^{n,+}\right)}{\Delta x} = \frac{f\left(Q_{i-\frac{1}{2}}^{n,+}\right) - f\left(Q_{i+\frac{1}{2}}^{n,-}\right)}{\Delta x}.$$

Substituting into (4.110) we have

$$(4.112) \qquad P_i^n(x,t) = Q_i^n + \frac{\Delta Q_i^n}{\Delta x}(x - x_i) + \frac{f\left(Q_{i-\frac{1}{2}}^{n,+}\right) - f\left(Q_{i+\frac{1}{2}}^{n,-}\right)}{\Delta x}(t - t^n).$$

Since the numerical flux is defined by

$$(4.113) \qquad f_{i+\frac{1}{2}} = \frac{1}{\Delta t}\int_{t^n}^{t^{n+1}} f\left(Q(x_{i+\frac{1}{2}}, t)\right) dt,$$

we have to compute the time integral above numerically. To achieve the second order of accuracy we have two possibilities

- The trapezoidal rule
$$(4.114)$$

$$\int_{t^n}^{t^{n+1}} f\left(Q(x_{i+\frac{1}{2}}, t)\right) dt \approx \frac{\Delta t}{2}\left(f\left(Q\left(x_{i+\frac{1}{2}}, t^n\right)\right) + f\left(Q\left(x_{i+\frac{1}{2}}, t^{n+1}\right)\right)\right).$$

- The mid–point rule

$$(4.115) \qquad \int_{t^n}^{t^{n+1}} f\left(Q(x_{i+\frac{1}{2}},t)\right) dt \approx \Delta t f\left(Q(x_{i+\frac{1}{2}}, t^{n+\frac{1}{2}})\right) = \Delta t f\left(Q_{i+\frac{1}{2}}^{n+\frac{1}{2}}\right),$$

Note that the mid–point integration is cheaper than the trapezoidal rule, hence it is used in the following.

Finally, the MUSCL (monotone upwind scheme for conservation laws) method of van Leer is given by

$$(4.116)\quad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x}\left(f_{i+\frac{1}{2}}\left(Q_{i+\frac{1}{2}}^{n+\frac{1}{2},-}, Q_{i+\frac{1}{2}}^{n+\frac{1}{2},+}\right) - f_{i-\frac{1}{2}}\left(Q_{i-\frac{1}{2}}^{n+\frac{1}{2},-}, Q_{i-\frac{1}{2}}^{n+\frac{1}{2},+}\right)\right).$$

It is second order accurate in space and time and is linearly stable up to Courant number one.

## 5. The MUSCL method on two–dimensional Cartesian grids

We consider the hyperbolic PDE

$$(4.117) \qquad \frac{\partial Q}{\partial t} + \frac{\partial f}{\partial x} + \frac{\partial g}{\partial y} = 0.$$

A finite volume scheme based on rectangular spatial control volumes $T_{ij} = [x_{i-\frac{1}{2}}; x_{i+\frac{1}{2}}] \times [y_{i-\frac{1}{2}}; y_{i+\frac{1}{2}}]$ becomes

$$(4.118) \qquad Q_{ij}^{n+1} = Q_{ij}^n - \frac{\Delta t}{\Delta x}\left(f_{i+\frac{1}{2},j} - f_{i-\frac{1}{2},j}\right) - \frac{\Delta t}{\Delta y}\left(g_{i,j+\frac{1}{2}} - g_{i,j-\frac{1}{2}}\right).$$

In order to compute the numerical fluxes across the element interfaces using the second order MUSCL scheme one proceeds as in one space dimension, obtaining:

(1) TVD reconstruction of the gradients (slopes) based on the solution at time $t^n$:

$$(4.119) \qquad \begin{aligned} \Delta Q_{ij}^x &= \text{minmod}(Q_{i+1,j}^n - Q_{i,j}^n, Q_{i,j}^n - Q_{i-1,j}^n), \\ \Delta Q_{ij}^y &= \text{minmod}(Q_{i,j+1}^n - Q_{i,j}^n, Q_{i,j}^n - Q_{i,j-1}^n). \end{aligned}$$

(2) The space–time polynomial in each cell $T_{ij}$ which is valid for the current time step is written in terms of a space–time Taylor series expanded about $x_i$, $y_i$ and $t^n$ as

$$(4.120)\quad Q_{ij}(x,y,t) = Q_{ij}^n + \Delta Q_{ij}^x/\Delta x(x - x_i) + \Delta Q_{ij}^y/\Delta y(y - y_i) + \partial_t Q_{ij}(t - t^n).$$

(3) The time derivative in the Taylor series is computed based on the PDE as

$$(4.121) \qquad \partial_t Q_{ij} = -f_x - g_y.$$

We now use a central finite difference with respect to the element barycenter and based on the reconstructed values *from within* the element $T_{ij}$ to approximate the spatial derivatives of the fluxes. This yields

$$(4.122) \qquad \partial_t Q_{ij} \approx -\frac{f(Q_{i+\frac{1}{2},j}^{n,-}) - f(Q_{i-\frac{1}{2},j}^{n,+})}{\Delta x} - \frac{g(Q_{i,j+\frac{1}{2}}^{n,-}) - g(Q_{i,j-\frac{1}{2}}^{n,+})}{\Delta y},$$

with the boundary extrapolated values at the old time $t^{n+1}$

$$(4.123) \qquad Q_{i\pm\frac{1}{2},j}^{n,\mp} = Q_{ij}^n \pm \frac{1}{2}\Delta Q_{ij}^x, \qquad Q_{i,j\pm\frac{1}{2}}^{n,\mp} = Q_{ij}^n \pm \frac{1}{2}\Delta Q_{ij}^y.$$

(4) Compute the boundary extrapolated values at the half time level $t^{n+\frac{1}{2}}$ as

$$(4.124)\quad Q_{i\pm\frac{1}{2},j}^{n+\frac{1}{2},\mp} = Q_{ij}^n \pm \frac{1}{2}\Delta Q_{ij}^x + \frac{1}{2}\partial_t Q_{ij}, \qquad Q_{i,j\pm\frac{1}{2}}^{n+\frac{1}{2},\mp} = Q_{ij}^n \pm \frac{1}{2}\Delta Q_{ij}^y + \frac{1}{2}\partial_t Q_{ij}.$$

(5) Insert these values into the numerical fluxes of the finite volume scheme (4.118) and update the solution to the new time $t^{n+1}$.

## 6. Multidimensional finite volume schemes on general meshes

The sketch of the general multidimensional case of the finite volume method for arbitrary unstructured meshes is given in the following. We consider the multi–dimensional nonlinear system of hyperbolic conservation laws

$$(4.125) \qquad \frac{\partial Q}{\partial t} + \frac{\partial f}{\partial x} + \frac{\partial g}{\partial y} + \frac{\partial h}{\partial z} = 0,$$

where $Q, f, g, h \in \mathbb{R}^m$ and $\vec{x} = (x, y, z) \in \mathbb{R}^3$. The PDE (4.125) in compact form results

$$(4.126) \qquad \frac{\partial Q}{\partial t} + \operatorname{div} F = 0, \quad \text{or} \quad \frac{\partial Q}{\partial t} + \nabla \cdot F = 0,$$

where $F$ denotes the flux tensor $F = (f, g, h)$. The PDE integrated over a general space-time control volume $T_i \times \left[t^n, t^{n+1}\right]$ gives the integral form of (4.126)

$$(4.127) \qquad \int_{T_i} \int_{t^n}^{t^{n+1}} \frac{\partial Q}{\partial t} dt dV + \int_{T_i} \int_{t^n}^{t^{n+1}} \nabla \cdot F dt dV = 0.$$

The volume integral is converted to a surface integral applying the Gauss–Ostrogradski theorem

$$(4.128) \qquad \int_{T_i} \left( Q \left( \vec{x}, t^{n+1} \right) - Q \left( \vec{x}, t^n \right) \right) dV + \int_{t^n}^{t^{n+1}} \int_{\partial T_i} F \cdot \vec{n} \, dt dS = 0,$$

where $\partial T_i$ is the boundary of the $i^{th}$ control volume and $\vec{n}_j$ denotes the outward pointing unit normal vector to the surface $\partial T_i$. We introduce the following definitions
(4.129)

$$Q_i^n = \frac{1}{|V_i|} \int_{T_i} Q \left( \vec{x}, t^n \right) dV, \quad F_{i,j+\frac{1}{2}} \cdot \vec{n}_j = \frac{1}{\Delta t} \frac{1}{|\partial T_i|} \int_{t^n}^{t^{n+1}} \int_{\partial T_i} F \cdot \vec{n}_j dt dV = 0,$$

where the index $j$ indicates the $j^{th}$ side of the control volume $T_i$, which is supposed to have a polygonal shape, i.e. the boundary is composed of piecewise linear faces, as plotted in Fig. 10, and $J$ is the number of faces. We obtain the following *exact* relation

$$(4.130) \qquad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{|T_i|} \sum_{j=1}^{J} |\partial T_j| \, F_{i,j+\frac{1}{2}} \cdot \vec{n}_j.$$

Inserting any classical numerical flux into the term $F_{i,j+\frac{1}{2}} \cdot \vec{n}_j$ will then give a finite volume scheme on general polygonal control volumes.

## 7. Rotational invariance

Newtonian mechanics is invariant with respect to a rotation of the coordinate system. As a consequence we expect that also the governing PDE of Newtonian mechanics are rotationally invariant. Let the angle of rotation be $\alpha$, then the vector which defines the first axis of the new coordinate system is defined as

$$(4.131) \qquad \vec{n} = \left( \begin{array}{c} \cos \alpha \\ \sin \alpha \end{array} \right).$$

FIGURE 10. A polygonal element of volume $V_i$ and total surface $\partial T_i$.



FIGURE 11. Rotation of the coordinate system by an angle $\alpha$.

The vector $\vec{x}$ of the old coordinates $x$ and $y$ is given in terms of the vector $\vec{x}'$ of the new coordinates $x'$ and $y'$ as

$$(4.132) \qquad \vec{x} = M\vec{x}', \qquad \text{with } M = \begin{pmatrix} \cos\alpha & -\sin\alpha \\ \sin\alpha & \cos\alpha \end{pmatrix}.$$

For a sketch of the coordinate rotation see 11. The same transformation applies to all vector quantities, such as velocity $\vec{v} = M\vec{v}'$, while the scalars (density, pressure, energy) do *not* change with the transformation. Let $T$ be the transformation matrix which expresses the coordinate transformation of the vector of conserved quantities. For scalar quantities, $T$ contains a row of zeros with a value of one on the diagonal, whereas for vector quantities, $T$ contains a double row of zeros with the transformation matrix $M$ on the corresponding diagonal block. We then have

$$(4.133) \qquad Q = TQ', \quad \text{and} \quad Q' = T^{-1}Q = T^T Q.$$

DEFINITION 2. *A system of PDE of the form*

$$(4.134) \qquad \frac{\partial Q}{\partial t} + \nabla \cdot F(Q) = 0, \quad \text{with } F(Q) = (f, g)$$

is said to be rotationally invariant if

$$(4.135) \qquad F(Q) \cdot \vec{n} = Tf(T^{-1}Q).$$

EXERCISE 5. *Prove that the Euler equations of compressible gasdynamics and the shallow water equations are rotationally invariant.*

## 8. The wave–propagation form of finite volume schemes

This particular formulation of conservative finite volume schemes is due to LeVeque [**55**] and gives an interesting alternative interpretation of the Riemann problems at the element interfaces. Given an explicit conservative method

$$(4.136) \qquad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x} \left( f_{i+\frac{1}{2}} - f_{i-\frac{1}{2}} \right),$$

with the Roe flux defined as

$$(4.137) \qquad f_{i+\frac{1}{2}} = \frac{1}{2} \left( f(Q_{i+1}^n) + f(Q_i^n) \right) - \frac{1}{2} |\tilde{A}_{i+\frac{1}{2}}| \left( Q_{i+1}^n - Q_i^n \right),$$

we have by adding and subtracting $f_i = f(Q_i^n)$

$$(4.138) \qquad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x} \left( f_{i+\frac{1}{2}} - f_i - (f_{i-\frac{1}{2}} - f_i) \right).$$

The right flux difference can be written as

$$(4.139) \qquad f_{i+\frac{1}{2}} - f_i = \frac{1}{2} \left( f(Q_{i+1}^n) - f(Q_i^n) \right) - \frac{1}{2} |\tilde{A}_{i+\frac{1}{2}}| \left( Q_{i+1}^n - Q_i^n \right)$$

and using the property of the Roe matrix we have

$$(4.140) \qquad f_{i+\frac{1}{2}} - f_i = \frac{1}{2} \tilde{A}_{i+\frac{1}{2}} \left( Q_{i+1}^n - Q_i^n \right) - \frac{1}{2} |\tilde{A}_{i+\frac{1}{2}}| \left( Q_{i+1}^n - Q_i^n \right),$$

or simply
$$(4.141)$$
$$f_{i+\frac{1}{2}} - f_i = \frac{1}{2} \left( \tilde{A}_{i+\frac{1}{2}} - |\tilde{A}_{i+\frac{1}{2}}| \right) \left( Q_{i+1}^n - Q_i^n \right) = \tilde{A}_{i+\frac{1}{2}}^- \left( Q_{i+1}^n - Q_i^n \right) := D_{i+\frac{1}{2}}^-.$$

A similar manipulation can be also performed on the left hand side, obtaining:
$$(4.142)$$
$$f_{i-\frac{1}{2}} - f_i = -\frac{1}{2} \left( \tilde{A}_{i-\frac{1}{2}} + |\tilde{A}_{i-\frac{1}{2}}| \right) \left( Q_i^n - Q_{i-1}^n \right) = \tilde{A}_{i-\frac{1}{2}}^+ \left( Q_i^n - Q_{i-1}^n \right) := D_{i-\frac{1}{2}}^+.$$

Inserting in the finite volume method we obtain the wave propagation form of LeVeque

$$(4.143) \qquad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x} \left( D_{i+\frac{1}{2}}^- + D_{i-\frac{1}{2}}^+ \right),$$

or with $\Delta Q_{i-\frac{1}{2}} = Q_i^n - Q_{i-1}^n$ and $\Delta Q_{i+\frac{1}{2}} = Q_{i+1}^n - Q_i^n$ we obtain the new cell average at time $t^{n+1}$ in terms of the right–moving waves that enter from the left and the left–moving waves that enter from the right, both expressed as a product of the positive or negative part of the Roe–matrix times the jump of the vector of conserved quantities at the interface:

$$(4.144) \qquad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x} \left( \tilde{A}_{i+\frac{1}{2}}^- \Delta Q_{i+\frac{1}{2}}^n + \tilde{A}_{i-\frac{1}{2}}^+ \Delta Q_{i-\frac{1}{2}}^n \right).$$

A graphical interpretation of this scheme is sketched in Fig. 12. Remember that in the method of Roe, the continuous rarefaction waves are replaced by rarefaction shocks. In fact, also the original method of Godunov [**37**] was expressed in the form illustrated by Fig. 12 and *not* in the usual conservative flux divergence form.
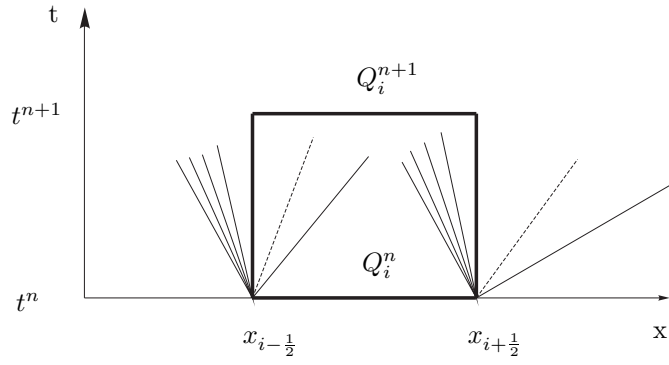
FIGURE 12. Graphical interpretation of the wave–propagation algorithm of LeVeque. The cell–average $Q_i$ is only modified by waves that *enter* the control volume from the boundaries.

**Part 2**

# Advanced Algorithms for the Discretization of Hyperbolic PDE

# Path–conservative finite volume schemes

In the previous part of this manuscript we have dealt with *conservative* hyperbolic systems. Actually for many physically relevant processes, it is possible to write the problem in conservative form, see for example the Euler or Navier–Stokes equations of compressible gas dynamics, the equations of magneto–hydrodynamics (MHD) or the equations governing nonlinear elasticity. However, there are also many physically relevant systems that can *not* be written in conservation form, in particular problems arising in multi–phase and multi–fluid flow physics. We therefore have to enrich our set of numerical techniques with schemes that are also able to handle more general hyperbolic systems written in non–conservative form

$$(5.1) \qquad \frac{\partial Q}{\partial t} + A(Q) \frac{\partial Q}{\partial x} = 0, \qquad A \in \mathbb{R}^{m \times m}, Q \in \mathbb{R}^m, x \in \mathbb{R}, t \in \mathbb{R}_0^+.$$

## 1. Some basic aspects of the theory of non–conservative systems

For classical (strong) solutions where $Q$ is differentiable there are no particular problems with PDE of the kind (5.1), since the derivatives are defined everywhere and the method of characteristics can be applied. The main problem arises with non–conservative PDE of the type (5.1) in the presence of shock waves (discontinuities), where the differential form of the equation does not hold any more. In the conservative case, the solution to the problem was given by the Rankine–Hugoniot conditions, that were derived from an integral form of the conservation law. For non–conservative systems, the theory of shock waves is much more complex and far less developed than for the conservative case. The main contribution to the theory of shock waves for non–conservative PDE is due to an article by Dal Maso, Le Floch and Murat [**58**], and whenever we refer to this theory in the following we will call it the DLM theory. According to the DLM theory, in the case of a non–conservative PDE, the jump relations across a discontinuity depend on the chosen *integration path* that is used to connect the left and the right state at the discontinuity in phase–space. Different integration paths result in different jump relations.

DEFINITION 3. *A path $\Psi$ in phase–space $\Omega \subset \mathbb{R}^m$ is a Lipschitz continuous function of a parameter $0 \leq s \leq 1$*

$$(5.2) \qquad \qquad \Psi : [0;1] \times \Omega \times \Omega \to \Omega,$$

*that must satisfy the following properties:*

$$(5.3) \qquad \Psi(0, Q_L, Q_R) = Q_L, \qquad \Psi(1, Q_L, Q_R) = Q_R, \qquad \forall Q_L, Q_R \in \Omega,$$

$$(5.4) \qquad |\frac{\partial \Psi}{\partial s}(s, Q_L, Q_R)| \leq k|Q_R - Q_L|, \qquad \forall Q_L, Q_R \in \mathcal{B} \subset \Omega,$$

$$(5.5) \qquad |\frac{\partial \Psi}{\partial s}(s, Q_L, Q_R) - \frac{\partial \Psi}{\partial s}(s, Q_l, Q_r)| \leq K(|Q_L - Q_l| + |Q_R - Q_r|),$$
$$\forall Q_L, Q_R, Q_l, Q_r \in \mathcal{B} \subset \Omega,$$

with $\mathcal{B} \subset \Omega$ being an arbitrary bounded subset of the phase–space $\Omega$.

EXAMPLE 32. *The simplest possible integration path is obviously the straight line segment, i.e.*

$$(5.6) \qquad \Psi = \Psi(s, Q_L, Q_R) = Q_L + s(Q_R - Q_L), \qquad 0 \le s \le 1.$$

EXERCISE 6. *Prove that the segment path (5.6) satisfies all the three properties given in (5.3)-(5.5).*

According to the DLM theory [58] weak solutions of (5.1) can be interpreted as a Borel measure and the classical Rankine–Hugoniot jump–relations are generalized to

$$(5.7) \qquad \sigma(Q_R - Q_L) = \int_0^1 A(\Psi(s, Q_L, Q_R)) \frac{\partial \Psi}{\partial s} ds = 0,$$

where $\sigma$ is the speed of the discontinuity between the two states $Q_L$ and $Q_R$. According to the DLM theory, a piecewise $\mathcal{C}^1$ function $Q$ is a *weak solution* of (5.1) if and only if it is a *classical* solution where it is $\mathcal{C}^1$ and if it satisfies (5.7) along the discontinuities.

It is easy to prove that in the case when $A(Q)$ is the Jacobian of a flux $f(Q)$, i.e. $A(Q) = \partial f(Q)/\partial Q$, the extended Rankine–Hugoniot conditions reduce to the classical ones, *independent* of the particular choice of the path:

$$(5.8)$$

$$\sigma(Q_R - Q_L) = \int_0^1 A(\Psi) \frac{\partial \Psi}{\partial s} ds = \int_0^1 \frac{\partial f(\Psi)}{\partial \Psi} \frac{\partial \Psi}{\partial s} ds = \int_{f(Q_L)}^{f(Q_R)} df = f(Q_R) - f(Q_R).$$

## 2. Path–conservative finite volume schemes

The first special case of a path–conservative scheme has been proposed by Toumi in [86], where he proposed a weak formulation of the method of Roe based on the integration along a particular path. The framework of path–conservative schemes has then been successively extended and analyzed by Parés [63] and Castro et al. [9, 10] in the finite volume context and by Rhebergen et al. [66] in the discontinuous Galerkin finite element context. The first better than second order accurate path–conservative schemes on general unstructured meshes in two and three space dimensions have been published in [21, 23, 28].

A general path–conservative finite volume scheme of first order of accuracy is written in the following way, which is similar to the wave–propagation formulation of LeVeque [55]:

$$(5.9) \qquad Q_i^{n+1} = Q_i^n - \frac{\Delta t}{\Delta x} \left( D_{i+\frac{1}{2}}^- + D_{i-\frac{1}{2}}^+ \right),$$

with the jump terms on the boundaries

$$(5.10) \qquad D_{i+\frac{1}{2}}^- = C_{i+\frac{1}{2}}^-(Q_{i+1}^n - Q_i^n), \qquad \text{and} \qquad D_{i-\frac{1}{2}}^+ = C_{i-\frac{1}{2}}^+(Q_i^n - Q_{i-1}^n).$$

The actual choice of the numerical scheme lies in the particular choice of the matrices $C_{i+\frac{1}{2}}^- = C_{i+\frac{1}{2}}^-(Q_i^n, Q_{i+1}^n)$ and $C_{i-\frac{1}{2}}^+ = C_{i-\frac{1}{2}}^+(Q_{i-1}^n, Q_i^n)$, which are both functions of the left and the right state at the cell interface. For any path–conservative scheme, the matrices $C_{i+\frac{1}{2}}^{\pm}$ must have the following property at the interface:

$$(5.11) \qquad C_{i+\frac{1}{2}}^- + C_{i+\frac{1}{2}}^+ = \tilde{A}_{i+\frac{1}{2}},$$

where $\tilde{A}_{i+\frac{1}{2}} = \tilde{A}(Q_i^n, Q_{i+1}^n)$ is the generalized Roe matrix with the following three features:

First, $\tilde{A}_{i+\frac{1}{2}}$ has to satisfy the generalized Rankine–Hugoniot conditions:

$$(5.12) \qquad \tilde{A}(Q_L, Q_R)(Q_R - Q_L) = \int\limits_0^1 A(\Psi(s, Q_L, Q_R)) \frac{\partial \Psi}{\partial s} ds.$$

From (5.12) it is obvious that for a conservative system, where $A(Q)$ is the Jacobian of a flux $f(Q)$, the generalized Roe property (5.12) reduces to the classical one

$$(5.13) \qquad \tilde{A}(Q_L, Q_R)(Q_R - Q_L) = f(Q_R) - f(Q_L).$$

Second, the generalized Roe matrix $\tilde{A}(Q_L, Q_R)$ must be hyperbolic, i.e. it must have real eigenvalues and a full set of linearly independent eigenvectors. Third, the generalized Roe matrix must satisfy the consistency condition

$$(5.14) \qquad \tilde{A}(Q, Q) = A(Q).$$

With these properties, it is then easy to prove the following

THEOREM 4. *A path–conservative finite volume scheme of the form* (5.9)-(5.12) *is exactly conservative in the case where the system matrix* $A(Q)$ *is the Jacobian of a flux* $f(Q)$, *i.e. when*

$$(5.15) \qquad A(Q) = \frac{\partial f(Q)}{\partial Q}.$$

PROOF. Assuming no contributions from the boundaries for simplicity, the conserved quantity $Q$ at the new time $t^{n+1}$ is

(5.16)
$$\sum_i \Delta x Q_i^{n+1} = \sum_i \Delta x Q_i^n - \Delta t \sum_i \left( C_{i+\frac{1}{2}}^-(Q_{i+1}^n - Q_i^n) + C_{i-\frac{1}{2}}^+(Q_i^n - Q_{i-1}^n) \right) =$$
$$\sum_i \Delta x Q_i^n - \Delta t \sum_i \left( C_{i-\frac{1}{2}}^- + C_{i-\frac{1}{2}}^+ \right) (Q_i^n - Q_{i-1}^n) =$$
$$\sum_i \Delta x Q_i^n - \Delta t \sum_i \tilde{A}_{i-\frac{1}{2}}(Q_i^n - Q_{i-1}^n) =$$
$$\sum_i \Delta x Q_i^n - \Delta t \sum_i \left( f(Q_i^n) - f(Q_{i-1}^n) \right) = \sum_i \Delta x Q_i^n,$$

since the generalized Rankine–Hugoniot conditions reduce to the classical ones (5.13) in this case and since the sum vanishes due to its telescopic property.   □

A particularly interesting general family of path–conservative schemes can be written in the following form:

$$(5.17) \qquad D_{i+\frac{1}{2}}^\pm = \frac{1}{2} \left( \tilde{A}_{i+\frac{1}{2}} \pm \Theta_{i+\frac{1}{2}} \right) \left( Q_{i+1}^n - Q_i^n \right),$$

where $\tilde{A}_{i+\frac{1}{2}}$ is the generalized Roe matrix and $\Theta_{i+\frac{1}{2}} > 0$ is the dissipation matrix of the scheme.

EXERCISE 7. *Prove that the family of path–conservative schemes given by expression* (5.17) *can be written in the classical flux–conservative form with numerical fluxes*

$$(5.18) \qquad f_{i+\frac{1}{2}} = D_{i+\frac{1}{2}}^- + f(Q_i^n), \qquad and \qquad f_{i-\frac{1}{2}} = f(Q_i^n) - D_{i-\frac{1}{2}}^+,$$

*for the conservative case* $A(Q) = \partial f(Q)/\partial Q$.

With different choices of $\Theta_{i+\frac{1}{2}}$ we can reproduce several classical schemes:

- The Lax–Friedrichs method.

$$(5.19) \qquad \Theta^{LF}_{i+\frac{1}{2}} = \frac{\Delta x}{\Delta t}\mathbf{I}.$$

- A Lax–Wendroff–type method.

$$(5.20) \qquad \Theta^{LW}_{i+\frac{1}{2}} = \frac{\Delta t}{\Delta x}\tilde{A}^2_{i+\frac{1}{2}}.$$

- The FORCE scheme [**80, 81**] and its extensions to non–conservative hyperbolic systems [**23, 8, 7**].

$$(5.21) \qquad \Theta^{FO}_{i+\frac{1}{2}} = \frac{1}{2}\left(\Theta^{LF}_{i+\frac{1}{2}} + \Theta^{LW}_{i+\frac{1}{2}}\right) = \frac{1}{2}\left(\frac{\Delta x}{\Delta t}\mathbf{I} + \frac{\Delta t}{\Delta x}\tilde{A}^2_{i+\frac{1}{2}}\right).$$

- The method of Rusanov [**70**].

$$(5.22) \qquad \Theta_{i+\frac{1}{2}} = \max(|\Lambda(Q^n_i)|, |\Lambda(Q^n_{i+1})|)\mathbf{I}.$$

- The method of Roe [**69**].

$$(5.23) \qquad \Theta_{i+\frac{1}{2}} = \left|\tilde{A}_{i+\frac{1}{2}}\right|.$$

For the segment path (5.6), this is identical with

$$(5.24) \qquad \Theta_{i+\frac{1}{2}} = \left|\tilde{A}_{i+\frac{1}{2}}\right| = \left|\int_0^1 A(\Psi(Q^n_i, Q^n_{i+1}, s))ds\right|.$$

The proof is left to the reader as exercise.

- It has been shown in [**28**] that for the case of a segment path (5.6) the method of Osher and Solomon [**62**] takes the form

$$(5.25) \qquad \Theta_{i+\frac{1}{2}} = \int_0^1 \left|A(\Psi(Q^n_i, Q^n_{i+1}, s))\right|ds,$$

which is very similar to the method of Roe, apart that the integration and the matrix absolute value operator are exchanged.

In the expressions above, $\mathbf{I}$ denotes the identity matrix.

## 3. Application to the shallow water equations with variable bottom

The shallow water equations with spatially variable but temporally fixed bottom topography read as

$$(5.26) \qquad \begin{aligned} \frac{\partial h}{\partial t} + \frac{\partial hu}{\partial x} &= 0, \\ \frac{\partial hu}{\partial t} + \frac{\partial}{\partial x}\left(hu^2 + \frac{1}{2}gh^2\right) &= -gh\frac{\partial b}{\partial x}, \\ \frac{\partial b}{\partial t} &= 0. \end{aligned}$$

The right hand side due to the spatially variable bottom can not be cast into conservation form, but the system (5.27) can be written easily under the general non–conservative form (5.1) with

$$(5.27) \qquad A(Q) = \begin{pmatrix} 0 & 1 & 0 \\ c^2 - u^2 & 2u & c^2 \\ 0 & 0 & 0 \end{pmatrix},$$

with the wave speed $c^2 = gh$.

**3.1. Eigenstructure.** The eigenvalue matrix is $\Lambda = \mathrm{diag}(u - c, 0, u + c)$ and the corresponding right eigenvector matrix and its inverse are

$$(5.28) \quad R = \begin{pmatrix} 1 & 1 & 1 \\ u - c & 0 & u + c \\ 0 & \frac{u^2}{c^2} - 1 & 0 \end{pmatrix}, \qquad R^{-1} = \frac{1}{2} \begin{pmatrix} \frac{c+u}{c} & -\frac{1}{c} & \frac{c}{c-u} \\ 0 & 0 & \frac{c^2}{u^2 - c^2} \\ \frac{c-u}{c} & \frac{1}{c} & \frac{c}{c+u} \end{pmatrix}.$$

**3.2. Well–balancing.** A very simple but physically very important steady state solution of (5.27) is the so–called still water or lake–at–rest solution, [**54**]:

$$(5.29) \qquad\qquad h + b = const. \qquad u = 0.$$

It simply states that if the free surface $\eta = h + b$ is constant and the velocity is zero, the solution must remain constant in time. Although this seems quite obvious from a physical point of view, a lot of research has been dedicated to this topic in the last decades in order to design numerical methods which satisfy the so-called C–property, see [**38, 54, 91, 88, 59, 60, 34, 3, 9, 63**] for more details. Castro et al. [**9**] have proven that path–conservative Roe–type schemes are *exactly* well–balanced for system (5.27) if the straight–line segment path is used. Also Osher–type path–conservative schemes [**28**] are exactly well–balanced for (5.27), and the other schemes of the family (5.17) can be made exactly well–balanced when using segment integration paths and an appropriate modification of the matrix $\Theta_{i+\frac{1}{2}}$, see [**8**].

EXERCISE 8. *Write a MATLAB program that implements the family of path–conservative methods* (5.17) *for the shallow water equations with variable bottom* (5.27). *For the Osher and the Roe method, use a numerical integration of the integrals along the segment path. Try different quadrature formulae, such as the trapezoidal rule, the Simpson rule or a three–point Gaussian quadrature rule with quadrature points* $s_1 = 1/2 - \sqrt{15}/10$, $s_2 = 1/2$, $s_3 = 1/2 + \sqrt{15}/10$ *and associated weights* $\omega_1 = 5/18$, $\omega_2 = 8/18$, $\omega_3 = 5/18$. *Solve the Riemann problem*

$$(5.30) \qquad\qquad Q(x, 0) = \begin{cases} Q_L, & \text{if} \quad x < 0, \\ Q_R, & \text{if} \quad x > 0, \end{cases}$$

*with the following initial conditions:*

- $Q_L = (2, 0, 0)^T$, *and* $Q_R = (1, 0, 1)^T$,
- $Q_L = (1, 0, 0)^T$, *and* $Q_R = (1e - 7, 0, 0)^T$,
- $Q_L = (1.46184, 0, 0)^T$, *and* $Q_R = (0.30873, 0, 0.2)^T$.

# Higher order WENO schemes

In the first part of this manuscript we have studied (and proven) the theorem of Godunov, which is a fundamental theorem for the construction of numerical methods for hyperbolic PDE. According to the Godunov theorem, it is not possible to devise any linear numerical scheme that is better than first order accurate and monotone. The only way to circumvent the theorem is the design of *nonlinear* scheme, since in this case the hypothesis of the theorem do not apply. The design of second order accurate and monotone schemes was possible by the use of *slope limiters*, which are nonlinear as they depend on the actual numerical solution of the problem. In this chapter, we study one possibility to construct even higher order accurate numerical methods, which can be achieved using nonlinear ENO [40] or WENO [48] interpolation.

## 1. Pointwise WENO reconstruction

In a finite volume scheme, we need to compute fluxes across the element interfaces. For this purpose, numerical flux functions are used, as presented in the first part of the manuscript. These numerical flux functions need two point values of the numerical solution at the element interface, one extrapolated to the interface from the left and another value extrapolated to the interface from the right. As a consequence, the WENO method of Jiang and Shu [48] produces a high order accurate point–wise reconstruction of the solution at the element interfaces $x_{i\pm\frac{1}{2}}$. The general ideas of the WENO scheme [48] are the following.

In order to obtain a $k$–th order accurate WENO method, called WENO $k$ in the following, we need piecewise reconstruction polynomials of degree $M = k-1$ for each cell $T_i = [x_{i-\frac{1}{2}}; x_{i+\frac{1}{2}}]$. To calculate the unknown coefficients of the reconstruction polynomials from the known cell averages $Q_j^n$ one needs a *reconstruction stencil* or *stencil*

$$(6.1) \qquad \mathcal{S}_i^M = \bigcup_{j=i-e}^{i+e} T_j,$$

composed of $k = 2e + 1$ elements, where $e$ is the extension of the stencil to the left and the right. The reconstruction stencil must always include the element $T_i$ itself. The resulting reconstruction polynomial has $k$ coefficients and is of degree $M = k - 1$. According to the relative position of the stencil elements with respect to the cell $T_i$ for which we want to do the reconstruction, a stencil is called *centered*, *left–sided* or *right–sided*. The reconstruction polynomial $P_i^M(x, t^n)$ of degree $M$ is obtained from the known cell averages $Q_j^n$ by requiring *integral conservation*, i.e. we require that

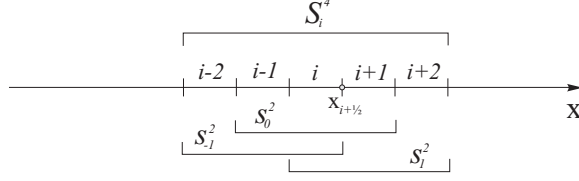$$(6.2) \qquad \frac{1}{\Delta x} \int_{T_j} P_i^M(x, t^n) dx = Q_j^n, \qquad \forall T_j \in \mathcal{S}_i^M.$$

FIGURE 1. Sketch of the WENO5 reconstruction procedure for interface $x_{i+\frac{1}{2}}$. The big stencil $S_4$ needed for the reconstruction of a fourth degree polynomial is divided into three smaller sub–stencils, on each of which a piecewise second degree polynomial is reconstructed. The point value $Q_{i+\frac{1}{2}}^n$ at the interface $x_{i+\frac{1}{2}}$ is then given by a suitable non–linear combination of the reconstruction polynomials obtained on the sub–stencils.

For the WENO method of order $k$ in one space dimension, we need *one big central* reconstruction stencil $S_i^M$ of $k = M + 1$ elements and $M/2 + 1$ small sub–stencils $s_m^{M/2}$ composed of $M/2 + 1$ elements to reconstruct several lower order polynomials of degree $M/2$. Here $m$ is the stencil–shift with respect to the central element $T_i$. The *linear* WENO reconstruction at the element interface $x_{i+\frac{1}{2}}$ is then given as a *linear combination* of the lower order reconstruction polynomials obtained from the substencils $s_r^{M/2}$ using the same integral conservation principle of (6.2), where the linear weights $\lambda$ are chosen in such a way that the linear combination of the lower order polynomials is *identical* with the one obtained via the reconstruction polynomial on the big stencil $S_i^M$. The weights $\lambda$ obviously depend on the position $x$ for which the reconstruction is to be done and for consistency reasons the weights must always sum to one. Furthermore, the weights $\lambda$ should be positive and must not depend on the solution $Q_j^n$.

EXAMPLE 33. *The third order accurate WENO3 method uses one big central stencil composed of the three elements $S_i^2 = \{T_{i-1}, T_i, T_{i+1}\}$ and two one–sided substencils $s_{-1}^1 = \{T_{i-1}, T_i\}$ and $s_1^1 = \{T_i, T_{i+1}\}$. Using the integral conservation principle, we obtain the second order reconstruction polynomial on the big stencil $S_i^2$ as*

(6.3)
$$P_i^2(x) = \frac{1}{2} Q_{i-1}^n + \frac{5}{6} Q_i^n - \frac{1}{6} Q_{i+1}^n + \left( Q_i^n - Q_{i-1}^n \right) \xi + \left( \frac{1}{2} Q_{i-1}^n - Q_i^n + \frac{1}{2} Q_{i+1}^n \right) \xi^2,$$

*with $x = x_{i-\frac{1}{2}} + \xi \Delta x$. For the two first order polynomials on the sub–stencils we obtain*

(6.4)
$$p_{-1}^1(x) = \frac{1}{2} Q_i^n + \frac{1}{2} Q_{i-1}^n + \left( Q_i^n - Q_{i-1}^n \right) \xi,$$

*and*

(6.5)
$$p_1^1(x) = \frac{3}{2} Q_i^n - \frac{1}{2} Q_{i+1}^n + \left( Q_{i+1}^n - Q_i^n \right) \xi.$$

*The two conditions for obtaining the linear weights $\lambda_{-1}$ and $\lambda_1$ are then*

(6.6)
$$\lambda_{-1} + \lambda_1 = 1,$$

*and*

(6.7)
$$\lambda_{-1} p_{-1}^1(x_{i+\frac{1}{2}}) + \lambda_1 p_1^1(x_{i+\frac{1}{2}}) = P_i^2(x_{i+\frac{1}{2}}).$$

*The resulting linear weights are $\lambda_{-1} = 1/3$ and $\lambda_{+1} = 2/3$.*

EXERCISE 9. *Compute the linear weights $\lambda$ for the WENO5 method when performing the reconstruction for the point $x_{i+\frac{1}{2}}$. This exercise can be either solved by hand or more conveniently at the aid of a computer algebra system such as MAPLE.*

In order to make the WENO scheme *nonlinear*, i.e. data–dependent, the reconstruction at point $x_{i+\frac{1}{2}}$ is obtained by using a *nonlinear* combination of the lower order reconstruction polynomials of the sub–stencils by substituting the linear weights with *nonlinear weights* $\omega$, which are defined as

$$(6.8) \qquad \omega_s = \frac{\tilde{\omega}_s}{\sum\limits_{s} \tilde{\omega}_s}, \qquad \tilde{\omega}_s = \frac{\lambda_s}{(\sigma_s + \epsilon)^r}.$$

Here, $\sigma_s$ denotes the so–called smoothness or oscillation indicator defined later, $\epsilon$ is a small number to avoid division by zero and $r$ is an exponent for which Jiang and Shu always choose $r = 2$. For the smoothness indicator Jiang and Shu propose

$$(6.9) \qquad \sigma_s = \sum_{l=1}^{M/2} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} h^{2l-1} \left( \frac{\partial^l}{\partial x^l} p_s^{M/2}(x) \right)^2 dx,$$

where the term $h^{2m-1}$ is used to remove scaling effects from the derivatives.

## 2. Polynomial WENO reconstruction

The original pointwise WENO reconstruction of Jiang and Shu [**48**] described for 1D in the previous section is rather difficult to generalize to unstructured triangular and tetrahedral meshes in two and three space dimensions because of the need to determine the optimal linear weights, see [**46, 72, 93**]. Therefore, we present a different idea in this section which can be extended very easily to the unstructured case. The alternative reconstruction procedure described here for the one-dimensional case follows directly from the general guidelines given in [**24, 25**] for unstructured triangular and tetrahedral meshes in two and three space dimensions. The polynomial WENO reconstruction operator of this section produces *entire polynomials* $w_i(x)$, as the ENO approach proposed by Harten *et al.* in [**40**]. However, we formally write our method like a WENO scheme [**48, 57**] with a particularly simple choice for the linear weights. The most important difference of our approach compared to the classical WENO schemes of Jiang and Shu [**48**] described in the next section is that standard WENO methods reconstruct *point values* at the Gaussian integration points instead of an entire polynomial valid inside each control volume $T_i = [x_{i-\frac{1}{2}}; x_{i+\frac{1}{2}}]$.
Reconstruction is again done for each element on a reconstruction stencil $\mathcal{S}_i^s$, which is given by the following union of the element $T_i$ and its neighbors $T_j$,

$$(6.10) \qquad \mathcal{S}_i^s = \bigcup_{j=i+s-e}^{i+s+e} T_j,$$

where $s$ is the stencil shift with respect to the central element $T_i$ and $e$ is the spatial extension of the stencil to the left and the right. A central reconstruction stencil is given by $s = 0$, an entirely left-sided stencil is given by $s = -e$ and an entirely right-sided stencil is given by $s = e$. In our approach, we always will use *three fixed* reconstruction stencils $\mathcal{S}_i^0$, $\mathcal{S}_i^{-e}$ and $\mathcal{S}_i^e$.
Given the cell average data $Q_i^n$ in all elements $Q_i$ we are looking for a spatial

reconstruction polynomial obtained from $\mathcal{S}_i^s$ at time $t^n$ of the form

$$(6.11) \qquad w_i^s(\xi, t^n) = \sum_{l=0}^{M} \phi_l(\xi)\hat{w}_l^s(t^n) := \phi_l(\xi)\hat{w}_l^s(t^n),$$

with the spatial reconstruction basis functions $\phi_l(\xi)$. The $\phi_l(\xi)$ span the space of piecewise polynomials of degree $M$, and as basis functions one can use for example monomials, or, better, rescaled Legendre polynomials on the unit interval, which form an orthogonal basis. For each element $T_i$ we use a reference coordinate $0 \leq \xi \leq 1$ given by $x = x_{i-\frac{1}{2}} + \xi\Delta x$. In the following, we will use standard tensor index notation, implying summation over indices appearing twice. The number of polynomial coefficients (degrees of freedom) is $k = M+1$, where $M$ is the degree of the reconstruction polynomial and $k$ is the spatial order of accuracy of the scheme in space. To compute the reconstruction polynomial $w_i(\xi, t^n)$ valid for element $T_i$ we require again *integral conservation* for all elements $T_j$ inside the stencil $\mathcal{S}_i^s$, i.e.

$$(6.12) \qquad \int_{T_j} w_i^s(\xi, t^n)d\xi = \int_{T_j} \phi_l(\xi)d\xi \cdot \hat{w}_l^s(t^n) = Q_j^n, \qquad \forall T_j \in \mathcal{S}_i^s.$$

Equation (6.12) yields a linear equation system for the unknown coefficients $\hat{w}_l^s(t^n)$ of the reconstruction polynomial on stencil $\mathcal{S}_i^s$ that can be easily solved.

To obtain the final non-oscillatory reconstruction polynomials for each $T_i$ at time $t^n$, we finally construct a data-dependent nonlinear combination of the polynomials $w_i^0(\xi, t^n)$, $w_i^{-k}(\xi, t^n)$ and $w_i^k(\xi, t^n)$ obtained from the central, left-sided and right-sided stencils as follows:

$$(6.13) \qquad w_i(\xi, t^n) = \hat{w}_l(t^n)\phi_l(\xi),$$

with

$$(6.14) \qquad \hat{w}_l(t^n) = \omega_0\,\hat{w}_l^0(t^n) + \omega_{-k}\,\hat{w}_l^{-k}(t^n) + \omega_k\,\hat{w}_l^k(t^n).$$

The nonlinear weights $\omega_s$ are given by the relations

$$(6.15) \qquad \omega_s = \frac{\tilde{\omega}_s}{\tilde{\omega}_0 + \tilde{\omega}_{-k} + \tilde{\omega}_k}, \quad \tilde{\omega}_s = \frac{\lambda_s}{(\sigma_s + \epsilon)^r}.$$

The oscillation indicators $\sigma_s$ are computed as for pointwise WENO reconstructions:

$$(6.16) \qquad \sigma_s = \sum_{l=1}^{M} \int_0^1 \left(\frac{\partial^l}{\partial \xi^l} w_i^s(\xi, t^n)\right)^2 d\xi.$$

The parameters $\epsilon$ and $r$ are constants for which one typically chooses $\epsilon = 10^{-14}$ and $r = 8$. For the linear weights $\lambda_s$ one chooses $\lambda_{-k} = \lambda_k = 1$ and a very large linear weight $\lambda_0$ on the central stencil, typically $\lambda_0 = 10^5$. It has been shown previously [48, 57] that the numerical results are quite insensitive to the WENO parameters $\epsilon$ and $r$ and also with respect to the linear weight on the central stencil $\lambda_0$, see [24]. The proposed reconstruction usually uses the accurate and linearly stable central stencil reconstruction in those regions of $\Omega$ where the solution is smooth because of the large linear weight $\lambda_0$. However, due to the strongly nonlinear dependence of the weights $\omega_s$ on the oscillation indicators $\sigma_s$, in the presence of discontinuities the smoother left- or right-sided stencils are preferred, as for standard ENO and WENO methods. For the nonlinear scalar case, the reconstruction operator described above can be directly applied to the cell averages $Q_i^n$ of the conserved quantity $u$. For nonlinear hyperbolic systems, the reconstruction should be done in *characteristic variables* [40] in order to avoid spurious oscillations that may appear when applying ENO or WENO reconstruction operators component-wise to nonlinear hyperbolic

systems. The result of the reconstruction procedure is a non-oscillatory spatial polynomial $w_i(\xi, t^n)$ defined at time $t^n$ inside each spatial element $T_i$.

The advantage of the polynomial WENO reconstruction is its straight–forward extension to general unstructured meshes. The drawback is that at a given order of accuracy $k$ the total stencil needed for the reconstruction is *bigger* than the one of the classical pointwise WENO scheme presented before.

EXAMPLE 34. *For the third order polynomial WENO reconstruction described in this section, we use the rescaled Legendre polynomials up to degree two as reconstruction basis functions:*

$$(6.17) \qquad \phi_0(\xi) = 1, \qquad \phi_1(\xi) = 2\xi - 1, \qquad \phi_2(\xi) = 1 - 6\xi + 6\xi^2.$$

*It can be easily shown that these functions are orthogonal on the unit interval. We obtain the following expansion coefficients for the left–sided stencil*

$$(6.18)$$
$$\hat{w}_0^{-1} = Q_i^n, \quad \hat{w}_1^{-1} = \frac{1}{4}Q_{i-2}^n - Q_{i-1}^n + \frac{3}{4}Q_i^n, \quad \hat{w}_2^{-1} = \frac{1}{12}Q_{i-2}^n - \frac{1}{6}Q_{i-1}^n + \frac{1}{12}Q_i^n,$$

*for the central stencil*

$$(6.19) \quad \hat{w}_0^0 = Q_i^n, \quad \hat{w}_1^0 = -\frac{1}{4}Q_{i-1}^n + \frac{1}{4}Q_{i+1}^n, \quad \hat{w}_2^0 = \frac{1}{12}Q_{i-1}^n - \frac{1}{6}Q_i^n + \frac{1}{12}Q_{i+1}^n,$$

*and for the right–sided stencil*

$$(6.20)$$
$$\hat{w}_0^{+1} = Q_i^n, \quad \hat{w}_1^{+1} = -\frac{3}{4}Q_i^n + Q_{i+1}^n - \frac{1}{4}Q_{i+2}^n, \quad \hat{w}_2^{+1} = \frac{1}{12}Q_i^n - \frac{1}{6}Q_{i+1}^n + \frac{1}{12}Q_{i+2}^n.$$

*The oscillation indicator is given by*

$$(6.21) \qquad\qquad\qquad \sigma_s = 156\left(\hat{w}_2^s\right)^2 + 4\left(\hat{w}_1^s\right)^2.$$

# The Discontinuous Galerkin finite element method

In high order finite volume schemes, only the cell averages $Q_i^n$ are stored for each element and are evolved in time. To achieve higher order of accuracy in space the high order polynomials must be obtained from the cell–averages via a reconstruction procedure or recovery procedure, see the previous chapter. In the discontinuous Galerkin finite element framework, the coefficients of higher order polynomials are *directly* evolved in time for each cell, *without* the need of using a reconstruction operator. This feature of DG schemes is in common with the classical finite element method (FEM). However, unlike in classical finite elements, the numerical solution given by a DG scheme is *discontinuous at element interfaces* and this discontinuity is resolved by the use of a *numerical flux function*, which is a common feature with shock–capturing finite volume schemes.

The DG method was invented by Reed and Hill [65] for the solution of neutron transport equations and was put on a solid mathematical basis by Cockburn and Shu in a well–known series of papers [14, 13, 12, 11, 16]. The literature on DG finite element schemes is rapidly growing at the moment since the method has several advantages with respect to classical finite volume schemes: The DG method reaches arbitrary order of accuracy in space on general unstructured meshes and allows easily for $hp$–adaptation, i.e. it allows for refinement and recoarsening of the mesh and for a dynamical adaptation of the polynomial degree of the numerical solution. For the DG method, Jiang and Shu have found a very elegant and general proof of *nonlinear stability* in $L_2$–norm, see [47]. In the following we give a brief sketch of the method for the one–dimensional scalar case and refer to the literature for details, see for example the nice review paper by Cockburn and Shu [17].

## 1. General algorithm

The governing PDE

$$(7.1) \qquad q_t + f(q)_x = 0$$

is multiplied by a test function $\phi = \phi(x)$ from the space of piecewise polynomials of degree $N$ and is integrated over a spatial control volume $T_i = [x_{i-\frac{1}{2}}; x_{i+\frac{1}{2}}]$. Integration by parts yields

$$(7.2) \qquad \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \phi q_t \, dx + \phi_{i+\frac{1}{2}}^- f_{i+\frac{1}{2}} - \phi_{i-\frac{1}{2}}^+ f_{i-\frac{1}{2}} - \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \phi_x f(q) dx = 0,$$

where $f_{i+\frac{1}{2}} = f_{i+\frac{1}{2}}(q_{i+\frac{1}{2}}^-, q_{i+\frac{1}{2}}^+)$ is a classical numerical flux function evaluated at the boundary–extrapolated values of the solution $q_{i+\frac{1}{2}}^\pm$ at the left and the right side of the element interface $x_{i+\frac{1}{2}}$, as used in the finite volume context. Likewise, $\phi_{i+\frac{1}{2}}^\pm$ is the boundary–extrapolated value of the test function. For the numerical solution

$q_h = q(x, t)$ we make the following ansatz:

$$(7.3) \qquad q_h(x, t) = \sum_{l=0}^{N} \hat{u}_l(t) \phi_l(x).$$

Inserting Eqn. (7.3) into (7.2) and using the same test functions $\phi$ as the ones used for the ansatz (7.3), we obtain for each element $T_i$ a system of $N+1$ ordinary differential equations for the unknown coefficients $\hat{u}_l(t)$ of the numerical solution $q_h$. This ODE system of the so–called *semi–discrete* form of the DG scheme is then typically integrated using a TVD Runge–Kutta scheme in time, see [**16, 17**]. However, also other time–discretizations are possible, such as ADER or Lax–Wendroff–type one–step time discretizations [**64, 26**]. It is very important to underline that explicit DG finite element schemes require a *very severe* CFL–type restriction of the time step. An estimate for the maximum Courant number of common RK–DG schemes is

$$(7.4) \qquad \text{CFL} < \frac{1}{2N+1}.$$

For an extension of the DG method to diffusion equations and PDE with higher order derivatives, which is not as straightforward as it seems, see [**4, 15, 92, 56, 42, 35, 1**].

## 2. The TVB limiter of Cockburn and Shu

A still open topic of research is the construction of appropriate *limiters* for high order DG finite element schemes. In [**13, 12**] Cockburn and Shu proposed a TVB (totally variation bounded) limiter, which depends on a constant $K > 0$, to preserve local extrema. They introduced a *modified minmod function*

$$(7.5) \qquad \mathrm{m}(a, b, c) = \begin{cases} a & \text{if } |a| < K\Delta x^2 \\ \mathrm{minmod}(a, b, c), & \text{else,} \end{cases}$$

with the classical minmod function being defined as

$$(7.6) \qquad \mathrm{minmod}(a, b, c) = \begin{cases} \min(a, b, c) & \text{if } \mathrm{sign}(a) = \mathrm{sign}(b) = \mathrm{sign}(c), \\ 0 & \text{else.} \end{cases}$$

The TVB limiter of Cockburn and Shu is then applied as follows: Compute the $L_2$–projection of the numerical solution $q_h$ into the space of piecewise polynomials of degree one. The result of this projection is called $v_h$. For hierarchical orthogonal basis functions such as the Legendre polynomials in 1D or the orthogonal basis functions on general elements in multiple dimensions defined in [**49, 19**], this projection is simply done by setting the higher order coefficients to zero. Then, compute

$$(7.7) \qquad u_{i+\frac{1}{2}}^- = \bar{v}_i + m(v_{i+\frac{1}{2}}^- - \bar{v}_i, \bar{v}_i - \bar{v}_{i-1}, \bar{v}_{i+1} - \bar{v}_i),$$

and

$$(7.8) \qquad u_{i-\frac{1}{2}}^+ = \bar{v}_i - m(\bar{v}_i - v_{i-\frac{1}{2}}^+, \bar{v}_i - \bar{v}_{i-1}, \bar{v}_{i+1} - \bar{v}_i),$$

where the bar denotes the cell average value. In the case of hierarchical orthogonal basis functions, such as the Legendre polynomials, the cell average is directly given as the first expansion coefficient, associated with the piecewise constant term. If $u_{i+\frac{1}{2}}^- = v_{i+\frac{1}{2}}^-$ and $u_{i-\frac{1}{2}}^+ = v_{i-\frac{1}{2}}^+$, then no limitation is applied to the numerical solution. Otherwise, set in each element

$$(7.9) \qquad q_h(x) = \bar{v}_i + (x - x_i) \cdot m(v_x, \frac{\bar{v}_i - \bar{v}_{i-1}}{\Delta x/2}, \frac{\bar{v}_{i+1} - \bar{v}_i}{\Delta x/2}).$$

## 3. Nonlinear $L_2$–stability

As we have seen in the first part of this manuscript, high order linear Finite Difference schemes produce spurious oscillations near discontinuities. The reason can be traced back to Godunov's theorem [37], which states that there is no monotone linear difference scheme of order greater than one. Linear means in this case that the coefficients of the scheme are independent of the numerical solution.

Obviously, also the DG method is not able to circumvent the theorem of Godunov, but it has been found in literature that even the linear version of the Discontinuous Galerkin finite element method without using limiters exhibits a peculiar robustness in the presence of discontinuities. This was already mentioned by Atkins and Shu [2]. To shed some more light onto this interesting fact, we recall in this section the analysis carried out by Jiang and Shu [47] and its essential results. Jiang and Shu were able to prove rigorously a discrete cell entropy inequality for the square entropy for the Discontinuous Galerkin method when applied to scalar nonlinear hyperbolic conservation laws. This is a very remarkable result, since it was derived only under the assumption of the existence of a monotone (entropy– satisfying) numerical flux and the analysis is valid even without the use of limiters. The analysis was extended by Yan and Shu for local Discontinuous Galerkin schemes for PDEs with higher order derivatives [92].

**3.1. Generalities and definitions.** We start with some general definitions and lemmata. The nonlinear scalar hyperbolic conservation laws under consideration have the form

$$(7.10) \qquad\qquad q_t + f(q)_x = 0.$$

DEFINITION 4. *The pair of functions $(U(q), F(q))$ is called an entropy pair to the hyperbolic conservation law (7.10), if the entropy $U$ fulfills*

$$(7.11) \qquad\qquad U'' > 0$$

*and the entropy flux $F$ satisfies*

$$(7.12) \qquad\qquad F'(q) = U'(q)f'(q).$$

The superscript $'$ denotes differentiation with respect to $q$.

LEMMA 2. *If $q$ is a continuously differentiable solution of the conservation law (7.10) then any entropy pair $(U(q), F(q))$ fulfills the equation*

$$(7.13) \qquad\qquad U(q)_t + F(q)_x = 0.$$

PROOF. If $q$ is a continuously differentiable solution of (7.10), then

$$(7.14) \qquad\qquad q_t + f'q_x = 0.$$

Multiplying by $U'$

$$(7.15) \qquad\qquad U'q_t + U'f'q_x = 0$$

and using (7.12) one obtains

$$(7.16) \qquad\qquad U_t + F'q_x = U_t + F_x = 0.$$

$\square$

LEMMA 3. *If $q$ is a weak entropy solution [51] of the conservation law (7.10) then any entropy pair $(U(q), F(q))$ fulfills the inequality*

$$(7.17) \qquad\qquad U(q)_t + F(q)_x \leq 0.$$

PROOF. The proof can be found in [51]          $\square$

DEFINITION 5. *We introduce the square entropy*

$$(7.18) \qquad U_2 = \frac{q^2}{2}$$

*and its associated entropy flux*

$$(7.19) \qquad F_2 = qf(q) - \int f(q)dq.$$

LEMMA 4. *The function pair $(U_2, F_2)$ is an entropy pair.*

PROOF. The first property of Definition 4 follows directly as

$$(7.20) \qquad U_2''(q) = 1 > 0.$$

For the flux $F_2$ we have

$$(7.21) \qquad F_2'(q) = f(q) + qf'(q) - f(q)$$

and with $U_2'(q) = q$ follows the second property:

$$(7.22) \qquad F_2'(q) = U_2'(q)f'(q).$$

$\square$

DEFINITION 6. *A Lipschitz continuous function $f_{i+\frac{1}{2}} = f_{i+\frac{1}{2}}(q^-, q^+)$ of two states $q^-$ and $q^+$ is called an e-flux [61] for the conservation law (7.10), if it has the following properties:*

$$(7.23) \qquad f_{i+\frac{1}{2}}(q, q) = f(q).$$

$$(7.24) \qquad \int\limits_{q^-}^{q^+} \left( f(q) - f_{i+\frac{1}{2}}(q^-, q^+) \right) dq \geq 0.$$

**3.2. Discrete cell entropy inequality.** Approximating the solution of (7.10) by $q \in V_h$ in a discrete function space $V_h \subset L^2$ and multiplying (7.10) by test functions $\Phi \in V_h$ from the same function space and integrating by parts, one gets the following variational formulation

$$(7.25) \qquad \int_{T_i} q_t \Phi \ dx + \Phi_{i+\frac{1}{2}}^- f_{i+\frac{1}{2}} - \Phi_{i-\frac{1}{2}}^+ f_{i-\frac{1}{2}} - \int_{T_i} \frac{\partial}{\partial x} \Phi f(q) \ dx.$$

In (7.25) the $T_i = [x_{i-\frac{1}{2}}; x_{i+\frac{1}{2}}]$ are the subintervals into which the computational domain is divided, and $f_{i+\frac{1}{2}} = f_{i+\frac{1}{2}}(q_{i+\frac{1}{2}}^-, q_{i+\frac{1}{2}}^+)$ is an e-flux between the element $T_i$ and its right neighbor $T_{i+1}$ and $f_{i-\frac{1}{2}} = f_{i+\frac{1}{2}}(q_{i-\frac{1}{2}}^-, q_{i-\frac{1}{2}}^+)$ denotes an e-flux between $T_i$ and its left neighbor $T_{i-1}$. $\Phi_{i+\frac{1}{2}}^- = \Phi(x_{i+\frac{1}{2}}^-)$ and $\Phi_{i-\frac{1}{2}}^+ = \Phi(x_{i-\frac{1}{2}}^+)$ denote the values of the test functions inside $T_i$ at the interface with the neighbor elements.

THEOREM 5. *The numerical solution $q \in V_h$ of the Discontinuous Galerkin scheme (7.25) fulfills the discrete cell–entropy inequality*

$$(7.26) \qquad \int_{T_i} \left( \frac{q^2}{2} \right)_t dx + \hat{F}_{i+\frac{1}{2}} - \hat{F}_{i-\frac{1}{2}} \leq 0.$$

*for the square entropy.*

PROOF. As $q$ and $\Phi$ are from the same function space, we can also use $q$ as a test function instead of $\Phi$. We get

$$(7.27) \qquad \int_{T_i} \left( \frac{q^2}{2} \right)_t dx + q_{i+\frac{1}{2}}^- f_{i+\frac{1}{2}} - q_{i-\frac{1}{2}}^+ f_{i-\frac{1}{2}} - \int_{T_i} q_x f(q) \ dx.$$

If we rewrite the last term as

$$(7.28) \qquad \int_{T_i} f(q) q_x \, dx = \int_{q_{i-\frac{1}{2}}^+}^{q_{i+\frac{1}{2}}^-} f(q) dq = g(q_{i+\frac{1}{2}}^-) - g(q_{i-\frac{1}{2}}^+)$$

with the definition

$$(7.29) \qquad g(q) = \int f(q) dq,$$

eqn. (7.27) takes the form

$$(7.30) \qquad \int_{T_i} \left(\frac{q^2}{2}\right)_t dx + \left(q_{i+\frac{1}{2}}^- f_{i+\frac{1}{2}} - g(q_{i+\frac{1}{2}}^-)\right) - \left(q_{i-\frac{1}{2}}^+ f_{i-\frac{1}{2}} - g(q_{i-\frac{1}{2}}^+)\right) = 0.$$

Without changing anything, this can be rewritten as

$$\int_{T_i} \left(\frac{q^2}{2}\right)_t dx \quad + \quad \left(q_{i+\frac{1}{2}}^- f_{i+\frac{1}{2}} - g(q_{i+\frac{1}{2}}^-)\right) - \left(q_{i-\frac{1}{2}}^- f_{i-\frac{1}{2}} - g(q_{i-\frac{1}{2}}^-)\right)$$
$$- \quad \left(q_{i-\frac{1}{2}}^+ f_{i-\frac{1}{2}} - g(q_{i-\frac{1}{2}}^+)\right) + \left(q_{i-\frac{1}{2}}^- f_{i-\frac{1}{2}} - g(q_{i-\frac{1}{2}}^-)\right) = 0.$$
$$(7.31)$$

Eqn. (7.31) finally can be rewritten more compactly as

$$(7.32) \qquad \int_{T_i} \left(\frac{q^2}{2}\right)_t dx + \hat{F}_{i+\frac{1}{2}} - \hat{F}_{i-\frac{1}{2}} + \hat{R}_{i-\frac{1}{2}} = 0,$$

with the definitions

$$(7.33) \qquad \hat{F}_{i+\frac{1}{2}} = q_{i+\frac{1}{2}}^- f_{i+\frac{1}{2}} - g(q_{i+\frac{1}{2}}^-),$$

$$(7.34) \qquad \hat{F}_{i-\frac{1}{2}} = q_{i-\frac{1}{2}}^- f_{i-\frac{1}{2}} - g(q_{i-\frac{1}{2}}^-),$$

$$(7.35) \qquad \hat{R}_{i-\frac{1}{2}} = - \left(q_{i-\frac{1}{2}}^+ f_{i-\frac{1}{2}} - g(q_{i-\frac{1}{2}}^+)\right) + \left(q_{i-\frac{1}{2}}^- f_{i-\frac{1}{2}} - g(q_{i-\frac{1}{2}}^-)\right).$$

The terms $\hat{F}_{i+\frac{1}{2}}$ and $\hat{F}_{i-\frac{1}{2}}$ are discrete entropy fluxes that are consistent with the continuous entropy flux of the square entropy (7.19).

The term $\hat{R}_{i-\frac{1}{2}}$ can be written more compactly as

$$(7.36) \qquad \hat{R}_{i-\frac{1}{2}} = \int_{q_{i-\frac{1}{2}}^-}^{q_{i-\frac{1}{2}}^+} \left(f(q) - f_{i+\frac{1}{2}}\left(q_{i-\frac{1}{2}}^-, q_{i-\frac{1}{2}}^+\right)\right) dq.$$

Because of property (7.24) of the monotone flux this term is always non-negative,

$$(7.37) \qquad \hat{R}_{i-\frac{1}{2}} \geq 0,$$

so from (7.32) and (7.37) finally follows

$$(7.38) \qquad \int_{T_i} \left(\frac{q^2}{2}\right)_t dx + \hat{F}_{i+\frac{1}{2}} - \hat{F}_{i-\frac{1}{2}} \leq 0,$$

which concludes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

COROLLARY 1. *The semi-discrete Discontinuous Galerkin scheme (7.25) is $L^2$ stable.*

PROOF. Summing up (7.38) over all elements $T_i$ in the domain $\Omega$ and imposing either periodic boundary conditions or zero fluxes at the domain boundary $\partial\Omega$ yields

$$(7.39) \qquad \int_\Omega \left(\frac{q^2}{2}\right)_t dx \leq 0.$$

$\square$

The remarkable result is that the cell entropy inequality (7.38) and the resulting $L^2$ stability (7.39) follow for *arbitrary* high order semi-discrete DG schemes and for *any* nonlinear scalar hyperbolic conservation law, only under the assumption of an e-flux $f_{i+\frac{1}{2}}$ and supposing that the test and approximation spaces in the variational formulation are the same. The proof of Jiang and Shu is short and elegant and it is not necessary to introduce limiters in order to obtain the cell entropy inequality. The amount of numerical dissipation which is introduced into the numerical solution by the DG scheme is controlled by the term $\hat{R}_{i-\frac{1}{2}}$ which is intrinsically linked to the jump in the variable $q$ at the cell interface. The dissipation becomes the greater the higher the jump.

# One–step time discretization based on a local space–time discontinuous Galerkin scheme

In this last chapter we present a one–step time discretization for high order finite volume and discontinuous Galerkin finite element methods first proposed in [**22**]. Again, we limit us to the scalar one–dimensional case. Here, we also consider nonlinear algebraic source terms, hence the governing PDE reads

$$(8.1) \qquad q_t + f(q)_x = S(q).$$

The time discretization presented in the following is to a certain extent similar to the one–step time discretization of the MUSCL scheme of van Leer [**87**] and the one used in the ENO method of Harten et al. [**40**]. In those approaches, the high order polynomials $w(x, t^n)$ are *locally evolved in time* (without coupling to the neighbor elements) using the governing partial differential equation. The result is a space–time polynomial $q(x, t)$ for each element. In the MUSCL method this evolution is done using a first order Taylor series in time, where the first time derivative is obtained by the PDE (8.1) and the spatial derivative of the flux is approximated by a central finite difference using the reconstruction polynomial within the element. For the ENO scheme [**40**] the space–time polynomial is obtained by a higher order Taylor series, where also the higher order time derivatives are replaced by pure space derivatives using the Cauchy–Kovalewski procedure. Another one–step time–discretization is automatically achieved within the ADER schemes of Titarev and Toro [**76, 85, 75, 84, 78**] who used the solution of the *generalized Riemann problem* (GRP) [**31, 6, 32, 5, 83**] at the cell interface as a building block for obtaining high order of accuracy in time. Also the ADER approach is essentially based on the use of the Cauchy–Kovalewski procedure. Unfortunately, this procedure becomes very cumbersome for classical nonlinear hyperbolic systems, see [**25**] for the 3D compressible Euler equations and [**74**] for the 2D magnetohydrodynamics (MHD) equations and is almost impossible to do for the general nonlinear hyperbolic systems case. Furthermore, it is not applicable to the case of stiff source terms. Therefore, we show here a different approach that is based on a *weak* formulation of the governing PDE in space–time which is very general and also applicable to stiff source terms.

Final aim of the approach described in the following is to solve an *element–local* Cauchy problem with initial data $w(x, t^n)$. The initial data stems either from a high order WENO reconstruction, or is directly available as $w(x, t^n) = q_h(x, t^n)$ from the discontinuous Galerkin finite element framework. For the solution of the element–local Cauchy problem, we take the PDE (8.1), multiply it with a space–time test function $\theta_k = \theta_k(x, t)$ and integrate over the space–time control volume $[x_{i-\frac{1}{2}}; x_{i+\frac{1}{2}}] \times [t^n; t^{n+1}]$:

$$(8.2) \qquad \int\limits_{t^n}^{t^{n+1}} \int\limits_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \theta_k \left( q_t + f(q)_x - S(q) \right) \, dx \, dt = 0.$$

Integration by parts of the time–derivative term yields

$$
(8.3) \quad
\int\limits_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \left( \theta_k(x,t^{n+1})q(x,t^{n+1}) - \theta_k(x,t^n)w(x,t^n) \right) dx - \int\limits_{t^n}^{t^{n+1}} \int\limits_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \frac{\partial}{\partial t}\theta_k q \, dx \, dt +
$$

$$
\int\limits_{t^n}^{t^{n+1}} \int\limits_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} \theta_k \left( f_x - S \right) dx \, dt = 0.
$$

Here, we have introduced the initial condition $w(x,t^n)$ at the initial time $t^n$. Introduction of the two operators

$$
(8.4)
$$

$$
[f,g]_{T_i}^t = \int\limits_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} f(x,t)g(x,t)dx, \quad \text{and} \quad \langle f,g \rangle_{T_i} = \int\limits_{t^n}^{t^{n+1}} \int\limits_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} f(x,t)g(x,t) \, dx \, dt = 0.
$$

yields the following short–hand notation of Eqn. (8.3):

$$
(8.5) \quad [\theta_k, q]_{T_i}^{t^{n+1}} - \left\langle \frac{\partial}{\partial t}\theta_k, q \right\rangle_{T_i} + \langle \theta_k, f_x - S \rangle_{T_i} = [\theta_k, w]_{T_i}^{t^n}.
$$

Using the same basis functions $\theta_k$ to expand the solution $q$, the flux $f$ and the source $S$ in (8.5)

$$
(8.6) \quad q_h(x,t) = \sum_l \theta_l(x,t)\hat{q}_l := \theta_l \hat{q}_l, \qquad f_h(x,t) = \theta_l \hat{f}_l, \qquad S_h(x,t) = \theta_l \hat{S}_l
$$

and insertion into (8.5) leads to

$$
(8.7) \quad [\theta_k, \theta_l]_{T_i}^{t^{n+1}} \hat{q}_l - \left\langle \frac{\partial}{\partial t}\theta_k, \theta_l \right\rangle_{T_i} \hat{q}_l + \left\langle \theta_k, \frac{\partial}{\partial x}\theta_l \right\rangle_{T_i} \hat{f}_l - \langle \theta_k, \theta_l \rangle_{T_i} \hat{S}_l = [\theta_k, \phi]_{T_i}^{t^n} \hat{w}_l,
$$

or in even more compact matrix notation

$$
(8.8) \quad K_{kl}^1 \hat{q}_l + K_{kl}^\xi \hat{f}_l - M_{kl}\hat{S}_l = F_{kl}^0 \hat{w}_l,
$$

with the matrix $K_{kl}^1 = F_{kl}^1 - K_{kl}^\tau$, the spatial stiffness matrix $K_{kl}^\xi$, the mass matrix $M_{kl}$ and the temporal flux matrix at time $t^n$ $F_{kl}^0$ being defined as

$$
(8.9) \quad
\begin{aligned}
K_{kl}^1 &= [\theta_k, \theta_l]_{T_i}^{t^{n+1}} - \left\langle \frac{\partial}{\partial t}\theta_k, \theta_l \right\rangle_{T_i}, \quad K_{kl}^\xi = \left\langle \theta_k, \frac{\partial}{\partial x}\theta_l \right\rangle_{T_i}, \\
M_{kl} &= \langle \theta_k, \theta_l \rangle_{T_i}, \quad F_{kl}^0 = [\theta_k, \phi]_{T_i}^{t^n}.
\end{aligned}
$$

If we choose a *nodal basis* $\theta_k$, then the degrees of freedom for the flux and the source can be simply evaluated as

$$
(8.10) \quad \hat{f}_l = f(\hat{q}_l), \qquad \text{and} \qquad \hat{S}_l = S(\hat{q}_l).
$$

If we use a *modal basis*, then these coefficients have to be computed via appropriate $L_2$–projection. Eqn. (8.8) constitutes a *small element–local* nonlinear algebraic equation system which can be easily solved by the following iteration procedure proposed in [20],

$$
(8.11) \quad K_{kl}^1 \hat{q}_l^{(m+1)} - M_{kl}\hat{S}_l^{(m+1)} = F_{kl}^0 \hat{w}_l - + K_{kl}^\xi \hat{f}_l^{(m)},
$$

which for the linear homogeneous case can be shown to converge to the exact solution in at most $M+1$ steps, where $M$ is the polynomial degree of the basis functions $\theta_k$. The result of the local space–time DG method is a space–time polynomial for each cell which can be directly used within a high–order finite volume or DG finite element scheme.

# Bibliography

[1] D.N. Arnold, F. Brezzi, B. Cockburn, and L.D. Marini. Unified analysis of discontinuous Galerkin methods for elliptic problems. *SIAM Journal on Numerical Analysis*, 39:1749–1779, 2002.

[2] H. Atkins and C.W. Shu. Quadrature-free implementation of the discontinuous Galerkin method for hyperbolic equations. *AIAA Journal*, 36:775–782, 1998.

[3] E. Audusse, F. Bouchut, M.O. Bristeau, R. Klein, and B. Perthame. A fast and stable well-balanced scheme with hydrostatic reconstruction for shallow water flows. *SIAM Journal on Scientific Computing*, 25:2050–2065, 2004.

[4] F. Bassi and S. Rebay. A high-order accurate discontinuous finite element method for the numerical solution of the compressible Navier-Stokes equations. *Journal of Computational Physics*, 131:267–279, 1997.

[5] M. Ben-Artzi and J. Falcovitz. A second-order godunov-type scheme for compressible fluid dynamics. *Journal of Computational Physics*, 55:1–32, 1984.

[6] A. Bourgeade, P. LeFloch, and P.A. Raviart. An asymptotic expansion for the solution of the generalized riemann problem. Part II: application to the gas dynamics equations. *Annales de l'institut Henri Poincaré (C) Analyse non linéaire*, 6:437–480, 1989.

[7] A. Canestrelli, M. Dumbser, A. Siviglia, and E.F. Toro. Well-balanced high-order centered schemes on unstructured meshes for shallow water equations with fixed and mobile bed. *Advances in Water Resources*, 33:291–303, 2010.

[8] A. Canestrelli, A. Siviglia, M. Dumbser, and E.F. Toro. A well-balanced high order centered scheme for nonconservative systems: Application to shallow water flows with fix and mobile bed. *Advances in Water Resources*, 32:834–844, 2009.

[9] M.J. Castro, J.M. Gallardo, and C. Parés. High-order finite volume schemes based on reconstruction of states for solving hyperbolic systems with nonconservative products. applications to shallow-water systems. *Mathematics of Computation*, 75:1103–1134, 2006.

[10] M.J. Castro, P.G. LeFloch, M.L. Muñoz-Ruiz, and C. Parés. Why many theories of shock waves are necessary: Convergence error in formally path-consistent schemes. *Journal of Computational Physics*, 227:8107–8129, 2008.

[11] B. Cockburn, S. Hou, and C. W. Shu. The Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws IV: the multidimensional case. *Mathematics of Computation*, 54:545–581, 1990.

[12] B. Cockburn, S. Y. Lin, and C.W. Shu. TVB Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws III: one dimensional systems. *Journal of Computational Physics*, 84:90–113, 1989.

[13] B. Cockburn and C. W. Shu. TVB Runge-Kutta local projection discontinuous Galerkin finite element method for conservation laws II: general framework. *Mathematics of Computation*, 52:411–435, 1989.

[14] B. Cockburn and C. W. Shu. The Runge-Kutta local projection P1-Discontinuous Galerkin finite element method for scalar conservation laws. *Mathematical Modelling and Numerical Analysis*, 25:337–361, 1991.

[15] B. Cockburn and C. W. Shu. The local discontinuous Galerkin method for time-dependent convection diffusion systems. *SIAM Journal on Numerical Analysis*, 35:2440–2463, 1998.

[16] B. Cockburn and C. W. Shu. The Runge-Kutta discontinuous Galerkin method for conservation laws V: multidimensional systems. *Journal of Computational Physics*, 141:199–224, 1998.

[17] B. Cockburn and C. W. Shu. Runge-Kutta discontinuous Galerkin methods for convection-dominated problems. *Journal of Scientific Computing*, 16:173–261, 2001.

[18] R. Courant, K. Friedrichs, and H. Lewy. über die partiellen differenzengleichungen der mathematischen physik. *Mathematische Annalen*, 100:32–74, 1928.

[19] M. Dubiner. Spectral methods on triangles and other domains. *Journal of Scientific Computing*, 6:345–390, 1991.

[20] M. Dumbser, D. Balsara, E.F. Toro, and C.D. Munz. A unified framework for the construction of one-step finite-volume and discontinuous Galerkin schemes. *Journal of Computational Physics*, 227:8209–8253, 2008.

[21] M. Dumbser, M. Castro, C. Parés, and E.F. Toro. ADER schemes on unstructured meshes for nonconservative hyperbolic systems: Applications to geophysical flows. *Computers and Fluids*, 38:1731–1748, 2009.

[22] M. Dumbser, C. Enaux, and E.F. Toro. Finite volume schemes of very high order of accuracy for stiff hyperbolic balance laws. *Journal of Computational Physics*, 227:3971–4001, 2008.

[23] M. Dumbser, A. Hidalgo, M. Castro, C. Parés, and E.F. Toro. FORCE schemes on unstructured meshes II: Non–conservative hyperbolic systems. *Computer Methods in Applied Mechanics and Engineering*, 2010.

[24] M. Dumbser and M. Käser. Arbitrary high order non-oscillatory finite volume schemes on unstructured meshes for linear hyperbolic systems. *Journal of Computational Physics*, 221:693–723, 2007.

[25] M. Dumbser, M. Käser, V.A Titarev, and E.F. Toro. Quadrature-free non-oscillatory finite volume schemes on unstructured meshes for nonlinear hyperbolic systems. *Journal of Computational Physics*, 226:204–243, 2007.

[26] M. Dumbser and C.D. Munz. Building blocks for arbitrary high order discontinuous Galerkin schemes. *Journal of Scientific Computing*, 27:215–230, 2006.

[27] M. Dumbser and E.F. Toro. On universal Osher–type schemes for general nonlinear hyperbolic conservation laws. *Communications in Computational Physics*. in press.

[28] M. Dumbser and E.F. Toro. A simple extension of the Osher Riemann solver to non-conservative hyperbolic systems. *Journal of Scientific Computing*. in press, DOI: 10.1007/s10915-010-9400-3.

[29] B. Einfeldt. On godunov-type methods for gas dynamics. *SIAM Journal on Numerical Analysis*, 25:294–318, 1988.

[30] B. Einfeldt, C. D. Munz, P. L. Roe, and B. Sjögreen. On godunov-type methods near low densities. *Journal of Computational Physics*, 92:273–295, 1991.

[31] P. Le Floch and P.A. Raviart. An asymptotic expansion for the solution of the generalized riemann problem. Part I: General theory. *Annales de l'institut Henri Poincaré (C) Analyse non linéaire*, 5:179–207, 1988.

[32] P. Le Floch and L. Tatsien. A global asymptotic expansion for the solution of the generalized riemann problem. *Annales de l'institut Henri Poincaré (C) Analyse non linéaire*, 3:321–340, 1991.

[33] J.E. Fromm. A method for reducing dispersion in convective difference schemes. *Journal of Computational Physics*, 3:176–189, 1968.

[34] P. Garcia-Navarro and M.E. Vázquez-Cendón. On numerical treatment of the source terms in the shallow water equations. *Computers & Fluids*, 29:951–979, 2000.

[35] G. Gassner, F. Lörcher, and C.D. Munz. A contribution to the construction of diffusion fluxes for finite volume and discontinuous Galerkin schemes. *Journal of Computational Physics*, 224:1049–1063, 2007.

[36] E. Godlewski and P. A. Raviart. *Numerical Approximation of Hyperbolic Systems of Conservation Laws*. Springer, 1996.

[37] S.K. Godunov. Finite difference methods for the computation of discontinuous solutions of the equations of fluid dynamics. *Mathematics of the USSR: Sbornik*, 47:271–306, 1959.

[38] J. M. Greenberg and A. Y. Le Roux. A wellbalanced scheme for the numerical processing of source terms in hyperbolic equations. *SIAM Journal on Numerical Analysis*, 33:1–16, 1996.

[39] A. Harten. High Resolution Schemes for Hyperbolic Conservation Laws. *J. Comput. Phys.*, 49:357–393, 1983.

[40] A. Harten, B. Engquist, S. Osher, and S. Chakravarthy. Uniformly high order essentially non-oscillatory schemes, III. *Journal of Computational Physics*, 71:231–303, 1987.

[41] A. Harten, P.D. Lax, and B. van Leer. On upstream differencing and godunov-type schemes for hyperbolic conservation laws. *SIAM Review*, 25(1):35–61, 1983.

[42] R. Hartmann and P. Houston. Symmetric interior penalty DG methods for the compressible navier–stokes equations I: Method formulation. *Int. J. Num. Anal. Model.*, 3:1–20, 2006.

[43] C. Hirsch. *Numerical Computation of Internal and External Flows Vol I: Fundamentals of Numerical Discretisation*. Wiley, 1988.

[44] C. Hirsch. *Numerical Computation of Internal and External Flows Vol II: Computational Methods for Inviscid and Viscous Flow*. Wiley, 1988.

[45] T.Y. Hou and P.G. LeFloch. Why nonconservative schemes converge to wrong solutions: error analysis. *Mathematics of Computation*, 62:497–530, 1994.

[46] C. Hu and C.W. Shu. Weighted essentially non-oscillatory schemes on triangular meshes. *Journal of Computational Physics*, 150:97–127, 1999.

[47] G. Jiang and C.W. Shu. On a cell entropy inequality for discontinuous Galerkin methods. *Mathematics of Computation*, 62:531–538, 1994.

[48] G.-S. Jiang and C.W. Shu. Efficient implementation of weighted ENO schemes. *Journal of Computational Physics*, 126:202–228, 1996.

[49] G. E. Karniadakis and S. J. Sherwin. *Spectral/hp Element Methods in CFD*. Oxford University Press, 1999.

[50] V. P. Kolgan. Application of the minimum-derivative principle in the construction of finite-difference schemes for numerical analysis of discontinuous solutions in gas dynamics. *Transactions of the Central Aerohydrodynamics Institute*, 3(6):68–77, 1972. in Russian.

[51] D. Kröner. *Numerical Schemes for Conservation Laws*. Wiley - Teubner, 1997.

[52] P.D. Lax. Weak solutions of nonlinear hyperbolic equations and their numerical computation. *Communications in Pure and Applied Mathematics*, 7:159–193, 1954.

[53] P.D. Lax and B. Wendroff. Systems of conservation laws. *Communications in Pure and Applied Mathematics*, 13:217–237, 1960.

[54] R. J. LeVeque. Balancing source terms and flux gradients in highresolution Godunov methods. *Journal of Computational Physics*, 146:346–365, 1998.

[55] R. J. LeVeque. *Finite Volume Methods for Hyperbolic Problems*. Cambridge University Press, 2002.

[56] D. Levy, C. W. Shu, and J. Yan. Local discontinuous Galerkin methods for nonlinear dispersive equations. *Journal of Computational Physics*, 196:751–772, 2004.

[57] X.D. Liu, S. Osher, and T. Chan. Weighted essentially non-oscillatory schemes. *Journal of Computational physics*, 115:200–212, 1994.

[58] G. Dal Maso, P.G. LeFloch, and F. Murat. Definition and weak stability of nonconservative products. *J. Math. Pures Appl.*, 74:483–548, 1995.

[59] S. Noelle, N. Pankratz, G. Puppo, and J.R. Natvig. Well-balanced finite volume schemes of arbitrary order of accuracy for shallow water flows. *Journal of Computational Physics*, 213:474–499, 2006.

[60] S. Noelle, Y.L. Xing, and C.W. Shu. High-order well-balanced finite volume WENO schemes for shallow water equation with moving water. *Journal of Computational Physics*, 226:29–58, 2007.

[61] S. Osher. Riemann solvers, the entropy condition and difference approximations. *SIAM Journal on Numerical Analysis*, 21:217–235, 1984.

[62] S. Osher and F. Solomon. Upwind difference schemes for hyperbolic conservation laws. *Math. Comput.*, 38:339–374, 1982.

[63] C. Parés. Numerical methods for nonconservative hyperbolic systems: a theoretical framework. *SIAM Journal on Numerical Analysis*, 44:300–321, 2006.

[64] J. Qiu, M. Dumbser, and C.W. Shu. The discontinuous Galerkin method with Lax-Wendroff type time discretizations. *Computer Methods in Applied Mechanics and Engineering*, 194:4528–4543, 2005.

[65] W.H. Reed and T.R. Hill. Triangular mesh methods for the neutron transport equation. Technical Report LA-UR-73-479, Los Alamos Scientific Laboratory, 1973.

[66] S. Rhebergen, O. Bokhove, and J.J.W. van der Vegt. Discontinuous Galerkin finite element methods for hyperbolic nonconservative partial differential equations. *Journal of Computational Physics*, 227:1887–1922, 2008.

[67] B. Riemann. Über die Fortpflanzung ebener Luftwellen von endlicher Schwingungsweite. *Göttinger Nachrichten*, 19, 1859.

[68] B. Riemann. Über die Fortpflanzung ebener Luftwellen von endlicher Schwingungsweite. *Abhandlungen der Königlichen Gesellschaft der Wissenschaften zu Göttingen*, 8:43–65, 1860.

[69] P.L. Roe. Approximate Riemann solvers, parameter vectors, and difference schemes. *Journal of Computational Physics*, 43:357–372, 1981.

[70] V. V. Rusanov. Calculation of Interaction of Non–Steady Shock Waves with Obstacles. *J. Comput. Math. Phys. USSR*, 1:267–279, 1961.

[71] E.N. Sarmin and L.A. Chudov. On the stability of the numerical integration of systems of ordinary differential equations arising in the use of the straight line method. *USSR Computational Mathematics and Mathematical Physics*, 3:1537–1543, 1963.

[72] J. Shi, C. Hu, and C.W. Shu. A technique of treating negative weights in WENO schemes. *Journal of Computational Physics*, 175:108–127, 2002.

[73] Y.I. Shokin. *The method of differential approximation*. Springer Verlag, 1983.

[74] A. Taube, M. Dumbser, D. Balsara, and C.D. Munz. Arbitrary high order discontinuous Galerkin schemes for the magnetohydrodynamic equations. *Journal of Scientific Computing*, 30:441–464, 2007.

[75] V.A. Titarev and E.F. Toro. ADER: Arbitrary high order Godunov approach. *Journal of Scientific Computing*, 17(1-4):609–618, December 2002.

[76] V.A. Titarev and E.F. Toro. ADER schemes for three-dimensional nonlinear hyperbolic sys-
     tems. *Journal of Computational Physics*, 204:715–736, 2005.

[77] E. F. Toro. *Riemann Solvers and Numerical Methods for Fluid Dynamics, Second Edition.*
     Springer–Verlag, 1999.

[78] E. F. Toro and V. A. Titarev. Derivative Riemann solvers for systems of conservation laws
     and ADER methods. *Journal of Computational Physics*, 212(1):150–165, 2006.

[79] E.F. Toro. *Riemann Solvers and Numerical Methods for Fluid Dynamics.* Springer, third
     edition, 2009.

[80] E.F. Toro and S. J. Billet. Centered TVD schemes for hyperbolic conservation laws. *IMA
     Journal of Numerical Analysis*, 20:44–79, 2000.

[81] E.F. Toro, A. Hidalgo, and M. Dumbser. FORCE schemes on unstructured meshes I: Con-
     servative hyperbolic systems. *Journal of Computational Physics*, 228:3368–3389, 2009.

[82] E.F. Toro, M. Spruce, and W. Speares. Restoration of the contact surface in the Harten-Lax-
     van Leer Riemann solver. *Journal of Shock Waves*, 4:25–34, 1994.

[83] E.F. Toro and V. A. Titarev. Solution of the generalized Riemann problem for advection-
     reaction equations. *Proc. Roy. Soc. London*, pages 271–281, 2002.

[84] E.F. Toro and V.A. Titarev. Very high order godunov-type schemes for nonlinear scalar
     conservation laws. In *Proceedings of ECCOMAS CFD Conference 2001*. ECCOMAS CFD
     Conference, 2001.

[85] E.F. Toro and V.A. Titarev. ADER schemes for scalar hyperbolic conservation laws with
     source terms in three space dimensions. *Journal of Computational Physics*, 202:196–215,
     2005.

[86] I. Toumi. A weak formulation of Roe's approximate Riemann solver. *Journal of Computa-
     tional Physics*, 102:360–373, 1992.

[87] B. van Leer. Towards the ultimate conservative difference scheme V: A second order sequel
     to Godunov's method. *Journal of Computational Physics*, 32:101–136, 1979.

[88] M.E. Vázquez-Cendón. Improved treatment of source terms in upwind schemes for the shal-
     low water equations in channels with irregular geometry. *Journal of Computational Physics*,
     148:497–526, 1999.

[89] D.H. Wagner. Equivalence of euler and lagrangian equations of gas dynamics for weak solu-
     tions. *Journal of Differential Equations*, 68:118–136, 1987.

[90] R.F. Warming and B.J. Hyett. The modified equation approach to the stability and accuracy
     analysis of finite-difference methods. *Journal of Computational Physics*, 14:159–179, 1974.

[91] K. Xu. A well-balanced gas-kinetic scheme for the shallow-water equations with source terms.
     *Journal of Computational Physics*, 178:533–562, 2002.

[92] J. Yan and C.W. Shu. A local discontinuous Galerkin method for KdV-type equations. *SIAM
     Journal on Numerical Analysis*, 40:769–791, 2002.

[93] Y.T. Zhang and C.W. Shu. Third order WENO scheme on three dimensional tetrahedral
     meshes. *Communications in Computational Physics*, 5:836–848, 2009.