



**HAL**  
open science

# Hyperbolic balance laws: residual distribution, local and global fluxes

Remi Abgrall, Mario Ricchiuto

► **To cite this version:**

Remi Abgrall, Mario Ricchiuto. Hyperbolic balance laws: residual distribution, local and global fluxes. 2022. hal-03508452v1

**HAL Id: hal-03508452**

**<https://inria.hal.science/hal-03508452v1>**

Preprint submitted on 3 Jan 2022 (v1), last revised 24 Jan 2022 (v2)

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

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

# Hyperbolic balance laws: residual distribution, local and global fluxes

Rémi Abgrall

Institute of Mathematics, University of Zürich, Switzerland,  
remi.abgrall@math.uzh.ch

Mario Ricchiuto

Team CARDAMOM, Inria Bordeaux Sud-Ouest, Talence, France,  
mario.ricchiuto@inria.fr

September 20, 2021

## Abstract

This review paper describes a class of scheme named "residual distribution schemes" or "fluctuation splitting schemes". They are a generalization of Roe's numerical flux [61] in fluctuation form. The so-called multidimensional fluctuation schemes have historically first been developed for steady homogeneous hyperbolic systems. Their application to unsteady problems and conservation laws has been really understood only relatively recently. This understanding has allowed to make of the residual distribution framework a powerful playground to develop numerical discretizations embedding some prescribed constraints. This paper describes in some detail these techniques, with several examples, ranging from the compressible Euler equations to the Shallow Water equations.

## 1 Introduction

We are interested in the numerical approximation of partial differential equations relevant in fluid dynamics. For the objectives of the present paper, we will focus on the Euler and Navier-Stokes equations on complex domains, as well as on the shallow water equations. These models are particular cases of the system of balance laws:

$$\frac{\partial \mathbf{u}}{\partial t} + \operatorname{div} \mathbf{f}(\mathbf{u}) = \operatorname{div} \mathbf{f}_v(\mathbf{u}, \nabla \mathbf{u}) + S(\mathbf{u}, \mathbf{x})$$

with initial and boundary conditions. The vector  $\mathbf{u}$  denotes a set of *conserved* variables, which are often (but not always) densities of conserved macroscopic quantities (mass, energy, etc). For the Euler and Navier-Stokes equations we have

$$\mathbf{u} = (\rho, \rho v, E)^T$$

where as usual,  $\rho$  is the mass density,  $\mathbf{v} \in \mathbb{R}^d$  ( $d = 1, 2, 3$ ) the velocity,  $E = e + \frac{1}{2}\rho\mathbf{v}^2$  is the total energy density with  $e = e(\rho, p)$  is the internal energy,  $p$  is the pressure. Here also,  $\mathbf{f}(\mathbf{u})$  is the inviscid flux

$$\mathbf{f}(\mathbf{u}) = \begin{pmatrix} \rho\mathbf{v} \\ \rho\mathbf{v} \otimes \mathbf{v} + p\text{Id} \\ \mathbf{v}(E + p) \end{pmatrix},$$

$\mathbf{f}_v$  is the viscous flux,  $S$  is a source term, and  $\text{Id}$  is the  $d \times d$  identity matrix.

In the case of the shallow water equations, we have  $\mathbf{u} = (H, H\mathbf{v})^T$  where  $H$  is the water height, the inviscid flux is

$$\mathbf{f}(\mathbf{u}) = \begin{pmatrix} H\mathbf{v} \\ H\mathbf{v} \otimes \mathbf{v} + p(H)\text{Id} \end{pmatrix}$$

with  $p(H) = \frac{g}{2}H^2$ , the viscous flux vanishes and the source term is

$$S(\mathbf{u}, \mathbf{x}) = \begin{pmatrix} \mathbf{0} \\ gH\nabla b(x) + c_F(\mathbf{u})H\mathbf{v} \end{pmatrix}$$

where  $b$  is the bottom topography, and  $c_F(\mathbf{u})$  a friction coefficient modelling the effects of the boundary layer on the sea-floor.

In this paper, we are mostly interested in the inviscid case, so we will drop the viscous term. See however [14, 9]. The general form we will discuss is

$$\frac{\partial \mathbf{u}}{\partial t} + \text{div } \mathbf{f}(\mathbf{u}) = S(\mathbf{u}, \mathbf{x}) \quad (1)$$

with initial and boundary conditions.

Finite volume methods, and more recently Discontinuous Galerkin (DG) methods, are very popular because they are known to be locally conservative. Thanks to Lax-Wendroff theorem, this is a very desirable property since it guaranty that "if all goes well", the limit solution is a weak solution of the problem. The extension of this notion to the non-homogenous case, and the construction finite volumes and DG schemes properly accounting for source terms is still a very open research subject, and over the years several interesting approaches have been proposed. Note that in this case system (1) admits non trivial solutions (steady and time dependent), and the consistency with these solutions is often a desired property for the schemes. Indeed, there is a close connection between the two: in the Lax-Wendroff theorem, one essential condition on the numerical flux is that of consistency. When all the arguments of the numerical flux are equal, we must recover the continuous flux. This is another way of saying that uniform solutions must be preserved. When source terms are present this is not necessarily the case. For example, the flat free surface state

$$H + b(x) = \eta_0 = \text{const} \quad (2)$$

with still water  $\mathbf{v} = 0$  is undoubtedly physically more relevant than constant  $\mathbf{u}$ . For channels with smooth surfaces, if friction is neglected the constant flux and constant energy steady state

$$H\mathbf{v} = q_0 = \text{const}, \quad \frac{\mathbf{v}^2}{2} + g(H + b(x)) = \mathcal{E}_0 = \text{const} \quad (3)$$

becomes the relevant one in general. This state is also compatible with the appearance of hydraulic jumps, across which the energy level is modified (see e.g. [21]). In many other applications however friction

cannot be neglected, and the most general form of steady state is obtained from the solution of

$$Hv = q_0 = \text{const}, \quad q_0^2/H + gH^2/2 + \int_{x_0}^x \left( gH \frac{\partial b}{\partial x} + c_F(q_0, H)q_0 \right) dx = q_0^2/H_0 + gH_0^2/2$$

If the bathymetry is linear, for example  $b(x) = b_0 - \xi_0 x$ , one can show that a constant state  $\mathbf{u}_0 = (H_0, q_0)$  satisfies the equilibrium [54]

$$-gH_0\xi_0 + c_F(q_0, H_0)q_0 = 0 \tag{4}$$

For more general definitions of  $b(x)$ , it is less apparent that a set of constant states can be associated to the steady equilibrium.

From the discrete approximation point of view, the problem is, how to modify a given numerical flux so that these states are preserved, possibly within machine accuracy. This is often an ad-hoc construction, and if one has a different problem depending on the application, and the flux correction needs to be re-written almost from scratch.

If one looks at the literature, there are other type of schemes. For example the stabilized variational methods using continuous finite elements such as the SUPG scheme [42], or the Galerkin scheme with jump stabilisation due to Burman et al.[24]. There are also the fluctuation splitting schemes that were initially designed by Roe and co-authors[67], and later extended to high order, steady and unsteady problems, as well as the shallow water equations [1, 2, 55, 71, 58, 11, 14, 9, 13, 30]. None of these scheme are formulated initially in terms of local fluxes, but they are working well. Indeed, in the numerical folklore, these schemes are often claimed not to be locally conservative, despite the contrary having been shown in several works [25, 41, 13, 4, 5]. The most interesting aspect for these is that treating (1) with or without source term involves no major modification, and no special tricks.

The purpose of this paper is to recall that when  $S = 0$ , residual distribution and continuous finite elements are locally conservative. In fact, we will also recall how to construct an *equivalent* flux formulation, and provide some explicit examples. When  $S \neq 0$ , we show how the schemes have been naturally extended to embed the source term. We explain how to link them to more recent flux based formulations, although this is not how they are designed, and their construction is much more natural. In one space dimension, we also recall that they have some relation to the so-called path conservative schemes (see [26] and references therein).

The format of this paper is as follows. First we recall the main discrete prototype we are interested in, written in a residual distribution form. In a second part, we show for steady problem that they have an equivalent flux formulation, and we explicitly construct the flux. In a third part, we extend this to unsteady problems. In a fourth part, we show (in 1D only) why these methods are agnostic to flux. Numerical examples are also given.

Throughout the paper, and for simplicity, we will not consider boundary conditions, even-though this is of course doable.

## 2 Geometrical notations

Let us fix here the main notation used for the mesh and related geometrical entities. The computational domain  $\Omega$ ,  $d = 1, 2, 3$ , is covered by a tessellation  $\mathcal{T}_h$ . We denote by  $\mathcal{E}_h$  the set of internal edges/faces of  $\mathcal{T}_h$ , and by  $\mathcal{F}_h$  the set of boundary faces. Mesh elements are generically denoted by  $K$ , while we use  $e$  for a face/edge  $e \in \mathcal{E}_h \cup \mathcal{F}_h$ . The mesh is assumed to be shape regular,  $h_K$  represents the diameter of the element  $K$ . Similarly, if  $e \in \mathcal{E}_h \cup \mathcal{F}_h$ ,  $h_e$  represents its diameter.

Throughout this paper, we follow Ciarlet's definition [29, 37] of a finite element approximation: we have a set of degrees of freedom  $\Sigma_K$  of linear forms acting on the set  $\mathbb{P}^k$  of polynomials of degree  $k$  such that the linear mapping

$$q \in \mathbb{P}^k \mapsto (\sigma_1(q), \dots, \sigma_{|\Sigma_K|}(q))$$

is one-to-one. The space  $\mathbb{P}^k$  is spanned by the basis function  $\{\varphi_\sigma\}_{\sigma \in \Sigma_K}$  defined by

$$\forall \sigma, \sigma', \sigma(\varphi_{\sigma'}) = \delta_{\sigma'}^{\sigma}.$$

We have in mind either Lagrange interpolations where the degrees of freedom are associated to points in  $K$ , or other type of polynomials approximation such as Bézier polynomials where we will also do the same geometrical identification. Considering all the elements covering  $\mathbb{R}^d$ , the set of degrees of freedom is denoted by  $\mathcal{S}$  and a generic degree of freedom by  $\sigma$ . We note that for any  $K$ ,

$$\forall \mathbf{x} \in K, \quad \sum_{\sigma \in \Sigma_K} \varphi_\sigma(\mathbf{x}) = 1.$$

For any element  $K$ ,  $\#K$  is the number of degrees of freedom in  $K$ .

The integer  $k$  is assumed to be the same for any element. We define

$$\mathcal{V}^h = \bigoplus_K \{\mathbf{v} \in L^2(K), \mathbf{v}|_K \in \mathbb{P}^k\}$$

where  $\mathbb{P}^k$  is the set of polynomials of degree less or equal to  $k$ . The solution will be sought for in a space  $V^h$  that is:

- Either  $V^h = \mathcal{V}^h$ . In that case, the elements of  $V^h$  can be discontinuous across internal faces/edges of  $\mathcal{T}_h$ . There is no conformity requirement on the mesh.
- Or  $V^h = \mathcal{V}_h \cap C^0(\Omega)$  in which case the mesh needs to be conformal.

Throughout the text, we need to integrate functions. This is done via quadrature formula, and the symbol  $\oint$  used in volume integrals

$$\oint_K v(\mathbf{x}) \, d\mathbf{x}$$

or boundary integrals

$$\oint_{\partial K} v(\mathbf{x}) \, d\gamma.$$

Note that the integration domain is uniquely defined by the limits of the integral, while the symbol  $\oint$  is used here to explicitly denote discrete integration via user defined quadrature formulas.

If  $e \in \mathcal{E}_h$ , represents any *internal* edge, i.e.  $e \subset K \cap K^+$  for two elements  $K$  and  $K^+$ , we define for any function  $\psi$  the jump  $[[\nabla\psi]] = \nabla\psi|_K - \nabla\psi|_{K^+}$ . Here the choice of  $K$  and  $K^+$  is important, and defines an orientation. Similarly,  $\{\mathbf{v}\} = \frac{1}{2}(\mathbf{v}|_K + \mathbf{v}|_{K^+})$ .

If  $\mathbf{x}$  and  $\mathbf{y}$  are two vectors of  $\mathbb{R}^q$ , for  $q$  integer,  $\langle \mathbf{x}, \mathbf{y} \rangle$  is their scalar product. In some occasions, it can also be denoted as  $\mathbf{x} \cdot \mathbf{y}$  or  $\mathbf{x}^T \mathbf{y}$ . We also use  $\mathbf{x} \cdot \mathbf{y}$  when  $\mathbf{x}$  is a matrix and  $\mathbf{y}$  a vector: it is simply the matrix-vector multiplication.

In section 4, we have to deal with oriented graph. Given two vertices of this graph  $\sigma$  and  $\sigma'$ , we write  $\sigma > \sigma'$  to say that  $[\sigma, \sigma']$  is a direct edge.

Section 5 mostly deals with one dimensional problems in  $\mathbb{R}$ . Here, a mesh is defined from a increasing sequence of points  $\{x_\sigma\}_{\sigma \in \mathbb{Z}}$ , the elements are the intervals

$$K_{\sigma+1/2} := [x_\sigma, x_{\sigma+1}].$$

The length of  $K_{\sigma+1/2}$  is  $\Delta_{\sigma+1/2}x = x_{\sigma+1} - x_\sigma$  and we set

$$\Delta_\sigma x = \frac{\Delta_{\sigma+1/2}x + \Delta_{\sigma-1/2}x}{2}.$$

### 3 Example of schemes and conservation

Let us provide several examples for the approximation of (1), to begin with, without source terms and in the steady case. They are

- The SUPG [42] variational formulation, with  $\mathbf{u}^h, \mathbf{v}^h \in V^h = \mathcal{V}^h \cap C^0(\mathbb{R}^d)$ :

$$a(\mathbf{u}^h, \mathbf{v}^h) := - \int_{\Omega} \nabla \mathbf{v}^h \cdot \mathbf{f}(\mathbf{u}^h) d\mathbf{x} + \sum_{K \subset \Omega} h_K \int_K [\nabla \mathbf{f}(\mathbf{u}^h) \cdot \nabla \mathbf{v}^h] \tau_K [\nabla \mathbf{f}(\mathbf{u}^h) \cdot \nabla \mathbf{u}^h] d\mathbf{x} \quad (5)$$

+ Boundary terms.

Here  $\tau_K$  is a positive parameter, or a positive definite matrix in the system case<sup>1</sup>.

- The Galerkin scheme with jump stabilization, see [24] for details. We have

$$a(\mathbf{u}^h, \mathbf{v}^h) := - \int_{\Omega} \nabla \mathbf{v}^h \cdot \mathbf{f}(\mathbf{u}^h) d\mathbf{x} + \sum_{e \subset \Omega} \theta_e h_e^2 \int_e [[\nabla \mathbf{v}^h]] \cdot [[\nabla \mathbf{u}^h]] d\gamma \quad (6)$$

+ Boundary terms.

Here,  $\mathbf{u}^h, \mathbf{v}^h \in V^h = \mathcal{V}^h \cap C^0(\Omega)$ , and  $\theta_e$  is a positive parameter.

- The discontinuous Galerkin formulation: we look for  $\mathbf{u}^h, \mathbf{v}^h \in V^h = \mathcal{V}^h$  such that

$$a(\mathbf{u}^h, \mathbf{v}^h) := \sum_{K \subset \Omega} \left( - \int_K \nabla \mathbf{v}^h \cdot \mathbf{f}(\mathbf{u}^h) d\mathbf{x} + \int_{\partial K} \mathbf{v}^h \cdot \hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}^h, \mathbf{u}^{h,+}) d\gamma \right) \quad (7)$$

In (7),  $\hat{\mathbf{f}}_{\mathbf{n}}$  is a numerical flux constant with the matrix  $\mathbf{f}$  vector  $\mathbf{n}$  product  $\mathbf{f}_{\mathbf{n}} := \mathbf{f} \cdot \mathbf{n}$ . The boundary integral is a sum of integrals on the faces of  $K$ , and here for any face of  $K$   $\mathbf{u}^{h,+}$  represents the approximation of  $\mathbf{u}$  on the other side of that face in the case of internal elements, and  $\mathbf{u}_b$  when that face is on  $\partial\Omega$ . Note that to fully comply with (9d), we should have defined for boundary faces  $\mathbf{u}^{h,+} = \mathbf{u}^h$ , and then (7) is rewritten as

$$a(\mathbf{u}^h, \mathbf{v}^h) := \sum_{K \subset \Omega} \left( - \int_K \nabla \mathbf{v}^h \cdot \mathbf{f}(\mathbf{u}^h) d\mathbf{x} + \int_{\partial K} \mathbf{v}^h \hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}^h, \mathbf{u}^{h,+}) d\gamma \right) \quad (8)$$

+  $\sum_{\Gamma \subset \partial\Omega}$  Boundary term for  $\Gamma$ .

---

<sup>1</sup>More precisely such that it is symmetric definite positive up to a symetrization matrix,  $A_0$  that, for fluid problems, is related to the Hessian of the entropy.

- A finite volume scheme. Even though the discontinuous Galerkin method boils down into a finite volume scheme (where the volumes are the mesh elements) for degree 0, let us provide a different example. The notations are defined in Figure 1. Again, we specialize ourselves to the

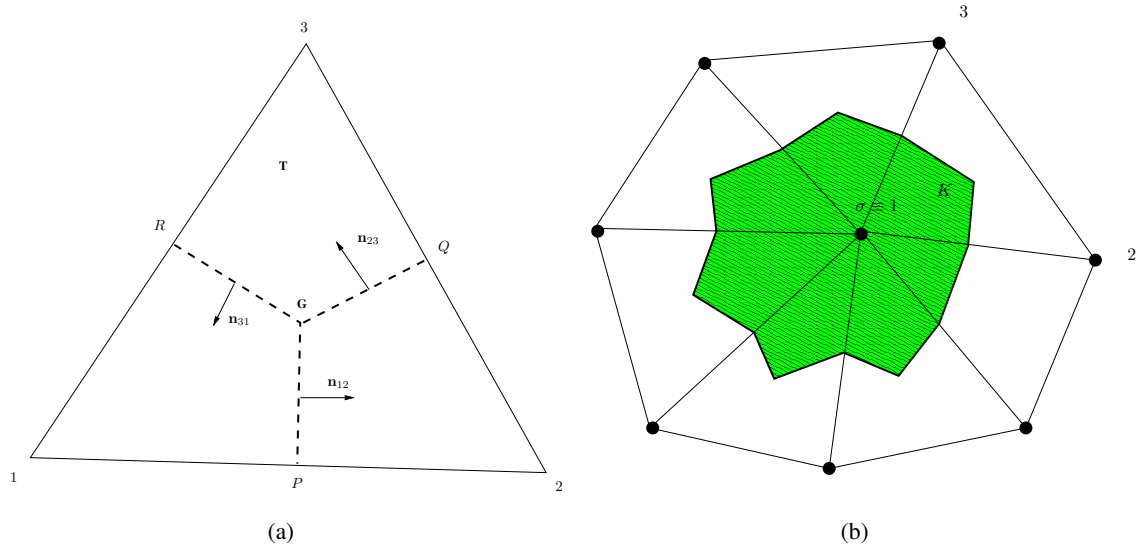


Figure 1: Notations for the finite volume schemes. On the left: definition of the control volume for the degree of freedom  $\sigma$ . The vertex  $\sigma$  plays the role of the vertex 1 on the left picture for the triangle  $K$ . The control volume  $C_\sigma$  associated to  $\sigma = 1$  is green on the right and corresponds to  $1PGR$  on the left. The vectors  $\mathbf{n}_{ij}$  are normal to the internal edges scaled by the corresponding edge length

case of triangular elements, but *exactly the same arguments* can be given for more general elements, provided a conformal approximation space can be constructed. This is the case for triangle elements, and we can take  $k = 1$ .

The control volumes in this case are defined as the median cell, see figure 1 and the scheme is

$$\sum_{\gamma \subset \partial C_\sigma} \hat{\mathbf{f}}_{\mathbf{n}_\gamma}(\mathbf{u}_\sigma, \mathbf{u}^+) = 0.$$

here we have taken a first order finite volume scheme, as it can be seen from the arguments of the numerical flux  $\hat{\mathbf{f}}_{\mathbf{n}}$ , however a high order extension with MUSCL extrapolation can equivalently be considered.

The interesting fact is that all these method can be rewritten in a unified manner, the residual distribution form. In order to integrate the steady version of (1) on a domain  $\Omega \subset \mathbb{R}^d$ , on each element  $K$  and any degree of freedom  $\sigma \in \mathcal{S}$  belonging to  $K$ , we define residuals  $\Phi_\sigma^K(\mathbf{u}^h)$ . Following [11, 14], they are

assumed to satisfy the following conservation relations: For any element  $K$ ,

$$\sum_{\sigma \in K} \Phi_{\sigma}^K(\mathbf{u}^h) = \Phi^K(\mathbf{u}^h) := \int_{\partial K} \hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}^h, \mathbf{u}^{h,+}) d\gamma \quad (9a)$$

where  $\Phi^K(\mathbf{u}^h)$  is often referred to as the ‘‘element residual’’. Note that if we denote by  $\mathbf{f}^h$  the polynomial flux approximation within the element of maximum degree w.r.t. which the quadrature formulas used in practice are exact, we can recast (9a) as

$$\sum_{\sigma \in K} \Phi_{\sigma}^K(\mathbf{u}^h) = \Phi^K(\mathbf{u}^h) = \int_{\partial K} (\hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}^h, \mathbf{u}^{h,+}) - \mathbf{f}_{\mathbf{n}}^h) d\gamma + \int_K \nabla \cdot \mathbf{f}^h d\mathbf{x} \quad (9b)$$

which allows to write the element residual as the integral of the PDE plus a boundary fluctuation.

In the case of a conformal mesh and with continuous elements, the conservation condition becomes

$$\sum_{\sigma \in K} \Phi_{\sigma}^K(\mathbf{u}^h) = \Phi^K(\mathbf{u}^h) = \int_{\partial K} \mathbf{f}_{\mathbf{n}}(\mathbf{u}^h) d\gamma = \int_K \nabla \cdot \mathbf{f}^h d\mathbf{x} \quad (9c)$$

where we recall that the difference between  $\mathbf{f}(\mathbf{u}^h)$  and  $\mathbf{f}^h$  is that the latter is the polynomial approximation within  $K$  of the highest possible degree for which the quadrature used is exact. In general,  $\mathbf{f}(\mathbf{u}^h(\mathbf{x})) \neq \mathbf{f}^h(\mathbf{x})$ .

The discretisation of (1) is achieved via: for any  $\sigma \in \mathcal{S}$ ,

$$\sum_{K \subset \Omega, \sigma \in K} \Phi_{\sigma}^K(\mathbf{u}^h) + \text{Boundary terms} = 0. \quad (9d)$$

Concerning the boundary terms, they can be very naturally embedded by appropriately embedding fluctuations on the boundary faces. We will omit this aspect for simplicity as it is mainly a technical detail.

Using the fact that the basis functions that span  $V_h$  have a *compact* support, then each scheme can be rewritten in the form (9d) with the following expression for the residuals:

- For the SUPG scheme (5), the residual are defined by

$$\begin{aligned} \Phi_{\sigma}^K(\mathbf{u}^h) &= \int_{\partial K} \varphi_{\sigma} \mathbf{f}(\mathbf{u}^h) \cdot \mathbf{n} d\gamma - \int_K \nabla \varphi_{\sigma} \cdot \mathbf{f}(\mathbf{u}^h) d\mathbf{x} \\ &\quad + h_K \int_K \left( \nabla_{\mathbf{u}} \mathbf{f}(\mathbf{u}^h) \cdot \nabla \varphi_{\sigma} \right) \tau_K \left( \nabla_{\mathbf{u}} \mathbf{f}(\mathbf{u}^h) \cdot \nabla \mathbf{u}^h \right) d\mathbf{x}. \end{aligned} \quad (10)$$

Note that in (10) we have made an abuse of language that we will make systematically: to comply with the Gauss theorem and the form of (5), we have written

$$\int_{\Omega} \nabla \varphi_{\sigma} \cdot \mathbf{f}(\mathbf{u}^h) d\mathbf{x}$$

where in the integral we have the ‘product’ of the vector  $\nabla \varphi_{\sigma}$  with the matrix  $\mathbf{f}(\mathbf{u}^h)$ . It has to be understood as

$$\int_{\Omega} \mathbf{f}(\mathbf{u}^h) \cdot \nabla \varphi_{\sigma} d\mathbf{x}.$$



- For the Galerkin scheme with jump stabilization (6), the residuals are defined by:

$$\Phi_{\sigma}^K(\mathbf{u}^h) = \int_{\partial K} \varphi_{\sigma} \mathbf{f}(\mathbf{u}^h) \cdot \mathbf{n} \, d\gamma - \int_K \nabla \varphi_{\sigma} \cdot \mathbf{f}(\mathbf{u}^h) \, d\mathbf{x} + \sum_{\substack{e \text{ faces} \\ \text{of } K}} \frac{\theta_e}{2} h_e^2 \int_e [[\nabla \mathbf{u}^h]] \cdot [[\nabla \varphi_{\sigma}]] \, d\gamma. \quad (11)$$

Here, since the mesh is conformal, any internal edge  $e$  (or face in 3D) is the intersection of the element  $K$  and another element denoted by  $K^+$ .

- For the discontinuous Galerkin scheme,

$$\Phi_{\sigma}^K(\mathbf{u}^h) = - \int_K \nabla \varphi_{\sigma} \cdot \mathbf{f}(\mathbf{u}^h) \, d\mathbf{x} + \int_{\partial K} \varphi_{\sigma} \hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}^h, \mathbf{u}^{h,+}) \, d\gamma. \quad (12)$$

- For the finite volume scheme the fact that the boundary of the control volume is closed implies that the sum of the outward normals vanishes. So, we can define

$$\begin{aligned} \Phi_{\sigma}^K(\mathbf{u}^h) &= \sum_{\gamma \subset (\partial C_{\sigma}) \cap K} (\hat{\mathbf{f}}_{\mathbf{n}_{\gamma}}(\mathbf{u}_{\sigma}, \mathbf{u}^+) - \mathbf{f}(\mathbf{u}_{\sigma}) \cdot \mathbf{n}_{\gamma}) \\ &= \sum_{\gamma \subset \partial(C_{\sigma} \cap K)} \hat{\mathbf{f}}_{\mathbf{n}_{\gamma}}(\mathbf{u}_{\sigma}, \mathbf{u}^+). \end{aligned} \quad (13)$$

We can then use the fact that on each facet separating nodes  $i, j$  local conservation implies  $\hat{\mathbf{f}}_{\mathbf{n}_{ij}}(\mathbf{u}_i, \mathbf{u}_j) + \hat{\mathbf{f}}_{\mathbf{n}_{ji}}(\mathbf{u}_j, \mathbf{u}_i) = 0$ , and recover the elemental conservation relation as follows:

$$\begin{aligned} \sum_{\sigma \in K} \Phi_{\sigma}^K(\mathbf{u}_h) &= \left( \hat{\mathbf{f}}_{\mathbf{n}_{12}}(\mathbf{u}_1, \mathbf{u}_2) - \hat{\mathbf{f}}_{\mathbf{n}_{13}}(\mathbf{u}_1, \mathbf{u}_3) - \mathbf{f}(\mathbf{u}_1) \cdot \mathbf{n}_{12} + \mathbf{f}(\mathbf{u}_1) \cdot \mathbf{n}_{31} \right) \\ &\quad + \left( \hat{\mathbf{f}}_{\mathbf{n}_{23}}(\mathbf{u}_2, \mathbf{u}_3) - \hat{\mathbf{f}}_{\mathbf{n}_{12}}(\mathbf{u}_2, \mathbf{u}_1) + \mathbf{f}(\mathbf{u}_2) \cdot \mathbf{n}_{12} - \mathbf{f}(\mathbf{u}_2) \cdot \mathbf{n}_{23} \right) \\ &\quad + \left( -\hat{\mathbf{f}}_{\mathbf{n}_{23}}(\mathbf{u}_3, \mathbf{u}_2) + \hat{\mathbf{f}}_{\mathbf{n}_{31}}(\mathbf{u}_3, \mathbf{u}_1) - \mathbf{f}(\mathbf{u}_3) \cdot \mathbf{n}_{23} + \mathbf{f}(\mathbf{u}_3) \cdot \mathbf{n}_{31} \right) \\ &= \mathbf{f}(\mathbf{u}_1) \cdot (\mathbf{n}_{12} - \mathbf{n}_{31}) + \mathbf{f}(\mathbf{u}_2) \cdot (-\mathbf{n}_{23} + \mathbf{n}_{31}) + \mathbf{f}(\mathbf{u}_3) \cdot (\mathbf{n}_{31} - \mathbf{n}_{23}) \\ &= \mathbf{f}(\mathbf{u}_1) \cdot \frac{\mathbf{n}_1}{2} + \mathbf{f}(\mathbf{u}_2) \cdot \frac{\mathbf{n}_2}{2} + \mathbf{f}(\mathbf{u}_3) \cdot \frac{\mathbf{n}_3}{2} \end{aligned}$$

where  $\mathbf{n}_j$  is the scaled inward normal of the edge opposite to vertex  $\sigma_j$ , i.e. twice the gradient of the  $\mathbb{P}^1$  basis function  $\varphi_{\sigma_j}$  associated to this degree of freedom. Thus, we can reinterpret the sum as the boundary integral of the Lagrange interpolant of the flux. The finite volume scheme is then a residual distribution scheme with residual defined by (13) and a total residual defined by

$$\Phi^K := \int_{\partial K} \mathbf{f}_{\mathbf{n}}^h \, d\gamma, \quad \mathbf{f}^h = \sum_{\sigma \in K} \mathbf{f}(\mathbf{u}_{\sigma}) \varphi_{\sigma}. \quad (14)$$

- The residual distribution formalism has also been used to build new schemes. A classical example is the nonlinear Lax-Friedrich's discretization built to satisfy both a high order

truncation error estimate of order  $\mathcal{O}(h^{p+1})$  for a polynomial approximation of degree  $p$ , and a positive-coefficient property [30, 13]. The scheme reads

$$\Phi_\sigma^K = \beta_\sigma^K \Phi^K$$

where the coefficients  $\beta_\sigma$  are designed in such a way that the scheme is both monotonicity preserving and formally  $k + 1$ -th order accurate if a polynomial approximation of degree  $k$  is used. This can be achieved in two steps as follows, see [11] for the technical details

1. First evaluate the Rusanov (or Local lax-Friedrich residuals),

$$\Phi_\sigma^{LF,K} = \frac{\Phi^K}{N_K} + \alpha_K (\mathbf{u}_\sigma - \bar{\mathbf{u}}^K)$$

where  $\alpha_K$  is larger than the maximum on  $K$  of  $\|\nabla \mathbf{f}^h\|$ ,  $N_K$  is the number of degree of freedom on  $K$  and  $\bar{\mathbf{u}}^K$  is the arithmetic average of the  $\mathbf{u}_\sigma$  for  $\sigma \in K$ .

2. Define  $x_\sigma$  as the ratio of  $\Phi_\sigma^{LF,K}$  by  $\Phi^K$ , and

$$\beta_\sigma = \frac{\max(x_\sigma, 0)}{\sum_{\sigma' \in K} \max(x_{\sigma'}, 0)}.$$

Since  $\sum_{\sigma \in K} x_\sigma = 1$ , we see that  $\sum_{\sigma' \in K} \max(x_{\sigma'}, 0) \geq 1$ , so that there is no problem in the definition of this quantity as long as  $\Phi^K \neq 0$ . If  $\Phi^K = 0$ , we can take any value since in the end the residual we are going to use is  $\Phi_\sigma^K = \beta_\sigma^K \Phi^K$ .

Other expressions for  $\beta_\sigma^K$  are feasible, but this is the one that is used in practice since it is very simple

This is not enough, as can be found in [11]: the solution appears wiggly, especially in the smooth part of the solution. This is *not* a problem of stability. This occurs because the scheme is over-compressing. One way to overcome this is to add some stabilizing/filtering term. For example, in several papers a streamline upwind/least square term has been added to  $\beta_\sigma^K \Phi^K$ , namely

$$\Phi_\sigma^{K\star} = \beta_\sigma^K \Phi^K + h_K \oint_K [\nabla \mathbf{f}^h(\mathbf{u}^u) \cdot \nabla \varphi_\sigma] \tau_K [\nabla \mathbf{f}^h(\mathbf{u}^u) \cdot \nabla \mathbf{u}^h] d\mathbf{x}$$

The resulting scheme is referred to later on in the paper as LLFs (the first "L" for Limited, the "s" for stabilized). In [8] is discussed the choice of minimal quadrature formula for the evaluation of the integral term. Adding the least square term destroys in principle the maximum preserving property of the method, however in practice it does not, this is why we call this essentially non oscillatory RD scheme: the least square term acts as a mild filter of the spurious modes. Another possible technique to achieve this filtering is to use jump terms as in (11).

In the case of system, one can extend the construction by using a characteristic decomposition of the residual, see again [11]. Last, any monotone first order residual can be used in the step 1 of the construction, not only the Local Lax-Friedrichs one.

All these residuals satisfy the relevant conservation relations, namely (9a), depending if we are dealing with element residuals or boundary residuals.

It can be shown, see [15] that a scheme defined by (9) satisfies a Lax-Wendroff like theorem: if the mesh is regular, if the numerical sequence is bounded in  $L^\infty$  and if a subsequence converges in  $L^2$  (for example) to a  $\mathbf{v} \in L^2$ , then this function is a weak solution of (1). A similar result holds on the entropy if an entropy inequality exists.

## 4 Flux formulation of Residual Distribution schemes

Conversely, we show in this section that any scheme (9d) also admits a flux formulation: *the method is also locally conservative*. In addition, we provide an *explicit* form of the flux. Local conservation is of course well known for the Finite Volume and discontinuous Galerkin approximations. It is much less understood for the continuous finite elements methods, despite the papers [42, 25]. This question of finding an equivalent flux formulation amounts to defining control volumes and flux functions.

We first have to adapt the notion of consistency. We define a multidimensional flux as follows:

**Definition 4.0.1.** *A multidimensional flux  $\hat{\mathbf{f}}_{\mathbf{n}} := \hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}_1, \dots, \mathbf{u}_N)$  is consistent if, when  $\mathbf{u}_1 = \mathbf{u}_2 = \dots = \mathbf{u}_N = \mathbf{u}$  then  $\hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}, \dots, \mathbf{u}) = \mathbf{f}(\mathbf{u}) \cdot \mathbf{n}$ .*

We proceed first with the general case and show the connection with elementary fact about graphs, and then provide several examples. The results of this section apply to any finite element method but also to discontinuous Galerkin methods. There is no need for the exact evaluation of integral formula (surface or boundary), so that these results apply to schemes as they are implemented.

We consider the general case, i.e when  $K$  is a polytope contained in  $\mathbb{R}^d$  with degrees of freedoms on the boundary of  $K$ . The set  $\mathcal{S}$  is the set of degrees of freedom. We construct a triangulation  $\mathcal{T}_K$ , i.e. a graph, of  $K$  whose vertices are exactly the elements of  $\mathcal{S}$ . Choosing an orientation of  $K$ , it is propagated on  $\mathcal{T}_K$ : the edges are oriented.

Inspired by (14), starting from the graph, we will construct control volumes and flux. Again inspired by the finite volume example, and since the shape of the control volumes is still unknown, we label the flux by the edges of the graph, and hence we slightly change notations. The problem is to find quantities  $\hat{\mathbf{f}}_{\sigma, \sigma'}$  for any edge  $[\sigma, \sigma']$  of  $\mathcal{T}_K$  such that:

$$\Phi_{\sigma} = \sum_{\text{edges } [\sigma, \sigma']} \hat{\mathbf{f}}_{\sigma, \sigma'} + \hat{\mathbf{f}}_{\sigma}^b \quad (15a)$$

with, since permuting two vertices amounts to change the orientation,

$$\hat{\mathbf{f}}_{\sigma, \sigma'} = -\hat{\mathbf{f}}_{\sigma', \sigma} \quad (15b)$$

and  $\hat{\mathbf{f}}_{\sigma}^b$  is the 'part' of  $\oint_{\partial K} \hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}^h, \mathbf{u}^{h,+}) d\gamma$  associated to  $\sigma$ . The control volumes will be defined by their normals so that we get consistency. The normal corresponding to the edge  $[\sigma, \sigma']$  is denoted by  $\mathbf{n}_{\sigma, \sigma'}$ .

Note that (15b) implies the conservation relation

$$\sum_{\sigma \in K} \Phi_{\sigma} = \sum_{\sigma \in K} \hat{\mathbf{f}}_{\sigma}^b. \quad (15c)$$

In this paper, we will take

$$\hat{\mathbf{f}}_{\sigma}^b = \oint_{\partial K} \varphi_{\sigma} \hat{\mathbf{f}}_{\mathbf{n}}(\mathbf{u}^h, \mathbf{u}^{h,+}) d\gamma, \quad (15d)$$

but other examples can be considered provided the consistency (15c) relation holds true, see [4]. Any edge  $[\sigma, \sigma']$  is either direct or, if not,  $[\sigma', \sigma]$  is direct. Because of (15b), we only need to know  $\hat{\mathbf{f}}_{\sigma, \sigma'}$  for direct edges. Thus we introduce the notation  $\hat{\mathbf{f}}_{\{\sigma, \sigma'\}}$  for the flux assigned to the direct edge whose extremities are  $\sigma$  and  $\sigma'$ . We can rewrite (15a) as, for any  $\sigma \in \mathcal{S}$ ,

$$\sum_{\sigma' \in \mathcal{S}} \varepsilon_{\sigma, \sigma'} \hat{\mathbf{f}}_{\{\sigma, \sigma'\}} = \Psi_{\sigma} := \Phi_{\sigma} - \hat{\mathbf{f}}_{\sigma}^b, \quad (16)$$

with

$$\varepsilon_{\sigma,\sigma'} = \begin{cases} 0 & \text{if } \sigma \text{ and } \sigma' \text{ are not on the same edge of } \mathcal{T}, \\ 1 & \text{if } [\sigma,\sigma'] \text{ is an edge and } \sigma \rightarrow \sigma' \text{ is direct,} \\ -1 & \text{if } [\sigma,\sigma'] \text{ is an edge and } \sigma' \rightarrow \sigma \text{ is direct.} \end{cases}$$

$\mathcal{E}^+$  represents the set of direct edges.

Hence the problem is to find a vector  $\hat{\mathbf{f}} = (\hat{\mathbf{f}}_{\{\sigma,\sigma'\}})_{\{\sigma,\sigma'\} \text{ direct edges}}$  such that

$$A\hat{\mathbf{f}} = \Psi$$

where  $\Psi = (\Psi_\sigma)_{\sigma \in \mathcal{S}}$  and  $A_{\sigma\sigma'} = \varepsilon_{\sigma,\sigma'}$ .

We have the following lemma which shows the existence of a solution. Its proof can be found in [4].

**Lemma 4.0.1.** *For any couple  $\{\Phi_\sigma\}_{\sigma \in \mathcal{S}}$  and  $\{\hat{\mathbf{f}}_\sigma^b\}_{\sigma \in \mathcal{S}}$  satisfying the condition (15c), there exists numerical flux functions  $\hat{\mathbf{f}}_{\sigma,\sigma'}$  that satisfy (15). Recalling that the matrix of the Laplacian of the graph is  $L = AA^T$ , we have*

1. *The rank of  $L$  is  $|\mathcal{S}| - 1$  and its image is  $(\text{span}\{\mathbf{1}\})^\perp$ . We still denote the inverse of  $L$  on  $(\text{span}\{\mathbf{1}\})^\perp$  by  $L^{-1}$ ,*
2. *With the previous notations, a solution is*

$$(\hat{\mathbf{f}}_{\{\sigma,\sigma'\}})_{\{\sigma,\sigma'\} \text{ direct edges}} = A^T L^{-1}(\Psi_\sigma)_{\sigma \in \mathcal{S}}. \quad (17)$$

This set of flux are consistent and we can estimate the normals  $\mathbf{n}_{\sigma,\sigma'}$ . In the case of a constant state, we have  $\Phi_\sigma = 0$  for all  $\sigma \in K$ . Let us assume that

$$\hat{\mathbf{f}}_\sigma^b = \mathbf{f}(\mathbf{u}^h) \cdot \mathbf{N}_\sigma \quad (18)$$

with  $\sum_{\sigma \in K} \mathbf{N}_\sigma = 0$ : this is the case for all the examples we consider. If  $\hat{\mathbf{f}}_\sigma^b$  is defined by (15d), we see that

$$\mathbf{N}_\sigma = \oint_{\partial K} \varphi_\sigma \mathbf{n} \, d\gamma.$$

The flux  $\mathbf{f}(\mathbf{u}^h)$  has components on the canonical basis of  $\mathbb{R}^d$ :  $\mathbf{f}(\mathbf{u}^h) = (f_1(\mathbf{u}^h), \dots, f_d(\mathbf{u}^h))$ , so that from (18), we get

$$\hat{\mathbf{f}}_\sigma^b = \sum_{i=1}^d f_i(\mathbf{u}^h) \mathbf{N}_\sigma^i.$$

Applying this to  $(\hat{\mathbf{f}}_{\sigma_1}^b, \dots, \hat{\mathbf{f}}_{\sigma_{\#K}}^b)$ , we see that the  $j$ -th component of  $\mathbf{n}_{\sigma,\sigma'}$  for  $[\sigma,\sigma']$  direct, must satisfy:

$$\text{for any } \sigma \in K, \mathbf{N}_\sigma^j = \sum_{[\sigma,\sigma'] \text{ edge}} \varepsilon_{\sigma,\sigma'} \mathbf{n}_{\sigma,\sigma'}^j$$

i.e.

$$(\mathbf{N}_{\sigma_1}^j, \dots, \mathbf{N}_{\sigma_{\#K}}^j)^T = A (\mathbf{n}_{\sigma,\sigma'}^j)_{[\sigma,\sigma'] \in \mathcal{E}^+}.$$

We can solve the system and the solution, with some abuse of language, is

$$(\mathbf{n}_{\sigma,\sigma'})_{[\sigma,\sigma'] \in \mathcal{E}^+} = A^T L^{-1}(\mathbf{N}_{\sigma_1}, \dots, \mathbf{N}_{\sigma_{\#K}})^T \quad (19)$$

This also defines the control volumes since we know their normals. We can state:

**Proposition 4.0.1.** *If the residuals  $(\Phi_\sigma)_{\sigma \in K}$  and the boundary fluxes  $(\hat{\mathbf{f}}_\sigma^b)_{\sigma \in K}$  satisfy (15c), and if the boundary fluxes satisfy the consistency relations (18), then we can find a set of consistent flux  $(\hat{\mathbf{f}}_{\sigma, \sigma'})_{[\sigma, \sigma']}$  satisfying (15). They are given by (17). In addition, for a constant state,*

$$\hat{\mathbf{f}}_{\sigma, \sigma'}(\mathbf{u}^h) = \mathbf{f}(\mathbf{u}^h) \cdot \mathbf{n}_{\sigma, \sigma'}$$

for the normals defined by (19).

We can state a couple of general remarks:

**Remark 4.0.1.**

1. *The flux  $\hat{\mathbf{f}}_{\sigma, \sigma'}$  depend on the  $\Psi_\sigma$  and not directly on the  $\hat{\mathbf{f}}_\sigma^b$ . We can design the fluxes independently of the boundary flux, and their consistency directly comes from the consistency of the boundary fluxes.*
2. *The residuals depends on more than 2 arguments. For stabilized finite element methods, or the non linear stable residual distribution schemes, see e.g. [42, 67, 11], the residuals depend on all the states on  $K$ . Thus the formula (17) shows that the flux depends on more than two states in contrast to the 1D case. In the finite volume case however, the support of the flux function is generally larger than the three states of  $K$ , think for example of an ENO/WENO method, or a simpler MUSCL ones.*
3. *The formula (17) make no assumption on the approximation space  $V^h$ : they are valid for continuous and discontinuous approximations. The structure of the approximation space appears only in the total residual.*

Let us give two examples, that will be valid for SUPG and the Galerkin scheme with stabilisation because the explicit form of the residual does not play any role.

Let  $K$  be a fixed triangle. The degrees of freedom (the vertices) will be denoted by  $\{\sigma\}_{\sigma \in K}$  or  $\{\sigma_i\}_{i=1,2,3}$  or their label in  $\{1, 2, 3\}$ . We are given a set of residues  $\{\Phi_\sigma^K\}_{\sigma \in K}$ , our aim here is to define a flux function such that relations similar to (13) hold true. The adjacency matrix is

$$A = \begin{pmatrix} 1 & 0 & -1 \\ -1 & 1 & 0 \\ 0 & -1 & 1 \end{pmatrix}.$$

A straightforward calculation shows that the matrix  $L = A^T A$  has eigenvalues 0 and 3 with multiplicity 2 with eigenvectors

$$R = \begin{pmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{-1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & 0 & \frac{-2}{\sqrt{6}} \end{pmatrix}$$

To solve  $A\hat{\mathbf{f}} = \Psi$ , we decompose  $\Psi$  on the eigenbasis:

$$\Psi = \alpha_2 R_2 + \alpha_3 R_3$$

where explicitly

$$\alpha_2 = \frac{1}{\sqrt{2}}(\Psi_1 - 2\Psi_2 + \Psi_3)$$

$$\alpha_3 = \sqrt{\frac{3}{2}}(\Psi_1 - \Psi_3)$$

so that

$$\hat{\mathbf{f}}_{12} = \frac{1}{3}(\Psi_1 - \Psi_2), \quad \hat{\mathbf{f}}_{23} = \frac{1}{3}(\Psi_2 - \Psi_3), \quad \hat{\mathbf{f}}_{32} = \frac{1}{3}(\Psi_3 - \Psi_1).$$

In order to describe the control volumes, we first have to make precise the normals  $\mathbf{n}_\sigma$  in that case. It is easy to see that in all the cases described above, we have

$$\mathbf{N}_\sigma = -\frac{\mathbf{n}_\sigma}{2}.$$

Then a short calculation shows that

$$\begin{pmatrix} \mathbf{n}_{12} \\ \mathbf{n}_{23} \\ \mathbf{n}_{31} \end{pmatrix} = \frac{1}{6} \begin{pmatrix} \mathbf{n}_1 - \mathbf{n}_2 \\ \mathbf{n}_2 - \mathbf{n}_3 \\ \mathbf{n}_3 - \mathbf{n}_1 \end{pmatrix}.$$

Using elementary geometry of the triangle, we see that these are the normals of the elements of the dual mesh. For example, the normal  $\mathbf{n}_{12}$  is the normal of  $PG$ , see figure 1.

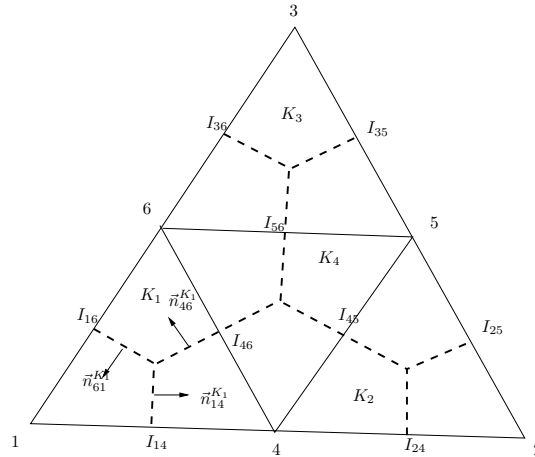


Figure 2: Geometrical elements for the  $\mathbb{P}^2$  case.  $I_{ij}$  is the mid-point between the vertices  $i$  and  $j$ . The intersections of the dotted lines are the centroids of the sub-elements.

In the case of a quadratic approximation a similar set of formula can be given. Again this will be valid

for SUPG and Galerkin with jumps. Using a similar method, we get (see figure 2 for some notations):

$$\begin{aligned}
\hat{\mathbf{f}}_{14} &= \frac{5\Psi_1}{12} - \frac{5\Psi_2}{36} - \frac{\Psi_3}{36} - \frac{7\Psi_4}{36} - \frac{\Psi_5}{12} + \frac{\Psi_6}{36} \\
\hat{\mathbf{f}}_{42} &= \frac{5\Psi_1}{36} - \frac{5\Psi_2}{12} + \frac{\Psi_3}{36} + \frac{7\Psi_4}{36} - \frac{\Psi_5}{36} + \frac{\Psi_6}{12} \\
\hat{\mathbf{f}}_{25} &= -\frac{\Psi_1}{36} + \frac{5\Psi_2}{12} - \frac{5\Psi_3}{36} + \frac{\Psi_4}{36} - \frac{7\Psi_5}{36} - \frac{\Psi_6}{12} \\
\hat{\mathbf{f}}_{53} &= \frac{\Psi_1}{36} + \frac{5\Psi_2}{36} - \frac{5\Psi_3}{12} + \frac{\Psi_4}{12} + \frac{7\Psi_5}{36} - \frac{\Psi_6}{36} \\
\hat{\mathbf{f}}_{36} &= -\frac{5\Psi_1}{36} - \frac{\Psi_2}{36} + \frac{5\Psi_3}{12} - \frac{\Psi_4}{12} + \frac{\Psi_5}{36} - \frac{7\Psi_6}{36} \\
\hat{\mathbf{f}}_{61} &= -\frac{5\Psi_1}{12} + \frac{\Psi_2}{36} + \frac{5\Psi_3}{36} - \frac{\Psi_4}{36} + \frac{\Psi_5}{12} + \frac{7\Psi_6}{36} \\
\hat{\mathbf{f}}_{64} &= \frac{\Psi_1}{9} - \frac{\Psi_3}{9} + \frac{2\Psi_4}{9} - \frac{2\Psi_5}{9} \\
\hat{\mathbf{f}}_{45} &= -\frac{\Psi_1}{9} + \frac{\Psi_2}{9} + \frac{2\Psi_5}{9} - \frac{2\Psi_6}{9} \\
\hat{\mathbf{f}}_{56} &= -\frac{\Psi_2}{9} + \frac{\Psi_3}{9} - \frac{2\Psi_4}{9} + \frac{2\Psi_6}{9}
\end{aligned}$$

Then we choose the boundary flux:

$$\hat{\mathbf{f}}_{\sigma}^b = \int_{\partial K} \varphi_{\sigma} \mathbf{f}(\mathbf{u}^h) \cdot \mathbf{n} \, d\gamma$$

and get:

$$\mathbf{N}_l = -\frac{\mathbf{n}_l}{6} \quad \text{if } l = 1, 2, 3$$

$$\mathbf{N}_4 = \frac{\mathbf{n}_3}{3} \quad \mathbf{N}_5 = \frac{\mathbf{n}_1}{3} \quad \mathbf{N}_6 = \frac{\mathbf{n}_2}{3}$$

The normals are given by:

$$\begin{aligned}
\mathbf{n}_{14} &= \frac{5\mathbf{N}_1}{12} - \frac{5\mathbf{N}_2}{36} - \frac{\mathbf{N}_3}{36} - \frac{7\mathbf{N}_4}{36} - \frac{\mathbf{N}_5}{12} + \frac{\mathbf{N}_6}{36} \\
\mathbf{n}_{42} &= \frac{5\mathbf{N}_1}{36} - \frac{5\mathbf{N}_2}{12} + \frac{\mathbf{N}_3}{36} + \frac{7\mathbf{N}_4}{36} - \frac{\mathbf{N}_5}{36} + \frac{\mathbf{N}_6}{12} \\
\mathbf{n}_{25} &= -\frac{\mathbf{N}_1}{36} + \frac{5\mathbf{N}_2}{12} - \frac{5\mathbf{N}_3}{36} + \frac{\mathbf{N}_4}{36} - \frac{7\mathbf{N}_5}{36} - \frac{\mathbf{N}_6}{12} \\
\mathbf{n}_{53} &= \frac{\mathbf{N}_1}{36} + \frac{5\mathbf{N}_2}{36} - \frac{5\mathbf{N}_3}{12} + \frac{\mathbf{N}_4}{12} + \frac{7\mathbf{N}_5}{36} - \frac{\mathbf{N}_6}{36} \\
\mathbf{n}_{36} &= -\frac{5\mathbf{N}_1}{36} - \frac{\mathbf{N}_2}{36} + \frac{5\mathbf{N}_3}{12} - \frac{\mathbf{N}_4}{12} + \frac{\mathbf{N}_5}{36} - \frac{7\mathbf{N}_6}{36} \\
\mathbf{n}_{61} &= -\frac{5\mathbf{N}_1}{12} + \frac{\mathbf{N}_2}{36} + \frac{5\mathbf{N}_3}{36} - \frac{\mathbf{N}_4}{36} + \frac{\mathbf{N}_5}{12} + \frac{7\mathbf{N}_6}{36} \\
\mathbf{n}_{64} &= \frac{\mathbf{N}_1}{9} - \frac{\mathbf{N}_3}{9} + \frac{2\mathbf{N}_4}{9} - \frac{2\mathbf{N}_5}{9} \\
\mathbf{n}_{45} &= -\frac{\mathbf{N}_1}{9} + \frac{\mathbf{N}_2}{9} + \frac{2\mathbf{N}_5}{9} - \frac{2\mathbf{N}_6}{9} \\
\mathbf{n}_{56} &= -\frac{\mathbf{N}_2}{9} + \frac{\mathbf{N}_3}{9} - \frac{2\mathbf{N}_4}{9} + \frac{2\mathbf{N}_6}{9}
\end{aligned}$$

The case of the discontinuous Galerkin can be worked out similarly.

## 5 Embedding source terms: well balancing and global fluxes

### 5.1 The one-dimensional case

We look now into the approximation of solutions to the steady limit of (1) in one dimension :

$$\frac{\partial \mathbf{f}(\mathbf{u})}{\partial x} = S(\mathbf{u}, x).$$

Despite the apparent simplicity of the problem, it is well known that some fundamental change of paradigm is required compared to conservation laws. In particular, the non-autonomous character of the problem, associated to the presence of the source term  $S(\mathbf{u}, x)$  requires a more general notion of consistency.

The examples provided in the introduction for the shallow water equations show that these non-trivial states can only in some cases be characterized by a set of physically relevant invariants. A possible way out to replace the notion of *consistency with constant states* is to introduce an (unknown) ‘‘source flux’’  $\mathbf{s}$  as

$$\mathbf{s}(x) = \int_{x_0}^x S(\mathbf{u}(x), x) dx$$

One can now argue that a more relevant notion of steady states is the one associated to a constant global flux

$$\mathbf{g} = \mathbf{f} - \mathbf{s} = \mathbf{g}_0 = \text{const}$$



Although several works have proposed explicit constructions of the local values of  $\mathbf{s}$  [31, 32, 27], this is essentially possible only in a 1D setting or by means of some dimension by dimension splitting. The main issue is how to construct schemes consistent with the notion of a constant global flux, without necessarily having its explicit knowledge.

This issue is dealt with very naturally in the residual distribution setting. Let us focus for the moment on continuous approximations on conformal meshes. The natural way to proceed is to generalize the notion of conservation defined by (9c) by including the whole PDE in it:

$$\sum_{\sigma \in K} \Phi_{\sigma}^K(\mathbf{u}^h) = \Phi^K(\mathbf{u}^h) = \int_K (\nabla \cdot \mathbf{f}^h - S^h) dx \quad (20)$$

where  $S^h$  is a discrete approximation of the source within the element compatible with a certain quadrature strategy. We will discuss this in aspect in some detail shortly. For the moment, let us consider the 1D residual distribution scheme, seeking the steady solution as the limit of

$$\Delta x_{\sigma} \frac{\mathbf{u}_{\sigma}^{n+1} - \mathbf{u}_{\sigma}^n}{\Delta t} + \sum_{K, \sigma \in K} \Phi_{\sigma}^K(\mathbf{u}^{h,n}) = 0$$

We focus on the  $\mathbb{P}^1$  case to begin with, but the extension to higher polynomials can be obtained similarly to what we discussed in section §4. Each node will receive two contributions,  $\Phi_{\sigma}^{\sigma-1/2}$  from the element on its left,  $K_{\sigma-1/2}$ , the second,  $\Phi_{\sigma}^{\sigma+1/2}$  from the element on its right,  $K_{\sigma+1/2}$ . The conservation relation (20) in one dimension and the  $\mathbb{P}^1$  setting simply writes

$$\Phi_{\sigma}^{\sigma+1/2} + \Phi_{\sigma+1}^{\sigma+1/2} = \mathbf{f}(\mathbf{u}_{\sigma+1}) - \mathbf{f}(\mathbf{u}_{\sigma}) - \int_{K_{\sigma+1/2}} S^h(x, \mathbf{u}^h) dx.$$

We have again adapted the notations to make them less heavy in the 1D case. We can proceed as follows. First we set

$$S_{\sigma+1/2} := \frac{1}{|K_{\sigma+1/2}|} \int_{K_{\sigma+1/2}} S^h(x, \mathbf{u}^h) dx = \frac{1}{\Delta_{\sigma+1/2} x} \int_{x_{\sigma}}^{x_{\sigma+1}} S^h(x, \mathbf{u}^h) dx.$$

Next we define

$$\mathbf{s}_{\sigma+1} = \mathbf{s}_{\sigma} + \Delta_{\sigma+1/2} x S_{\sigma+1/2}(x, \mathbf{u}^h)$$

with  $\mathbf{s}_{\sigma_0} = 0$  for a given but arbitrary  $\sigma_0$ . We then set  $\mathbf{g}_{\sigma+1} = \mathbf{f}_{\sigma+1} - \mathbf{s}_{\sigma+1}$ , and recast the iterations as

$$\Delta_{\sigma} x \frac{\mathbf{u}_{\sigma}^{n+1} - \mathbf{u}_{\sigma}^n}{\Delta t} + (\mathbf{g}_{\sigma} + \Phi_{\sigma}^{K_{\sigma+1/2}}) - (\mathbf{g}_{\sigma} - \Phi_{\sigma}^{K_{\sigma-1/2}}) = 0$$

Finally we set

$$\hat{\mathbf{g}}_{\sigma \pm 1/2} := \mathbf{g}_{\sigma} \pm \Phi_{\sigma}^{\sigma \pm 1/2}$$

which is a consistent numerical global flux. Note however that  $\hat{\mathbf{g}}_{\sigma \pm 1/2}$  is never explicitly built in the residual distribution approach!

To give a few examples, let us start from the centered splitting

$$\Phi_{\sigma}^{\sigma \pm 1/2} = \frac{1}{2} \Phi^{\sigma \pm 1/2} = \frac{1}{2} (\Delta \mathbf{f}_{\sigma \pm 1/2} - \Delta_{\sigma+1/2} x S_{\sigma+1/2})$$

We can easily show that this splitting leads to an equivalent global finite volume flux

$$\hat{\mathbf{g}}_{\sigma\pm 1/2} := \mathbf{g}_\sigma \pm \frac{1}{2} \Phi^{K_{\sigma\pm 1/2}} = \mathbf{f}_\sigma - \mathbf{s}_\sigma + \frac{\Delta \mathbf{f}_{\sigma\pm 1/2}}{2} - \frac{\Delta \mathbf{s}_{\sigma\pm 1/2}}{2} = \frac{\mathbf{g}_\sigma + \mathbf{g}_{\sigma\pm 1}}{2}$$

Similarly, the Galerkin scheme can be shown to be equivalent to the finite volume scheme with global numerical flux given by

$$\hat{\mathbf{g}}_{\sigma\pm 1/2}^{\text{Gal}} = \frac{\mathbf{g}_\sigma + \mathbf{g}_{\sigma\pm 1}}{2} \pm \int_{K_{\sigma\pm 1/2}} \left( \varphi_\sigma - \frac{1}{2} \right) \left( \frac{\partial \mathbf{f}^h}{\partial x} - S^h \right) dx ,$$

and for the SUPG we have

$$\hat{\mathbf{g}}_{\sigma\pm 1/2}^{\text{SUPG}} = \frac{\mathbf{g}_\sigma + \mathbf{g}_{\sigma\pm 1}}{2} \pm \int_{K_{\sigma\pm 1/2}} \left[ \left( \varphi_\sigma - \frac{1}{2} \right) + \nabla_{\mathbf{u}} \mathbf{f} \frac{\partial \varphi_\sigma}{\partial x} \tau_{K_{\sigma\pm 1/2}} \right] \left( \frac{\partial \mathbf{f}^h}{\partial x} - S^h \right) dx ,$$

In general, we can follow for example section 4, consider test functions  $\{\omega_\sigma\}$  defining a partition of unity for conservation purposes, and set

$$\Phi_\sigma^K = \int_K \omega_\sigma \left( \nabla \cdot \mathbf{f}^h - S^h \right) dx , \quad (21)$$

This scheme is equivalent in 1D to the finite volume global flux method defined by

$$\hat{\mathbf{g}}_{\sigma\pm 1/2}^{\text{RD}} = \frac{\mathbf{g}_\sigma + \mathbf{g}_{\sigma\pm 1}}{2} \pm \int_{K_{\sigma\pm 1/2}} \left( \omega_\sigma - \frac{1}{2} \right) \left( \frac{\partial \mathbf{f}^h}{\partial x} - S^h \right) dx , \quad (22)$$

In one space dimension, all these schemes are compatible with the discrete steady state

$$\frac{\partial \mathbf{f}^h}{\partial x} = S^h \iff \forall \sigma \quad \mathbf{g}_\sigma(\mathbf{u}, S(\mathbf{u}, x)) = \mathbf{g}_0 = \text{const} \quad (23)$$

which is here the only relevant consistency condition.

*Second order at steady state.* It is important to remark the following: the residual distribution numerical flux (22) is a compact consistent flux (in the sense of (23)) which takes as inputs *unreconstructed* states:

$$\hat{\mathbf{g}}_{\sigma\pm 1/2}^{\text{RD}} = \hat{\mathbf{g}}_{\sigma\pm 1/2}^{\text{RD}}(\mathbf{u}_\sigma, \mathbf{u}_{\sigma\pm 1}; x_\sigma, x_{\sigma\pm 1})$$

Despite of this fact, the residual formulation provides a framework to design the flux in a way guaranteeing at least second-order truncation at steady state, without any gradient reconstruction. This can be shown for steady balance laws following e.g. [30] by estimating the truncation error defined as (see also [56], appendix B)

$$\begin{aligned} \epsilon &:= \left\| \sum_\sigma v(x_\sigma) \sum_{K, \sigma \in K} \Phi_\sigma^K(\mathbf{w}_{\text{ex}}^h) \right\| \\ &= \left\| \int_\Omega v^h \left( \frac{\partial \mathbf{f}_{\text{ex}}^h}{\partial x} - S_{\text{ex}}^h \right) + \sum_K \sum_{\sigma, \sigma' \in K} \frac{v(x_\sigma) - v(x_{\sigma'})}{2} \int_K (\omega_\sigma - \varphi_\sigma) \left( \frac{\partial \mathbf{f}_{\text{ex}}^h}{\partial x} - S_{\text{ex}}^h \right) \right\| \end{aligned}$$

with  $v(x)$  any smooth compactly supported test function, with  $\mathbf{w}_{\text{ex}}$  a regular enough steady solution,  $\Phi_{\sigma}^K(\mathbf{w}_{\text{ex}}^h)$  the residual distribution (21) evaluated when nodally replacing the numerical solution with samples of the exact one. The analysis shows that the main design rules for second order fluxes of the form (22) are the boundedness of  $\omega_{\sigma}$  and the formal second order of the spatial approximations of the flux  $\mathbf{f}^h$ , and of the source  $S^h$ , which are readily obtained by means of e.g. linear interpolation between two neighbouring states.

The most classical particular case is the upwind fluctuation splitting of Roe [62], obtained in the  $P^1$  case by setting

$$\omega_{\sigma}|_{\sigma_{\pm 1/2}} := \frac{1 \mp \text{sign}(\widetilde{\nabla_{\mathbf{u}}\mathbf{f}}_{\sigma_{\pm 1/2}})}{2}$$

where the sign of a matrix is defined as usual via its eigen-decomposition, and where following [62]  $\widetilde{\nabla_{\mathbf{u}}\mathbf{f}}$  denotes the exact linearization of the flux Jacobian verifying the conservation condition

$$\widetilde{\nabla_{\mathbf{u}}\mathbf{f}}_{\sigma_{\pm 1/2}} \Delta \mathbf{u}_{\sigma_{\pm 1/2}} = \Delta \mathbf{f}_{\sigma_{\pm 1/2}}.$$

Note that in 1D this linearization establishes a direct link between the cell conservation relation (9c) and the linearized non-conservative form of the PDE. This allows to mention another known particular case, when the initial differential problem contains non-conservative terms

$$\frac{\partial \mathbf{f}(\mathbf{u})}{\partial x} + B(\mathbf{u}) \frac{\partial \mathbf{u}}{\partial x} = S(\mathbf{u}, x).$$

In this case, one cannot simple apply the definition of conservation according to the principles introduced so far. In the residual distribution setting this can be handled by embedding the non-conservative term in the cell residual, so that (20) becomes in 1D

$$\sum_{\sigma \in K} \Phi_{\sigma}^K(\mathbf{u}^h) = \Phi^K(\mathbf{u}^h) = \int_K \left( \frac{\partial \mathbf{f}^h}{\partial x} - S^h + B(\mathbf{u}^h) \frac{\partial \mathbf{u}^h}{\partial x} \right) dx$$

The approximation of the last term has been reduced in the residual distribution setting simply to a quadrature problem for a given (linear) variation of  $\mathbf{u}^h$  (see e.g. [66, 60] §9.5, and [70, 69]). This is exactly what is done in path-conservative finite volume (see [26] and references therein), when the path chosen to connect the left and right states at a cell interface is linear. In particular, if  $A = \nabla_{\mathbf{u}}\mathbf{f} + B$ , path conservative finite volumes are equivalent to the residual scheme obtained with

$$\omega_{\sigma}|_{\sigma_{\pm 1/2}} := \frac{1 \mp \text{sign}(A_{\sigma_{\pm 1/2}})}{2}, \quad \Phi^K = \Delta \mathbf{f}_K - |K| S_K + B_K \Delta \mathbf{u}_K$$

where  $B_K$  can be evaluated for example with a one point quadrature over the element. The approximation and quadrature choices made above to evaluate

$$(B^h \frac{\partial \mathbf{u}^h}{\partial x})(x) = B(\mathbf{u}^h(x)) \frac{\partial \mathbf{u}^h(x)}{\partial x}$$

correspond to the choice of the path in the finite volume context. Assuming  $\mathbf{u}^h(x)$  to linearly join two states is one possibility. One could also have  $\mathbf{u}^h(x) = \mathbf{u}(\mathbf{v}^h(x), f(x))$  with  $\mathbf{v} = \mathbf{v}(\mathbf{u})$  some array of physical states (assumed to be a  $C^1$  invertible function of  $\mathbf{u}$  and to vary linearly), and  $f$  a given field. Many other

choices are possible. Note that this does not solve the issues raised by the non-conservative nature of the system, namely the fact that the classical characterization of weak solutions and the Lax Wendroff theorem cannot be applied. This leaves all the uncertainties on the right form for a numerical scheme, see [10] for a counter example. However, we also remark that the RD framework can help in correcting schemes that discretise a non conservative form of a system in conservation form to account for certain constraints: see [6] for an example involving multiphase flows.

## 5.2 Multiple dimensions, beyond second order, and other extensions

The discussion provided allows to systematically design, by means of a residual based approach, well balanced fluxes with a genuine second order truncation without the need of any reconstruction. We consider here several extensions, with focus on the multidimensional steady case:

$$\nabla \cdot \mathbf{f}(\mathbf{u}) + \mathbf{B}(\mathbf{u}) \cdot \nabla \mathbf{u} = S(\mathbf{u}, \mathbf{x}), \quad (24)$$

although we will not dwell too much on the issues related to the the presence of the non-conservative term for the reasons stated above.

The main recipe behind the method considered is already contained in equation (21). As in the 1D case, without loss of generality we will assume that the discrete unknowns are obtained as the steady limit of the pseudo-time iteration

$$|C_\sigma| \frac{\mathbf{u}_\sigma^{n+1} - \mathbf{u}_\sigma^n}{\Delta t} + \sum_{K, \sigma \in K} \Phi_\sigma^K(\mathbf{u}^{h,n}) \quad (25)$$

with the conservation/consistency constraint that

$$\begin{aligned} \Phi^K &:= \sum_{\sigma \in K} \Phi_\sigma^K(\mathbf{u}^h) = \oint_{\partial K} \mathbf{f}_n^h d\gamma - \int_K S^h d\mathbf{x} + \int_K (\mathbf{B} \cdot \nabla \mathbf{u})^h d\mathbf{x} \\ &= \int_K \{ \nabla \cdot \mathbf{f}^h - S^h + (\mathbf{B} \cdot \nabla \mathbf{u})^h \} d\mathbf{x} \end{aligned} \quad (26)$$

where as before  $\mathbf{f}^h$  is the polynomial flux approximation of the highest degree for which the quadrature employed is exact, while both  $S^h$  and  $(\mathbf{B} \cdot \nabla \mathbf{u})^h$  are appropriately defined continuous approximations of the source and non-conservative terms, consistent with the quadrature strategy adopted.

*Global fluxes.* Without specifying the form of  $\Phi_\sigma^K$  we could repeat the construction of section §4 and abstractly provide definitions of local fluxes embedding all the terms of the PDE. Differently from the 1D case however, in multiple dimensions the presence of the source term  $S$ , makes it quite unclear how to define consistency in a genuinely multidimensional setting.

Concerning the non-conservative term, the choice of the approximation/quadrature for the term  $(\mathbf{B} \cdot \nabla \mathbf{u})^h$  can be seen as choosing the manifold along which solutions can evolve. In this sense one could speak of *manifold-conservative* approach. As for path-conservative schemes, the authors remain skeptical as to how much specifying this notion would allow to side-step the fact that the classical definition of weak solution does not apply here. As in 1D, several choices are possible, the most obvious being here to take

$$(\mathbf{B} \cdot \nabla \mathbf{u})^h = \mathbf{B}(\mathbf{u}^h) \cdot \nabla \mathbf{u}^h$$

and evaluate the integral of this term by means of some quadrature formula. In the remainder of the paper, we will omit this term as none of the examples considered contain it.

*Consistency for general smooth steady solutions.* The examples provided in section §3 can all be cast as a particular case of the general prototype

$$\Phi_{\sigma}^K(\mathbf{u}^h) = \int_K \omega_{\sigma} \{ \nabla \cdot \mathbf{f}^h - S^h \} d\mathbf{x} + \oint_{\partial K} [\mathcal{L}(\varphi_{\sigma})] \cdot [\tau_{\mathcal{L}} \mathcal{L}(\mathbf{u}^h)] d\gamma \quad (27)$$

with  $\mathcal{L}(\cdot)$  some linear differential operator. The error analysis recalled in section §5.1 can be used in this more general setting. In particular, given a smooth exact solution  $\mathbf{w}$  we define

$$\begin{aligned} \epsilon(\mathbf{w}^h) := & \left\| \sum_{\sigma} \sum_{K, \sigma' \in K} v(\mathbf{x}_{\sigma}) \Phi_{\sigma}^K(\mathbf{w}^h) \right\| = \left\| \int_{\Omega} v^h \{ \nabla \cdot \mathbf{f}^h - S^h \} d\mathbf{x} \right. \\ & + \sum_K \sum_{\sigma, \sigma' \in K} \frac{v(\mathbf{x}_{\sigma}) - v(\mathbf{x}_{\sigma'})}{N_K} \int_K (\omega_{\sigma} - \varphi_{\sigma}) \{ \nabla \cdot \mathbf{f}^h - S^h \} d\mathbf{x} \\ & \left. + \sum_K \sum_{\sigma, \sigma' \in K} \frac{v(\mathbf{x}_{\sigma}) - v(\mathbf{x}_{\sigma'})}{N_K} \oint_{\partial K} [\mathcal{L}(\varphi_{\sigma})] \cdot [\tau_{\mathcal{L}} \mathcal{L}(\mathbf{u}^h)] d\gamma \right\| \end{aligned} \quad (28)$$

Simple approximation arguments can be used to show that [13] the above prototype has a consistency of order  $O(h^{p+1})$  as soon as the underlying polynomial approximation is of degree  $p$ , and provided that  $\omega_{\sigma}$  is uniformly bounded (w.r.t. solution, mesh size, and problem data), and that the  $\tau_{\mathcal{L}}$  scales appropriately. For  $\mathcal{L} = \nabla$ , the appropriate scaling is  $\tau_{\mathcal{L}} = O(h^2)$  as in the Galerkin with jump stabilization (11). This generalizes the compact second order construction discussed in the previous section to the multidimensional case, and to higher degree approximations. The estimate essentially allows to recover the underlying finite element approximation error. One can however do more if some knowledge of the exact solution is embedded in this approximation.

*Super-consistency exact preservation of steady invariants.* Interesting results can be shown when the source term depends on some given data, say a given field  $f(\mathbf{x})$  as for example the bathymetry in the shallow water equations, or some geometrical parametrization when considering the solution of the differential problem on a manifold (see e.g. [64] and references therein). We are in particular interested in exact steady solutions characterized by the existence of a set of invariants  $\mathbf{v} = \mathbf{v}(\mathbf{u}, f)$  constant throughout the spatial domain. Several examples have been provided in the introduction. Assuming a sufficient smoothness of  $f$ , of the solution, and of the mapping  $(\mathbf{v}, f) \mapsto \mathbf{u}(\mathbf{v}, f)$ , we can write

$$\nabla \cdot \mathbf{f}(\mathbf{u}) = (\nabla_{\mathbf{u}} \mathbf{f} \nabla_{\mathbf{v}} \mathbf{u}) \cdot \nabla \mathbf{v} + (\nabla_{\mathbf{u}} \mathbf{f} \nabla_f \mathbf{u}) \cdot \nabla f = \nabla_{\mathbf{v}} \mathbf{f} \cdot \nabla \mathbf{v} + \nabla_f \mathbf{f} \cdot \nabla f$$

Solutions characterized by the invariance relation  $\mathbf{v} = \mathbf{v}_0 = \text{const.} \forall \mathbf{x}$ , satisfy

$$\nabla_f \mathbf{f}(\mathbf{v}_0, f) \cdot \nabla f + S(\mathbf{v}_0, f) = 0. \quad (29)$$

This shows that for these solutions the flux and source dependence on the data, and the approximation and quadrature of the latter will play a crucial role. For smooth/simple enough problems, the above relation can be reproduced quite accurately in the residual context. An interesting result can be obtained by analyzing the error (28) when the approximation is written directly for the steady invariants, and thus  $\epsilon(\mathbf{w}^h) = \epsilon(\mathbf{v}^h, f^h) = \epsilon(\mathbf{v}_0, f^h)$ . For schemes of the form (27) with approximation/quadrature choices consistent with exactness for  $\mathbf{v}$  constant, the following is shown in [53, 54].

**Proposition 5.0.1** (Steady invariants and superconsistency). *Under standard regularity assumptions on the mesh, provided the test function  $\omega_\sigma$  in (27) is uniformly bounded w.r.t.  $h$ ,  $\mathbf{u}_h$ , element residuals, data of the problem, and provided  $\tau_{\mathcal{L}}$  is  $O(h^{2d_{\mathcal{L}}})$  with  $d_{\mathcal{L}}$  the highest derivative order of the operator  $\mathcal{L}$ , then:*

- for exact integration scheme (25)-(27) with  $\mathbf{u}^h = \mathbf{u}(\mathbf{v}^h, f)$ ,  $\mathbf{f}^h = \mathbf{f}^h(\mathbf{v}^h, f)$ , and  $S^h = S^h(\mathbf{v}^h, f)$  preserves exactly the equilibrium (29);
- for approximate integration, assuming that a flux quadrature exact for approximate polynomial fluxes of degree  $p_f$  is used, and a source quadrature exact for approximate polynomial sources of degree  $p_v$ , and assuming that  $f \in H^{p+1}$  with  $\nabla f \in H^p$ , and  $p > \min(p_f, p_v)$ , then the scheme is superconsistent w.r.t. (29), and in particular, its consistency is of order  $r = \min(p_f + 2, p_v + 3)$ .

Independently of the details, the meaning of this result is that if the field  $f$  and its derivatives can be approximated by a smooth enough function given analytically, and if the approximation is done in terms of steady invariants instead of conserved variables, than the consistency of the scheme is determined by the quadrature strategy and it is in particular independent on the order of the underlying approximation

A few remarks are in order. The numerical results will provide examples indicating that numerical convergence w.r.t. the order of the quadrature formulas is indeed observed in practice, at least for simple cases. However, exact preservation is possible for some important and physically relevant examples. We can mention at least two for the shallow water equations :

1. Lake at rest state (2). Exact preservation has been guaranteed by choosing the same approximation for  $h$  and  $b$ , and performing the quadrature of the hydrostatic terms  $\omega_{sigma}(\nabla(gh^2/2) + gh\nabla b)$  either exactly, or using the chain rule  $\omega_{sigma}(\nabla(gh^2/2) + gh\nabla b) = \omega_{sigma}gh\nabla\eta$  [56, 57];
2. Constant slope equilibrium (4). This is a case compatible with constant  $\mathbf{u}$ , for which any consistent quadrature becomes exact. However, exact preservation is guaranteed only due to the residual formulation in which the different source terms are simultaneously integrated [54].

Another remark concerns the smoothness of  $f$ . The proposition above is built upon estimates involving approximation estimates, and related quadrature error formulas over elements. This suggests that one can construct higher order approximations with less quadrature points either by means of a clever mesh generation step, embedding regions containing jumps in  $f$  or in its derivatives as mesh edges/points, or by means of an adaptive quadrature strategy, avoiding the use of quadrature formulas across such discontinuities.

## 6 Time dependent problems

### 6.1 Preliminaries: global fluxes, time derivative and mass matrices

To fix some basic concepts, we start by the simplest problem: the 1D advection equation

$$\frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} = 0$$

The most classical discretization we can apply is the upwind scheme

$$\Delta_{\sigma x} \frac{du_\sigma}{dt} = -(\hat{f}_{\sigma+1/2} - \hat{f}_{\sigma-1/2}), \quad \hat{f}_{\sigma+1/2} = \begin{cases} f_\sigma & \text{if } a > 0 \\ f_{\sigma+1} & \text{if } a < 0 \end{cases} .$$

This scheme is known to be only first order in space, and typically high order approximations are obtained by replacing the values of  $f_\sigma = f(u_\sigma)$  and  $u_\sigma$  by evaluations of appropriately reconstructed polynomials on either side of the interfaces  $\sigma \pm 1/2$ .

Let us make another experiment instead: we set  $S := -\partial u / \partial t$ , and apply an upwind global flux method. We can proceed as in section §5.1, and formally define a global flux consistent with

$$g = f - \int_{x_0}^x S$$

We can now for example define upwind fluxes as

$$\hat{g}_{\sigma+1/2} = \begin{cases} g_\sigma & \text{if } a > 0 \\ g_{\sigma+1} & \text{if } a < 0 \end{cases},$$

having set

$$g_{\sigma+1} = g_\sigma + \Delta f_{\sigma+1/2} - \int_{x_\sigma}^{x_{\sigma+1}} S \approx g_\sigma + \Delta f_{\sigma+1/2} + \frac{|K_{\sigma+1/2}|}{2} \left( \frac{\partial u_\sigma}{\partial t} + \frac{\partial u_{\sigma+1}}{\partial t} \right).$$

The resulting scheme reads

$$\hat{g}_{\sigma+1/2} - \hat{g}_{\sigma-1/2} = 0,$$

or equivalently, using the definition of the numerical flux, and rearranging terms:

$$\frac{g_{\sigma+1} - g_{\sigma-1}}{2} + \frac{\text{sgn}(a)}{2} (g_\sigma - g_{\sigma-1}) - \frac{\text{sgn}(a)}{2} (g_{\sigma+1} - g_\sigma) = 0$$

The definition of the global flux given above leads, after some manipulations, to the following semi-discrete evolution scheme

$$\Delta_\sigma x \frac{d\hat{u}_\sigma}{dt} + \hat{f}_{\sigma+1/2} - \hat{f}_{\sigma-1/2} = 0 \quad (30)$$

where  $\hat{f}_{\sigma\pm 1/2}$  are exactly those of the first order upwind scheme, while the nodal approximation of the time derivative is now defined as

$$\begin{aligned} \Delta_\sigma x \frac{d\hat{u}_\sigma}{dt} := & \frac{|K_{\sigma-1/2}|}{4} \left( \frac{\partial u_{\sigma-1}}{\partial t} + \frac{\partial u_\sigma}{\partial t} \right) + \frac{|K_{\sigma+1/2}|}{4} \left( \frac{\partial u_{\sigma+1}}{\partial t} + \frac{\partial u_\sigma}{\partial t} \right) \\ & + \frac{\text{sgn}(a)}{2} \frac{|K_{\sigma+1/2}|}{2} \left( \frac{\partial u_{\sigma+1}}{\partial t} + \frac{\partial u_\sigma}{\partial t} \right) \\ & - \frac{\text{sgn}(a)}{2} \frac{|K_{\sigma-1/2}|}{2} \left( \frac{\partial u_{\sigma-1}}{\partial t} + \frac{\partial u_\sigma}{\partial t} \right) \end{aligned} \quad (31)$$

Quite interestingly, this method can be checked (e.g. with a truncated Taylor series analysis) to have a second order truncation error in space without the need of any polynomial reconstruction. This is not related to error compensation on a uniform mesh, but to the improved balance of the different terms for linear data within each cell. This simple example shows how the notion of a global flux can be applied to other types of terms in the PDE. In the case of the time derivative, the global flux approach leads to the appearance of a mass matrix.

As we have shown previously, there is a direct between the upwind finite volume method and residual based schemes which can be summarized into the equality

$$\hat{f}_{\sigma+1/2} - \hat{f}_{\sigma-1/2} = \sum_{K, \sigma \in K} \int_K \omega_\sigma \frac{\partial f^h}{\partial x} dx$$

with  $f^h$  piecewise linear. There are at least two definitions of the test function  $\omega_\sigma$  which give back the upwind scheme, namely

$$\omega_{\sigma|K_{\sigma\pm 1/2}} = \varphi_\sigma + a \frac{\partial \varphi_\sigma}{\partial x} \tau_{\sigma+1/2}, \quad \tau_{\sigma+1/2} = \frac{|K_{\sigma+1/2}|}{2|a|}$$

and

$$\omega_{\sigma|K_{\sigma\pm 1/2}} = \frac{1 \mp \text{sgn}(a)}{2}$$

with  $\varphi_\sigma$  the linear finite element test functions. The method (30)-(31) can be obtained in a much more natural and elegant way as a particular case of a residual method, and in particular of the one corresponding to the second definition of  $\omega_\sigma$  above. The first definition, provides an even better variant with a truncation error which improves to an order  $\Delta x^3$  (see e.g. [59]) for uniform meshes ! The benefit of this idea is to allow high order of accuracy with the most compact stencil. Its drawback is that it requires inverting the mass matrix. This analogy has also been used in other context to generate compact high order finite difference schemes associated to a variational form [45].

The next sections discuss how to generalize this idea to multiple dimensions and, more importantly, how the issue of inverting the mass matrix has been side-stepped.

## 6.2 Generalization

As the last section has shown, following a finite element strategy for (1) for the unsteady case will always lead to a formulation of the form

$$M(\mathbf{u}^{n+1} - \mathbf{u}^n) + \Delta t \delta \mathbf{F} = \mathcal{S},$$

where  $M$  is a mass matrix,  $\delta \mathbf{F}$  contains all the spatial approximation terms, and  $\mathcal{S}$  the approximation of the source term. It is possible, depending on the formulation, that several instances of  $\mathbf{u}$  appear. One of the biggest problem is the mass matrix.

Things are different for the classical formulations of finite volume and discontinuous Galerkin schemes, which lead to diagonal or block diagonal matrices because of the locality of the approximation of  $\mathbf{u}$ . These are small (however dense) matrices which can be inverted locally on each element, and are independent on the mesh connectivity. This is probably one of the keys of the success of these methods, especially for genuinely hyperbolic and evolutionary problems. In the case of continuous approximation, the story is not as simple. For example, the SUPG method will lead to a mass matrix that may evolve in time. This is also the case of the RD schemes developed in [9, 1, 55, 58, 71]. The Galerkin method with jump stabilisation does not have this problem, but nevertheless, we still have a sparse positive definite matrix to invert. One of the strategies followed in the past has been to work on highly implicit variants of the schemes, trying to cover this computational overhead with the possibility of using large time steps. Unfortunately, despite the excellent results, the schemes obtained in this way are relatively cumbersome to code [12, 34, 58, 39].



Moreover, the advantage of using large time steps, very useful for viscous flows and problems with large stiffness, is less obvious for wave propagation problems, even on non-uniform meshes [39, 65, 40].

In [55] is explained how to approximate the solution in time, but without having to invert a mass matrix. In this reference, the method is explained for piecewise linear element (and triangular element). The method was further extended to any order (and any type of simplex) in [3].

In practice, for the steady version of (1), each of the known schemes can be written using test functions. This is clear for the SUPG scheme, where the test functions are defined in each element and are possibly discontinuous across element. Please note that in the non linear case, the test functions will depend on  $\mathbf{u}$ . The same is true for the schemes of [11, 14], except that the scheme will be non linear even for a linear problem in order to enforce non oscillatory constraints. In the case of Galerkin method with jump, one can also reinterpret the method in this way, thanks to the use of a lifting operator allowing to embed the jump terms in a numerical flux. Hence, in all cases, we write

$$\Phi_{\sigma}^K(\mathbf{u}^h) = \int_K \omega_{\sigma} \operatorname{div} \mathbf{f}(\mathbf{u}^h) d\mathbf{x}$$

where  $\omega_{\sigma}$  is the test function associated to the element  $K$  and the degree of freedom  $\sigma$ . For example, for the SUPG method, this is

$$\omega_{\sigma} = \varphi_{\sigma} + h_K \nabla_{\mathbf{u}} \mathbf{f}(\mathbf{u}^h) \cdot \nabla \varphi_{\sigma} \tau_K.$$

Integrating (1), and using, for simplicity of exposure, the mid point rule in time, will lead to

$$\int_{\Omega} \omega_{\sigma} (\mathbf{u}^{n+1} - \mathbf{u}^n) + \frac{\Delta t}{2} \left( \int_{\Omega} \omega_{\sigma} \operatorname{div} \mathbf{f}(\mathbf{u}^{n+1}) d\mathbf{x} + \int_{\Omega} \omega_{\sigma} \operatorname{div} \mathbf{f}(\mathbf{u}^n) d\mathbf{x} \right) = 0, \quad (32)$$

that despite its complexity, we still can rewrite in a form similar to (9) with:

$$\Phi_{\sigma}^K(\mathbf{u}, \mathbf{v}) = \int_K \omega_{\sigma} (\mathbf{u} - \mathbf{v}) + \frac{\Delta t}{2} \left( \int_K \omega_{\sigma} \operatorname{div} \mathbf{f}(\mathbf{u}) d\mathbf{x} + \int_K \omega_{\sigma} \operatorname{div} \mathbf{f}(\mathbf{v}) d\mathbf{x} \right)$$

and  $\mathbf{u} = \mathbf{u}^{n+1}$ ,  $\mathbf{v} = \mathbf{u}^n$ .

What makes this complex is the term in time. If we consider the simpler set of residual, for  $|C_{\sigma}| > 0$  to be defined,

$$\psi_{\sigma}^K(\mathbf{u}, \mathbf{v}) = |C_{\sigma}^K| (\mathbf{u}_{\sigma} - \mathbf{v}_{\sigma}) + \frac{\Delta t}{2} \left( \int_K \omega_{\sigma} \operatorname{div} \mathbf{f}(\mathbf{u}) d\mathbf{x} + \int_K \omega_{\sigma} \operatorname{div} \mathbf{f}(\mathbf{v}) d\mathbf{x} \right),$$

then the scheme defined from this is easily solvable: we know  $\mathbf{u}^n$ , we can get  $\mathbf{u}^{n+1}$  explicitly. The key question are (i) how to define the lumping parameter  $|C_{\sigma}^K|$ , and (ii) how to combine the schemes defined by this two set of residual in order to get an approximation the solution  $\mathbf{u}^{n+1}$  given by (32) with the *same* accuracy.

In [55], in the case of a  $\mathbb{P}^1$  approximation and triangle element, it is shown that if  $|C_{\sigma}^K| = \frac{|K|}{3}$  and if

we define  $\mathbf{u}^{n+1}$  using a predictor corrector algorithm as:

$$\begin{aligned}
|C_\sigma|(\mathbf{u}_\sigma^{(1)} - \mathbf{u}_\sigma^n) &= -\Delta t \int_\Omega \omega_\sigma \operatorname{div} \mathbf{f}(\mathbf{u}^n) d\mathbf{x} \\
|C_\sigma|(\mathbf{u}_\sigma^{(2)} - \mathbf{u}_\sigma^{(1)}) &= -\frac{\Delta t}{2} \left( \int_\Omega \omega_\sigma \operatorname{div} \mathbf{f}(\mathbf{u}^{(1)}) d\mathbf{x} + \int_\omega \omega_\sigma \operatorname{div} \mathbf{f}(\mathbf{u}^n) d\mathbf{x} \right) \\
&\quad - \int_\Omega \omega_\sigma (\mathbf{u}^{(1)} - \mathbf{u}^n) \\
\mathbf{u}^{n+1} &= \mathbf{u}^{(2)}
\end{aligned} \tag{33}$$

then we have a second order scheme in time, with similar stability properties as the original steady scheme.

The extension to higher than second order and general simplex has been done in [3]. The main idea is to notice that (33) can be reinterpreted as a defect correction method: if one wants to solve  $L^{(2)}(U) = 0$ , and if one has a second operator, called  $L^{(1)}$ , such that in some norm,

$$\|(L^{(1)}(U) - L^{(2)}(U)) - (L^{(1)}(V) - L^{(2)}(V))\| \leq \Delta \|U - V\|, \tag{34a}$$

if in addition  $L^{(1)}$  satisfies a coercivity relation,

$$\alpha \|U - V\| \leq \|L^{(1)}(U) - L^{(1)}(V)\|, \tag{34b}$$

and finally, if  $L^{(2)}(U) = 0$  has a unique solution  $U^\star$ , then the solution  $U^{(p)}$  of the iterative scheme:

$$\begin{aligned}
&U^{(0)} \text{ given} \\
\text{do for } p \geq 0 \quad &L^{(1)}(U^{(p+1)}) = L^{(1)}(U^{(p)}) - L^{(2)}(U^{(p)}).
\end{aligned} \tag{35}$$

Of course the question is to know a good stopping criteria. The answer is given by the following: it can be shown [3] that

$$\|U^{(p)} - U^\star\| \leq \left(\frac{\Delta}{\alpha}\right)^p \|U^{(0)} - U^\star\|. \tag{36}$$

Here, if the coefficients  $|C_\sigma^K|$  are chosen such that

$$\sum_{\sigma \in K} |C_\sigma^K| \mathbf{u}_\sigma = \int_K \mathbf{u}(x) d\mathbf{x},$$

i.e

$$|C_\sigma^K| = \int_K \varphi_\sigma d\mathbf{x},$$

then the conditions (34) are met and then a CFL-like condition  $\frac{\Delta}{\alpha} \approx \Delta t$ , we can interpret (36), after  $p$  iterations, as

$$\|U^{(p)} - U^\star\| \approx C \Delta t^p.$$

This means that if the un-lumped formulation is of order  $p$ , then the lumped with the algorithm (35) one will provide a solution with the same accuracy after  $p$  iteration only. In the case of piecewise linear elements,  $L^2$  is defined from the residuals  $\Phi_\sigma^K$  and  $L^1$  is defined from  $\Psi_\sigma^K$ .

The problem is that often  $\int_K \varphi_\sigma d\mathbf{x}$  is not positive: this is the case for quadratic Lagrange interpolant in triangles. A possible remedy to this is to use basis functions that are positive as for example Bézier polynomials [3]. Note that the linear polynomials are also Bézier polynomials of degree 1.

In practice, we split the time interval  $[t_n, t_{n+1}]$  with  $p$  sub-time steps  $t_n = t_{p,0} < t_{p,1} = t_n + \alpha_1 \Delta t < \dots < t_{p,p-1} = t_n + \alpha_{p-1} \Delta t < t_{p,p} = t_{n+1} = t_n + \Delta t$ , the vector  $\mathbf{u}$  contains the approximations of  $\mathbf{u}$  for the sub-time steps, i.e.  $\mathbf{u} = (\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_p)$  with  $\mathbf{u}_j \approx \mathbf{u}(\cdot, t_{p,j})$ . Then we write

1. Set  $\mathbf{u}^{(0)} = (\mathbf{u}^n, \mathbf{u}^n, \dots, \mathbf{u}^n)$ : we initialise the vector with the state at time  $t_n$ ,
2. Do for  $l = 1, \dots, p$ , do for  $k = 1, \dots, p$

$$|C_\sigma|(\mathbf{u}_k^{(l+1)} - \mathbf{u}_k^{(l)}) + \underbrace{\sum_{K, \sigma \in K} \left( \int_K \omega_\sigma (\mathbf{u}_k^{(l)} - \mathbf{u}^n) d\mathbf{x} \right)}_{(I)} + \underbrace{\oint_{t_n}^{t_{p,k}} \left( \int_K \omega_\sigma \operatorname{div} \mathbf{f}(\mathbf{u}^{(p)}) d\mathbf{x} \right)}_{(II)} = 0$$

(I) is evaluated by spatial quadratures and (II) by space-time quadratures: we write

$$\oint_{t_n}^{t_{p,k}} \left( \int_K \omega_\sigma \operatorname{div} \mathbf{f}(\mathbf{u}) d\mathbf{x} \right) = \Delta t \sum_{r=0}^p \theta_r^k \int_K \omega_\sigma \operatorname{div} \mathbf{f}(\mathbf{u}_r^{(l)}) d\mathbf{x}$$

where  $\theta_r^k$  is the integral over  $[0, \alpha_k]$  of the (time) Lagrange interpolant at the points  $\{\alpha_0 = 0, \alpha_1, \dots, \alpha_{p-1}, \alpha_p = 1\}$ . and  $|C_\sigma| = \sum_{K, \sigma \in K} |C_\sigma^K|$ .

3.  $\mathbf{u}^{n+1} = \mathbf{u}^{(p+1)}$

The integer  $p$  is equal to the expected order of accuracy. The procedure is explicit.

Using the results of section 4, it is easy to see that the sum of (I) and (II) is a sum of flux: If the quadrature formula in time needs  $p$  steps, then we extrude the element  $K$  adding  $p$  layers, this induces natural graph that is connected, and hence the discussion of section 4 can be repeated.

### 6.3 Unsteady problems and well balanced on dynamic meshes

To complete the presentation of the time dependent case, we go back to balance laws. The schemes developed in these pages have been initially designed having in mind general adaptive meshes. In the time dependent case, it is thus natural to envision some dynamic adaptation method. The literature is filled with promising adaptive techniques, see e.g. [46, 35, 18] and references therein for a (non-comprehensive) review. We want underline here an important aspect allowing to generalize some of the concepts introduced in section §5. In particular, when dynamic meshes are employed a fundamental step is the operator allowing to map the solution from one mesh to another. In steady computations the error possibly introduced during the projection from one mesh to another may be lost at convergence, provided the boundary conditions are not affected by this aspect. In time dependent simulations however, the remap may pollute the evolution of the solution, and affect both its accuracy and stability, as much as the underlying discretization method.

Design criteria for the remap are thus consistency, conservation, monotonicity preservation, etc: exactly the same criteria applied when solving the main PDE problem ! For balance laws, if the source term depends on external data, these must also be projected, and a well balanced condition may be on the shopping list of desired properties. To fix ideas, we will consider here projection methods based on

some kind of Lagrangian, or rather Arbitrary Lagrangian Eulerian (ALE) remap (see e.g. [52, 18] and references therein). To start with, we recast (1) with  $S = S(\mathbf{u}, f(\mathbf{x}))$  in ALE form:

$$\left. \frac{\partial(J\mathbf{u})}{\partial t} \right|_{\mathbf{x}} + J \operatorname{div}(\mathbf{f}(\mathbf{u}) - \mathbf{w}\mathbf{u}) = J S(\mathbf{u}, f(\mathbf{x})) \quad (37)$$

with the additional definitions/relations

$$\begin{aligned} \left. \frac{d\mathbf{x}(t)}{dt} \right|_{\mathbf{x}} &= \mathbf{w}, \quad \mathbf{x}(0) = \mathbf{X} \\ \left. \frac{\partial J}{\partial t} \right|_{\mathbf{x}} - J \operatorname{div} \mathbf{w} &= 0 \\ \left. \frac{\partial f}{\partial t} \right|_{\mathbf{x}} - \mathbf{w} \cdot \nabla f &= 0 \end{aligned} \quad (38)$$

with  $J$  the determinant of the Jacobian of the mapping  $M : \mathbf{X} \mapsto \mathbf{x}$ , namely

$$J := \det\left(\frac{\partial \mathbf{x}}{\partial \mathbf{X}}\right).$$

Note that the numerical discretization essentially provides a discrete equivalent of (37). The relations (38), which are true and exact on the continuous level, are not explicitly solved numerically but must be seen as constraints to embed as much as possible in the discretization. As in [52, 18] we assume that mesh operations can be represented by some continuous deformation operator, so that the projection from one mesh to the other boils down somehow to mimic (37), with constraints (38). This allows to update the list of design criteria for the schemes:

- *Discrete Geometric Conservation.* This is essentially a discrete analog of the second relation in (38) and represents the conservation of volume along the mapping  $M$ ;
- *Mass Conservation.* Without loss of generality, we can assume that the first relation in (37) is homogenous (so  $S_1 = 0$ ) and represents mass conservation:

$$\left. \frac{\partial(J\rho)}{\partial t} \right|_{\mathbf{x}} + J \nabla \cdot (\mathbf{v}\rho - \mathbf{w}\rho) = 0.$$

Note that one of most classical translations of geometric conservation is to check that uniform flows are preserved by the discretization [68]. This corresponds to the fact that for  $\rho$  and  $\mathbf{v}$  constant the last mass conservation equation becomes the second in (38), and is also equivalent to the standard notion of consistency with respect to constant states;

- *Well balancedness* As already remarked in section §5, consistency for constant states is not necessarily applicable to balance laws, as only non-constant, data dependent, steady states are admissible. So well balanced is in contradiction with some of the above properties, and notably conservation;
- *ALE remap.* The last relation in (38) represents the ALE time derivative of the data  $f$ , which is also often written in conservative form by combining it with geometric conservation:

$$\left. \frac{\partial(Jf)}{\partial t} \right|_{\mathbf{x}} - J \operatorname{div}(\mathbf{w}f) = 0.$$

The ALE time derivative is essentially an advection operator. Its approximation poses very similar questions of consistency, accuracy, stability, and bounded variations for the data, as the approximation of (37). Moreover, for problems in which the data play a key role (e.g. topography inundation and coastal risk assessment), the deterioration of such data due to the approximation error may introduce an unacceptable uncertainty in the predictions. So the use of standard techniques to approximate this advection problem, as e.g. proposed in [72] to guarantee both well balancedness and mass conservation, may be very delicate. For example, to avoid diverging from reality the last reference proposes to periodically re-initialize the data, which implies losing the consistency with the ALE projection which may cost well balancedness, or mass conservation.

Ideally, all of the above properties should be satisfied. However, as remarked well balanced is in general in contradiction with e.g. mass conservation, and the ALE projection may be at odds with the preservation of the accuracy of the data involved. We consider here a (physically relevant) example to better highlight this issue, and show a possible solution in the context of residual distribution.

**Example: Lake at rest solutions on moving meshes.** We consider the shallow water equations in ALE form with an arbitrary (non constant) bathymetry  $b(\mathbf{x})$ . As discussed above, the classical characterization of the DGCL based on the preservation of the state  $\mathbf{u} = \mathbf{u}_0 = \text{const}$  cannot be used here, as constant states are not solutions of the problem due to the presence of the source. To better fix ideas, we will recast (38) as follows

$$\underbrace{J \left( \frac{\partial \mathbf{u}}{\partial t} \Big|_X - \mathbf{w} \cdot \nabla \mathbf{u} \right)}_{H_1} + \underbrace{\mathbf{u} \left( \frac{\partial J}{\partial t} \Big|_X - J \nabla \cdot \mathbf{w} \right)}_{H_2} + \underbrace{J \left( \nabla \cdot \mathbf{f} - S \right)}_{H_3} = 0$$

As amply discussed in section §5, we have a general framework to devise well-balanced Eulerian discretization methods do embed integral (or even local) versions of  $H_3 = 0$ . Previous work has shown how to extend this framework to an ALE setting for both explicit and implicit time integration [20, 33, 40, 49], embedding discretely the constraint of the geometric conservation law  $H_2 = 0$ . Unfortunately by their nature Eulerian methods are unable to embed the condition  $H_1 = 0$ , which will be polluted, also in correspondence of steady exact solutions which would be exactly represented on fixed meshes.

A possible way out of this limitation for solutions admitting a set of steady invariants  $\mathbf{v}$ , is that the ALE formulation, and thus the mesh projection, should be performed using  $\mathbf{v}$  as main variable. To explain we consider the lake at rest state, but other cases can be treated in a similar way. In this case, we can set  $\mathbf{v} = [\eta, h\mathbf{v}] = [H + b, H\mathbf{v}]$ , and steady states are characterized by  $\mathbf{v} = \mathbf{v}_0 = [\eta_0, 0]$ , and thus  $h = h(\mathbf{x}) = \eta_0 - b(\mathbf{x})$ . In the continuous case we can invoke the fact that the bathymetry satisfies the ALE remap (last in (38)), namely

$$\underbrace{\frac{\partial b}{\partial t} \Big|_X - \mathbf{w} \cdot \nabla b}_{H_4=0} = 0$$

This can be used to modify the ALE formulation and write it directly in terms of  $\mathbf{v}$ :

$$\frac{\partial (J\mathbf{v})}{\partial t} \Big|_X + J \nabla \cdot (\mathbf{f}(\mathbf{v}, b) - \mathbf{w}\mathbf{v}) + JS = 0 \quad (39)$$

This formulation is equivalent to

$$J \underbrace{\left( \frac{\partial \mathbf{v}}{\partial t} \Big|_X - \mathbf{w} \cdot \nabla \mathbf{v} \right)}_{H_1+H_4} + \mathbf{v} \underbrace{\left( \frac{\partial J}{\partial t} \Big|_X - J \nabla \cdot \mathbf{w} \right)}_{H_2} + J \underbrace{\left( \nabla \cdot \mathbf{f} + S \right)}_{H_3} = 0$$

We can now use any Eulerian scheme which is well balanced and compatible with geometric conservation, and we will be able to ensure that all the terms in the above sumamtion will be zero if  $\mathbf{v}$  is constant.

For completeness, we recall that the formulation (39), which is referred to in [18] as to the well balanced form of the equations, is similar to the pre-balanced form of the Shallow Water equations of [63] which uses a modified definition of the flux and source terms (cf. [18] for details).

*The problem of mass conservation.* We now consider the additional constraint of achieving discrete conservation of the total water mass in the domain. We integrate in space and in time the mass conservation equation in well-balanced form (39)

$$\int_{\Omega(t)} \eta(\mathbf{x}(t), t) d\mathbf{x} - \int_{\Omega_X} \eta(X, 0) d\mathbf{x} + \int_0^t \int_{\partial\Omega(t)} (H\mathbf{v} - \eta\mathbf{w}) \cdot \mathbf{n} ds dt = 0$$

Let  $V(t) = \int_{\Omega(t)} H d\mathbf{x}$  be the total volume of water in the domain at time  $t$ , and define  $B(t) = \int_{\Omega(t)} b d\mathbf{x}$ . We can rewrite the above conservation statement as

$$H(t) - H(0) + \int_0^t \int_{\partial\Omega(t)} H(\mathbf{v} - \mathbf{w}) \cdot \mathbf{n} ds dt = - \left( B(t) - B(0) - \int_0^t \int_{\partial\Omega(t)} b\mathbf{w} \cdot \mathbf{n} ds dt \right) \quad (40)$$

which states that, modulo the boundary conditions, we have conservation over the full domain provided that we satisfy geometry consrvation and the bathymetry satisfies the ALE remap given by the last in (38), namely if

$$B(t) - B(0) - \int_0^t \int_{\partial\Omega(t)} b\mathbf{w} \cdot \mathbf{n} ds dt = 0$$

As already said, some work in literature propose indeed to evolve the bathymetry according to the ALE remap (see e.g. [72]). As discussed earlier, the uncertainty on the topography associated to the error introduced by this approach may not be acceptable in many applications (e.g. assessment of coastal risks). Combining the ALE remap with some periodic reinitialization of the data would end up breaking mass conservation unless a more clever fix is sought. A possible one has been suggested in [18, 19], and is recalled hereafter.

Assume for simplicity that the domain boundaries are not moving, or that  $\mathbf{w} \cdot \mathbf{n}$  is verified. We can write the mass error at time  $t$  as

$$E_{mass} = H(t) - H(0) + \int_0^t \int_{\partial\Omega} H\mathbf{v} \cdot \mathbf{n} ds dt = B(0) - B(t)$$

We now remark that the two quantities on the right hand side are in principle equal, as they are both approximations of the integral of  $b(\mathbf{x})$  over the domain. If the domain boundaries are not moving, this quantity should remain constant in time. In practice however, these two integrals will be evaluated on a moving mesh. This means that, even if both the domain of integration and the data being integrated are constant, the quadrature points used will move, so the result will not be the same. To be more precise, the

evaluation of  $B(t)$  will be expressed by a sum which will depend on the time update of the scheme. For example, for the explicit approach discussed in section §6.2, we will have

$$B(t) = \sum_{\sigma} b_i(t) |C_{\sigma}(t)| \quad (41)$$

with  $b_i = b(\mathbf{x}_i(t))$ . The idea proposed in [18, 19] is to replace  $b(\mathbf{x}_i(t))$  in the last expression by some mapped value, allowing to minimize the overall mass error, and exploiting as much as possible the actual bathymetric data. In particular, one way to achieve this is to set

$$\tilde{b}_{\sigma} := \frac{1}{|C_{\sigma}|} \int_{C_i(t)} b(\mathbf{x}(t)) dx \approx \sum_{f=1}^{N_q} \omega_q b(\mathbf{x}_q(t)) \quad (42)$$

where the right hand side defines a high order accurate quadrature formula of the real (initial/reference) data over the current cell. The increase in accuracy of the quadrature within each moving cell, allows to compensate for the movement of the cells themselves. In practice, for most problems quadrature formulas exact for degree 2 polynomials are enough to keep the mass error to machine zero levels.

Please refer to [18, 19] to for the extension to problems with dry areas and on curvilinear coordinates.

## 7 Examples

### 7.1 Some examples of compressible flows simulations

We present two cases: the first one is the well known DMR test case by Colella and Woodward: it is the interaction of a Mach 10 shock wave in a quiescent media with wedge angle of  $30^\circ$ . The result is displayed in figure 3 for a cubic (Bézier approximation) and the quality of results is comparable to what can be found in the litterature for a similar resolution (estimated as  $100^2$  for a Cartesian mesh). The second one can be

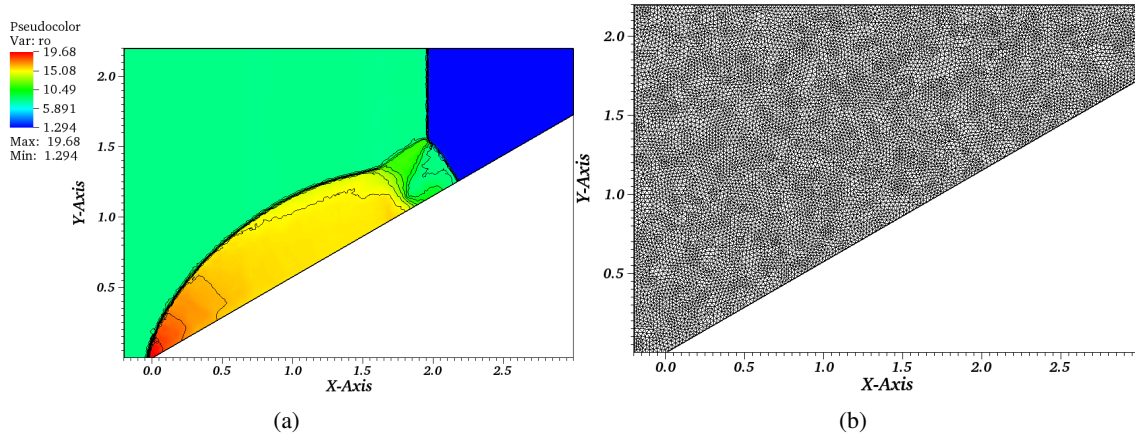


Figure 3: DMR at  $T = 0.2$ , solution on a triangular mesh with 19248 elements.

seen as a 2D version of the Shu-Osher case where we have the interaction between a shock wave and a density wave. The conditions are thus: at  $\mathbf{x}$  in the disc of center  $(0, 0)$  and radius 6,

$$\begin{pmatrix} \rho \\ \vec{u} \\ p \end{pmatrix} = \begin{cases} (3.857143, 2.69369 \frac{\mathbf{x}}{\|\mathbf{x}\|}, 10.333333)^T & \text{if } \|\mathbf{x}\| \leq 1 \\ (1 + 2 \sin(5\|\mathbf{x}\|), \vec{0}, 1)^T & \text{if } 1 < \|\mathbf{x}\| \leq 4 \\ (1 + 2 \sin(5 \times 4), \vec{0}, 1)^T & \text{if else.} \end{cases}$$

The mesh, the initial solution, an intermediate solution and the final one at  $t = 1.8$  are displayed on figure 4. They are obtained with a third order accurate time scheme and quadratic Bézier approximation. It can be observed that the sine wave is not damped, though the resolution of the mesh is not very fine (all degrees of freedom are represented, there are 8985 dofs). The schemes are a combination of the Galerkin scheme with jump stabilisation and the LLFs with quadratic/cubic approximation, see [7] for more details.

## 7.2 Shallow water and the lake at rest state state

We comment on some results of the shallow water equations originally appeared in [57, 18]. The test considered is a quite classical a perturbation of the lake at rest state, with a bathymetry defined by a smooth exponential hump. We refer to the original papers, and references therein, for details. The scheme used is the non-linear limited Lax-Friedrich's residual distribution, described in section §3, with appropriate modifications of the mass matrix and stabilization operators to handle both smooth and discontinuous flows, while accounting for well balanced and wet/dry transitions (see [54] for the explicit method on fixed meshes). The same polynomial approximation is used for the conservative variables  $\mathbf{u}^h$ , and for the topography  $b^h$ . The quadrature strategy is exact on lake at rest solutions (cf. section §5). The main difference between the two references is that in [57] no special care is taken in handling the time derivative, and a fully implicit (in time) approach is used, based on a second order trapezoidal method. On the contrary, in [18] the authors have combined the error correction method discussed in this paper, with a well balanced ALE formulation on moving adaptive meshes. The high order discrete remapping of the initial topographic data recalled in section §5 is used to preserve mass conservation on moving meshes, within the errors of the quadrature formulas used in the remap.

We report in the top and mid-rows of figure 5 the contours of the free surface level obtained respectively with the implicit scheme (top), and with the explicit adaptive approach (middle). The results allow to compare the evolution of the initial perturbation, in both cases unspoiled by any unphysical oscillations in the free stream relation, which is the main interest of constructing well balanced schemes. The bottom row shows the moving adaptive meshes produced with technique proposed in [18] which follow very closely the wave pattern.

Despite of the fact that the contours of the implicit method are slightly crisper than those of the explicit one, the water levels obtained are very similar. This is confirmed by the cuts along the centreline, reported on figure 6. The plots in this figure show that the adaptive computations on a relatively coarse mesh (roughly 12k nodes) compare very favourably in terms of the peaks and troughs of the free surface with those of the implicit scheme on a finer grid (roughly 20k nodes), and with a reference solution of obtained with the explicit scheme on a fine mesh (roughly 50k nodes). The water levels on the unadapted coarse mesh are also reported to show the substantial benefit of adaptation.

On the left in figure 7 we report the total mass error for different quadrature strategies in the bathymetric remap. For this smooth case the error levels are already small with second order quadrature ( $r$  in the figure denotes the degree of exactness of the integration formulas). For higher order formulas, the error remains



practically at machine zero. Finally, the table on the right reports CPU times for the explicit computations, which show a computational gain of almost 35% in time for the adaptive method compared to the fine mesh results. This figure could be further improved with some local or partitioned time marching method. To conclude on this aspect, we refer to [54] for similar comparisons between the explicit and implicit residual schemes on fixed meshes, shown that for this type of problems, for a given resolution the explicit approach can be 3 to 6 time faster.

### 7.3 Shallow water and moving equilibria

We consider here three applications involving moving equilibria in shallow water flows with constant total energy (3). The first two applications involve a very classical situations with a  $C^0$  bathymetry defined as

$$b(x) = \begin{cases} \frac{1}{5} \left( 1 - \frac{(x-10)^2}{4} \right) & \text{if } x \in [8, 12] \\ 0 & \text{otherwise} \end{cases} .$$

The prescribed values of total energy and mass flux are :

$$[\mathcal{E}_0, q_0] = [22.06605m^2/s^2, 4.42m^2/s].$$

The spatial domain has a horizontal length of  $25m$ .

The first test, see figure 8, consists of perturbing the 1D steady state corresponding to the above choices within the slice  $x \in [6.5, 7.5]$ . The perturbation is added to the free surface, and has a magnitude of  $0.05m$ . We use a structured triangulation containing with spacing  $25/200$ , and periodic boundary conditions in the vertical direction. The scheme used is the implicit LLFs scheme of [57], with approximation done in terms of the steady invariants  $\mathbf{v} = [\mathcal{E}, h\mathbf{v}]$ , and with a piecewise linear approximation of the bathymetry. We refer to [51, 54] for the practical implementation of this choice, which requires non-linear iterations to locally invert the mapping  $(\mathbf{v}, b) \mapsto \mathbf{u}$ .

The results are in excellent agreement, at least qualitatively with similar results of the litterature, using dG or WENO schemes.

We then consider a similar test, but on the 2D domain  $[0, 25]^2$  (see [54] concerning the relevance of (3) in 2D). We add a 2D perturbation to the same one-dimensional steady state obtained with the choices described above. As before, a perturbation of  $0.05m$  is added to the free surface  $\eta$ , only this time in the subdomain  $[6.5, 7.5] \times [12, 13]$ . We report on figure 9 snapshots of the perturbation  $\delta\eta := \eta - \eta_{\text{steady}}(x)$ , with  $\eta_{\text{steady}}(x)$  the free surface level corresponding to the exact 1D steady solution. We compare the results obtained on a structured triangulation with spacing  $25/200$  with the straightforward application of the LLFs scheme, denoted by LLFs( $\mathbf{u}^h$ ) exactly well balanced only for the lake at rest state), with the same scheme based on the approximation of the steady invariants (denoted by LLFs( $\mathbf{v}^h$ )). In both cases the standard linear approximation of the bathymetry  $b^h$  is used.

Also in in 2D the LLFs( $\mathbf{v}^h$ ) provides a perfect evolution of the perturbation with no visible spurious effects. The LLFs( $\mathbf{u}^h$ ) scheme shows such perturbations. However, we stress very strongly that on such meshes these are of the order of  $10^{-4}m$ , thus several orders of magnitude smaller than the initial one added to the free surface. Of course, these would become relevant had we reduced the magnitude of the initial perturbation. This is a very nice result for obtained with the ‘‘straight out of the box’’ application of the

residual distribution method. This kind of result is not obtained by a straightforward application of finite volumes or dG schemes.

Finally, we consider a genuinely 2D configuration obtained by replacing the structured triangulation with an unstructured one. The bathymetry is now defined by a series of  $C^1$  sinusoidal ribs (see [54] for details). We compare the LLFs( $\mathbf{v}^h$ ) scheme with approximation in steady invariants, and analytical data  $b(x)$  and  $b'(x)$  used in the residual evaluations, and the standard LLFs( $\mathbf{u}^h$ ) method using piecewise linear solution and data. Unstructured triangulations are used. The first scheme fits the hypotheses of Proposition 2 (see section §5.2), and we can see on the leftmost picture of figure 10 that indeed convergence to the exact solution w.r.t quadrature accuracy is obtained.

Then, a perturbation is added to the initial total energy level, corresponding to an increase in free surface level of the order of  $0.05m$  (see [54] for details). The evolution of the free surface perturbation  $\delta\eta$  obtained by removing the exact steady solution from the results is reported in the middle column of figure 10. The LLFs( $\mathbf{v}^h$ ) with analytical expressions for the topography provides a perfectly clean evolution of the perturbation. A somewhat noisy result is obtained with the standard method, with spurious effects still relatively small compared to the actual perturbation, which would be hardly obtainable with other approaches. However, this shows how sensitive the results may be to the mesh, and that the multidimensional case may still require some improvements.

## 7.4 Shallow water with dry areas

This is a very classical benchmark reproducing the experiments of [23] on the runup of solitary waves on a conical island. We refer to the above reference, and to [18] for details. We report here results on moving adaptive meshes based on the use of the Lax-Friedrichs' based distribution which allows to have control on the non-negativity of the water depth (cf. [54, 18]). Figure 11 shows a reference solution obtained on a relatively fine mesh (uniformly refined in the interaction region).

Solution contours and adapted meshes obtained with the ALE adaptive method discussed in [18] are shown on figure 12. To check the improvement brought by adaptation we report on figure 13 the time series of the water elevation at the front and rear of the island (top row), and the CPU times. The coarse mesh (dashed blue curves) clearly fails to provide correct values of the maximum runup levels, which are much better predicted with the adaptive ALE approach (dashed red curves), with CPU times of less than 30% those of the fine mesh.

The bottom-right plot on the same figure proves the effectiveness of the mass conservative bathymetry remap of [18] also in presence of dry areas.

## 8 Conclusion and outlook

In this paper, we have described the numerical framework known as Residual Distribution. This formalism encompasses most, if not all, the known schemes (at least known by the authors), both on structured and unstructured meshes. In particular, we have shown that despite not being formulated in terms of numerical fluxes, one can easily provide explicit definitions of local numerical fluxes, with a consistency defined in the classical sense. This result applies to continuous finite elements functional approximation, for example, showing that they are locally conservative. Examples of explicit formula for the flux have been provided in the multi-dimensional case. For non-homogenous problems, these methods are by their nature well-balanced. We have in particular shown that the classical weighted residual formulation

boils down to a global flux method in one space dimension, allowing to embed a more general notion of consistency better suited for balance laws. These properties naturally generalize to arbitrary order of accuracy via appropriate choices for the underlying polynomial approximation and quadrature formulas. Other extensions of the method have not been discussed here, as e.g. the consistent treatment of higher order derivatives [14, 9, 47, 59, 48] and the treatment of non-conservative models [6, 17].

As the reader might have seen, there are some similarities between the RD schemes and the wave propagation algorithm by R. LeVeque, see e.g. [44], where the emphasis is put on the upwind nature of the algorithm as well as a non linear stabilisation on the flux. There is also similarities with the methods developed by A. Lerat and his co-authors, see e.g. [43, 36]. They are residual methods, formally they look like the SUPG scheme with the difference that they introduce a non linear stabilisation to deal with the flow discontinuities. There are more differences with the recently developed active flux method however, see [22, 50, 28]: the approximation of the solution is completely different, and the evolution operator uses in depth the structure of the exact one. These schemes are probably in their infancy and further development will certainly occur, see [38, 28, 22, 16].

The topics discussed are the basis for many future possible developments aiming in particular at further generalizing some of the properties discussed, as e.g. the notion of well balanced in multiple dimensions, as well as further combining them with adaptive strategies.

## Acknowledgements

We acknowledge the work of many PhD and post-docs students: M. Mezine, C. Tavé, A. Larat, G. Baurin, D. de Santis, L. Arpaia, A. Filippini, P. Bacigaluppi, D. Torlo, and S. Tokareva. Discussions with S. Karni and P. Roe (both from the University of Michigan) are also warmly acknowledged.

RA has been partially financed by SNF grant 200020\_175784 and an International Chair of Inria.

## References

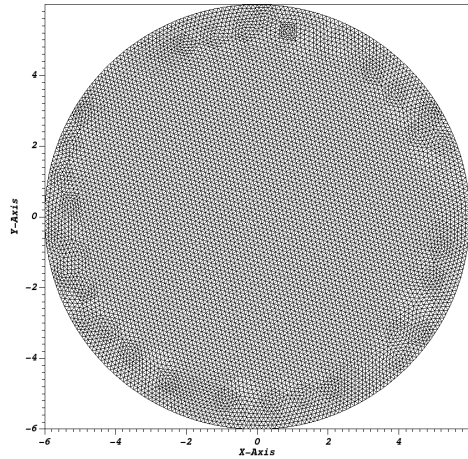
- [1] R. Abgrall. Toward the ultimate conservative scheme: Following the quest. *J. Comput. Phys.*, 167(2):277–315, 2001.
- [2] R. Abgrall. Essentially non-oscillatory residual distribution schemes for hyperbolic problems. *J. Comput. Phys.*, 214(2):773–808, 2006.
- [3] R. Abgrall. High order schemes for hyperbolic problems using globally continuous approximation and avoiding mass matrices. *J. Sci. Comput.*, 73(2-3):461–494, 2017.
- [4] R. Abgrall. Some remarks about conservation for residual distribution schemes. *Comput. Methods Appl. Math.*, 18(3):327–351, 2018.
- [5] R. Abgrall. The Notion of Conservation for Residual Distribution Schemes (or Fluctuation Splitting Schemes), with Some Applications. *Commun. Appl. Math. Comput.*, 2:341–368, 2020.
- [6] R. Abgrall, P. Bacigaluppi, and S. Tokareva. A high-order nonconservative approach for hyperbolic equations in fluid dynamics. *Computers & Fluids*, 169:10 – 22, 2018. Recent progress in nonlinear numerical methods for time-dependent flow & transport problems.

- [7] R. Abgrall, P. Bacigaluppi, and S. Tokareva. “A *Posteriori*” limited high order and robust schemes for transient simulations of fluid flows in gas dynamics. *Journal of Computational Physics*, 2020. in revision.
- [8] R. Abgrall, G. Baurin, P. Jacq, and M. Ricchiutto. Some examples of high order simulations in parallel of inviscid flows on unstructured and hybrid meshes by residual distribution schemes. *Computers and Fluids*, 61:6–13, 2012.
- [9] R. Abgrall and D. de Santis. Linear and non-linear high order accurate residual distribution schemes for the discretization of the steady compressible Navier-Stokes equations. *J. Comput. Phys.*, 283:329–359, 2015.
- [10] R. Abgrall and S. Karni. A comment on the computation of non-conservative products. *J. Comput Phys.*, 229(8):2759–276, 2010.
- [11] R. Abgrall, A. Larat, and M. Ricchiutto. Construction of very high order residual distribution schemes for steady inviscid flow problems on hybrid unstructured meshes. *J. Comput. Phys.*, 230(11):4103–4136, 2011.
- [12] R. Abgrall and M. Mezzine. Construction of second-order accurate monotone and stable residual distribution schemes for unsteady flow problems. *J. Comput. Phys.*, 188:16–55, 2003.
- [13] R. Abgrall and M. Ricchiutto. High order methods for CFD. In Rene de Borst Erwin Stein and Thomas J.R. Hughes, editors, *Encyclopedia of Computational Mechanics, Second Edition*. John Wiley and Sons, 2017.
- [14] R. Abgrall, M. Ricchiutto, and D. de Santis. High-order preserving residual distribution schemes for advection-diffusion scalar problems on arbitrary grids. *SIAM J. Sci. Comput.*, 36(3):A955–A983, 2014.
- [15] R. Abgrall and P. L. Roe. High-order fluctuation schemes on triangular meshes. *J. Sci. Comput.*, 19(1-3):3–36, 2003.
- [16] Rémi Abgrall. A combination of residual distribution and the active flux formulations or a new class of schemes that can combine several writings of the same hyperbolic problem: application to the 1d euler equations, 2020. <https://arxiv.org/abs/2011.12572>.
- [17] Rémi Abgrall, Paola Bacigaluppi, and Barbara Re. On the simulation of multicomponent and multiphase compressible flows, 2020. <https://arxiv.org/abs/2006.01630>.
- [18] L. Arpaia and M. Ricchiutto. r-adaptation for Shallow Water flows: conservation, well balancedness, efficiency. *Computers & Fluids*, 160:175–203, 2018.
- [19] L. Arpaia and M. Ricchiutto. Well balanced residual distribution for the ALE spherical shallow water equations on moving adaptive meshes. *Journal of Computational Physics*, 405, 2020. Article 109173.
- [20] L. Arpaia, M. Ricchiutto, and R. Abgrall. An ALE formulation for explicit Runge-Kutta residual distribution. *J. Sci. Comput.*, 190(34):1467–1482, 2014.
- [21] J. Balbás and S. Karni. A central scheme for shallow water flows along channels with irregular geometry. *ESAIM: Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique*, 43(2):333–351, 2009.
- [22] Wasilij Barsukow, Jonathan Hohm, Christian Klingenberg, and Philip L. Roe. The active flux scheme on Cartesian grids and its low Mach number limit. *J. Sci. Comput.*, 81(1):594–622, 2019.

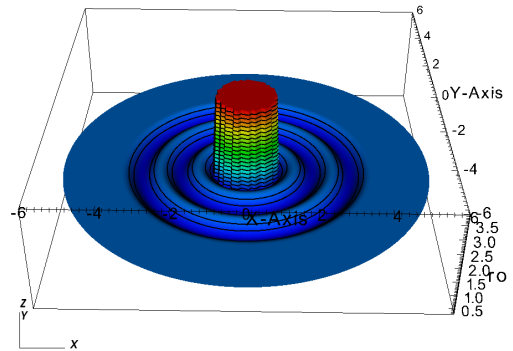
- [23] M.J. Briggs, C.E. Synolakis, G.S.Harkins, and D.R. Green. Laboratory experiments of tsunami runup on a circular island. *Pure Appl. Geophys.*, 144:569–593, 1995.
- [24] E. Burman and P. Hansbo. Edge stabilization for Galerkin approximation of convection-diffusion-reaction problems. *Comput. Methods Appl. Mech. Engrg*, 193:1437–1453, 2004.
- [25] E. Burman, A. Quarteroni, and B. Stamm. Interior penalty continuous and discontinuous finite element approximations of hyperbolic equations. *J. Sci. Comput.*, 43(3):293–312, 2010.
- [26] M.J. Castro, T. Morales de Luna, and C. Parés. Chapter 6 - Well-Balanced Schemes and Path-Conservative Numerical Methods. In R. Abgrall and C.-W. Shu, editors, *Handbook of Numerical Methods for Hyperbolic Problems*, volume 18 of *Handbook of Numerical Analysis*, pages 131 – 175. Elsevier, 2017.
- [27] Y. Cheng, A. Chertock, M. Herty, A. Kurganov, and T. Wu. A new approach for designing moving-water equilibria preserving schemes for the shallow water equations. *J.Sci.Comp.*, 80:538–554, 2019.
- [28] Erik Chudzik, Christiane Helzel, and David Kerkmann. The Cartesian grid active flux method: linear stability and bound preserving limiting. *Appl. Math. Comput.*, 393:125501, 19, 2021.
- [29] P. Ciarlet. *The finite element method for elliptic problems*. North-Holland, Amsterdam, 1978.
- [30] H. Deconinck and M. Ricchiuto. Residual Distribution Schemes: Foundations and Analysis. In *Encyclopedia of Computational Mechanics, second edition*. John Wiley and Sons, 2017.
- [31] A.I. Delis and T. Katsaounis. Relaxation schemes for the shallow water equations. *International Journal for Numerical Methods in Fluids*, 41:695–719, 2003.
- [32] A.I. Delis and T. Katsaounis. Numerical solution of the two-dimensional shallow water equations by the application of relaxation methods. *Applied Mathematical Modelling*, 29:754–783, 2005.
- [33] J. Dobes and H. Deconinck. A second order space-time residual distribution method for solving compressible flow on moving meshes. 43rd AIAA Aerospace Sciences Meeting and Exhibit, June 2012. <https://arc.aiaa.org/doi/abs/10.2514/6.2005-493>.
- [34] J. Dobes, M. Ricchiuto, and H. Deconinck. Implicit space-time residual distribution method for unsteady laminar viscous flow. *Computers & Fluids*, 34(4-5):593 – 615, 2005.
- [35] R. Donat, M.C. Martí, A. Martínez-Gavara, and P. Mulet. Well-balanced adaptive mesh refinement for shallow water flows. *Journal of Computational Physics*, 257:937 – 953, 2014.
- [36] Xi Du, Christophe Corre, and Alain Lerat. A third-order finite-volume residual-based scheme for the 2D Euler equations on unstructured grids. *J. Comput. Phys.*, 230(11):4201–4215, 2011.
- [37] A. Ern and J.L. Guermond. *Theory and practice of finite elements*, volume 159 of *Applied Mathematical Sciences*. Springer verlag, 2004.
- [38] Christiane Helzel, David Kerkmann, and Leonardo Scandurra. A new ADER method inspired by the active flux method. *J. Sci. Comput.*, 80(3):1463–1497, 2019.
- [39] M. Hubbard and M. Ricchiuto. Discontinuous upwind residual distribution: A route to unconditional positivity and high order accuracy. *Computers and Fluids*, 46(1):263 – 269, 2011.
- [40] M. Hubbard, M. Ricchiuto, and D. Sarmany. Space-time residual distribution on moving meshes. *Computers & Mathematics with Applications*, 79:1561–1589, 2020.

- [41] Th. J.R. Hughes, G. Engel, L. Mazzei, and M. G. Larson. The Continuous Galerkin Method Is Locally Conservative. *Journal of Computational Physics*, 163(2):467 – 488, 2000.
- [42] T.J.R. Hughes, L.P. Franca, and M. Mallet. A new finite element formulation for CFD: I. symmetric forms of the compressible Euler and Navier-Stokes equations and the second law of thermodynamics. *Comp. Meth. Appl. Mech. Engrg.*, 54:223–234, 1986.
- [43] A. Lerat. An efficient high-order compact scheme for the unsteady compressible Euler and Navier-Stokes equations. *J. Comput. Phys.*, 322:365–386, 2016.
- [44] R. J. LeVeque. Wave propagation algorithms for multi-dimensional hyperbolic systems. *J. Comput. Phys.*, 131(2):327–353, 1997.
- [45] H. Li, S. Xie, and X. Zhang. A high order accurate bound-preserving compact finite difference scheme for scalar convection diffusion equations. *SIAM Journal on Numerical Analysis*, 56:3308–3345, 2018.
- [46] A. Loseille. Chapter 10 - unstructured mesh generation and adaptation. In Rémi Abgrall and Chi-Wang Shu, editors, *Handbook of Numerical Methods for Hyperbolic Problems*, volume 18 of *Handbook of Numerical Analysis*, pages 263 – 302. Elsevier, 2017.
- [47] A. Mazaheri and H. Nishikawa. Improved second-order hyperbolic residual-distribution scheme and its extension to third-order on arbitrary triangular grids. *Journal of Computational Physics*, 300:455 – 491, 2015.
- [48] A. Mazaheri, M. Ricchiuto, and H. Nishikawa. A first-order hyperbolic system approach for dispersion. *Journal of Computational Physics*, 321:593 – 605, 2016.
- [49] C. Michler and H. Deconinck. An arbitrary lagrangian eulerian formulation for residual distribution schemes on moving grids. *Computers and Fluids*, 32(1):59–71, 2001.
- [50] Hiroaki Nishikawa and Philip L. Roe. Third-order active-flux scheme for advection diffusion: hyperbolic diffusion, boundary condition, and Newton solver. *Comput. & Fluids*, 125:71–81, 2016.
- [51] S. Noelle, Y. Xing, and C.-W. Shu. High order well-balanced finite volume weno schemes for shallow water equation with moving water. *J. Comput. Phys.*, 226:29–58, 2007.
- [52] B. Re, C. Dobrzynsk, and A. Guardone. An interpolation-free ALE scheme for unsteady inviscid flows computations with large boundary displacements over three-dimensional adaptive grids. *Journal of Computational Physics*, 340:26 – 54, 2017.
- [53] M. Ricchiuto. On the C-property and Generalized C-property of Residual Distribution for the Shallow Water Equations. *Journal of Scientific Computing*, 48:304–318, 2011.
- [54] M. Ricchiuto. An explicit residual based approach for shallow water flows. *J. Comput. Phys.*, 80:306–344, 2015.
- [55] M. Ricchiuto and R. Abgrall. Explicit Runge-Kutta residual distribution schemes for time dependent problems: second order case. *J. Comput. Phys.*, 229(16):5653–5691, 2010.
- [56] M. Ricchiuto, R. Abgrall, and H. Deconinck. Application of conservative residual distribution schemes to the solution of the shallow water equations on unstructured meshes,. *J. Comput. Phys.*, 222:287–331, 2007.
- [57] M. Ricchiuto and A. Bollermann. Stabilized residual distribution for shallow water simulations. *J. Comput. Phys*, 228(4):1071–1115, 2009.

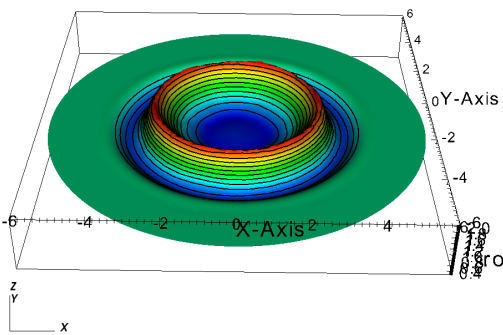
- [58] M. Ricchiuto, Á. Csík, and H. Deconinck. Residual distribution for general time-dependent conservation laws. *J. Comput. Phys.*, 209(1):249–289, 2005.
- [59] M. Ricchiuto and A.G. Filippini. Upwind residual discretization of enhanced boussinesq equations for wave propagation over complex bathymetries. *J. Comput. Phys.*, 271:306–341, 2014.
- [60] M. Ricchiuto, D. T. Rubino, J. Witteveen, and H. Deconinck. A residual distributive approach for one-dimensional two-fluid models and its relation to godunov finite volume schemes. In *ASTAR International Workshop on Advanced Numerical Methods for Multidimensional Simulation of Two phase Flow*, Garching, Germany, 2003.
- [61] P.L. Roe. Approximate Riemann solvers, parameter vectors, and difference schemes. *J. Comput. Phys.*, 43:357–372, 1981.
- [62] P.L. Roe. Upwind differencing schemes for hyperbolic conservation laws with source terms. In Claude Carasso, Denis Serre, and Pierre-Arnaud Raviart, editors, *Nonlinear Hyperbolic Problems*, pages 41–51, Berlin, Heidelberg, 1987. Springer Berlin Heidelberg.
- [63] B.D. Rogers, A. Borthwick, and P. Taylor. Mathematical balancing of flux gradient and source terms prior to using Roe’s approximate Riemann solver. *J.Comput.Phys.*, 192:422–451, 2003.
- [64] J.A. Rossmanith, D.S. Bale, and R.J. LeVeque. A wave propagation algorithm for hyperbolic systems on curved manifolds. *J.Comput.Phys.*, 199:61–662, 2004.
- [65] D. Sarmany, M. Hubbard, and M. Ricchiuto. Unconditionally stable space-time discontinuous residual distribution for shallow water flows. *Journal of Computational Physics*, 253:86–113, 2013.
- [66] H. Staedtke, G. Franchello, B. Worth, U. Graf, P. Romstedt, A. Kumbaro, J. Garcia-Cascales, H. Paillere, H. Deconinck, M. Ricchiuto, B. Smith, F. De Cachard, E.F. Toro, E. Romenski, and S. Mimouni. Advanced three-dimensional two-phase flow simulation tools for application to reactor safety (astar). *Nuclear Engineering and Design*, 235(2):379 – 400, 2005.
- [67] R. Struijs, H. Deconinck, and P.L. Roe. Fluctuation splitting schemes for the 2D Euler equations. VKI-LS 1991-01, 1991. Computational Fluid Dynamics.
- [68] P.D. Thomas and C.K Lombard. Geometric conservation law and its application to flow computations on moving grids. *AIAA Journal*, pages 1030–1037, 1979.
- [69] E. Valero, J. de Vicente, and G. Alonso. The application of compact residual distribution schemes to two-phase flow problems. *Computers and Fluids*, 38(10):1950 – 1968, 2009.
- [70] E. Valero, M. Ricchiuto, and G. Degrez. Two-phase flow computations using a two-fluid model and fluctuation splitting. *Trends in Numerical and Physical Modeling for Industrial Two-Phase Flows*, Cargese, France, September 2000.
- [71] N. Villedieu, T. Quintino, M. Ricchiuto, and H. Deconinck. Third order residual distribution schemes for the Navier-Stokes equations. *J. Comput. Phys.*, 230(11):4301–4315, 2011.
- [72] F. Zhou, G. Chen, Y. Huang, J. Z. Yang, and H. Feng. An adaptive moving finite volume scheme for modeling flood inundation over dry and complex topography. *Water Resour. Res.*, 49:1914–1928, 2013.



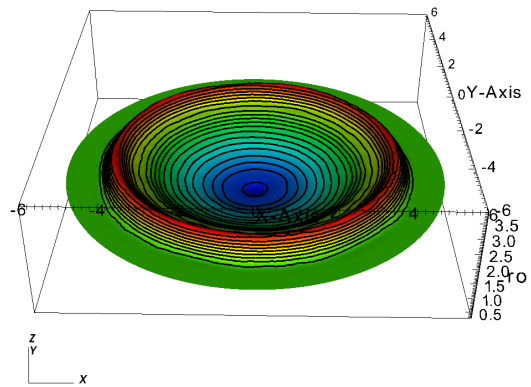
(a)



(b)



(c)



(d)

Figure 4: (a): mesh, all the degrees of freedom are represented, (b): initial condition, (c) Solution at an intermediate time, (d) Solution at  $t = 1.8$ . The CFL is 0.125.



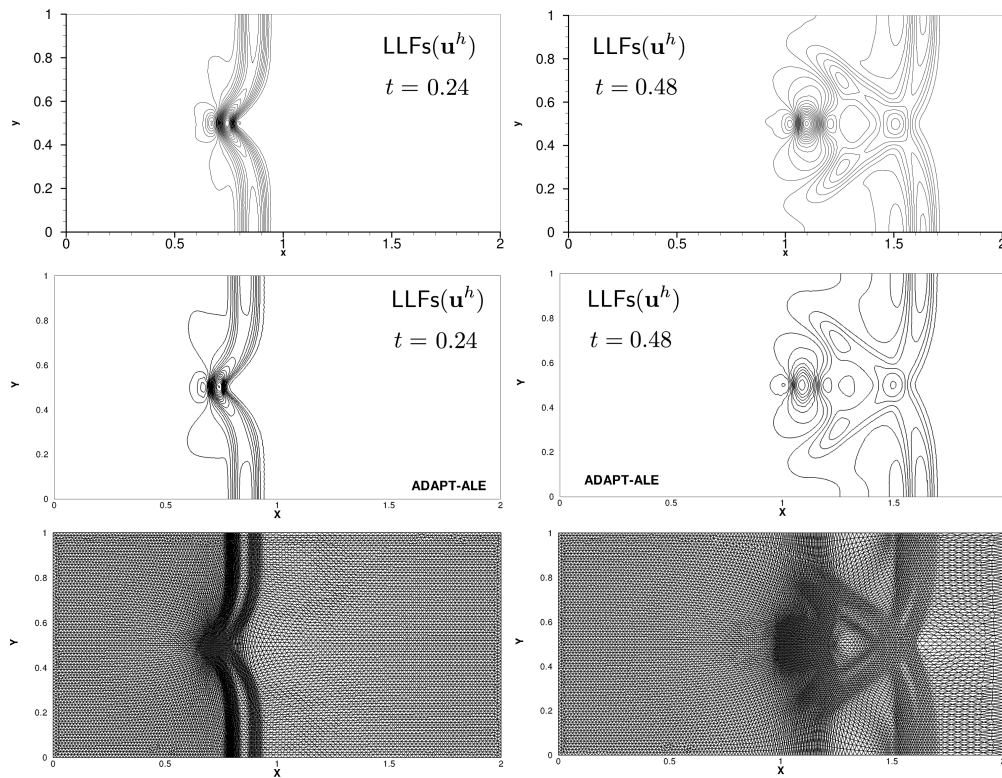


Figure 5: Perturbation of the lake at rest over a smooth topography: free surface evolution. Top: implicit scheme of [57]. Middle and bottom : explicit adaptive moving mesh approach by [18]

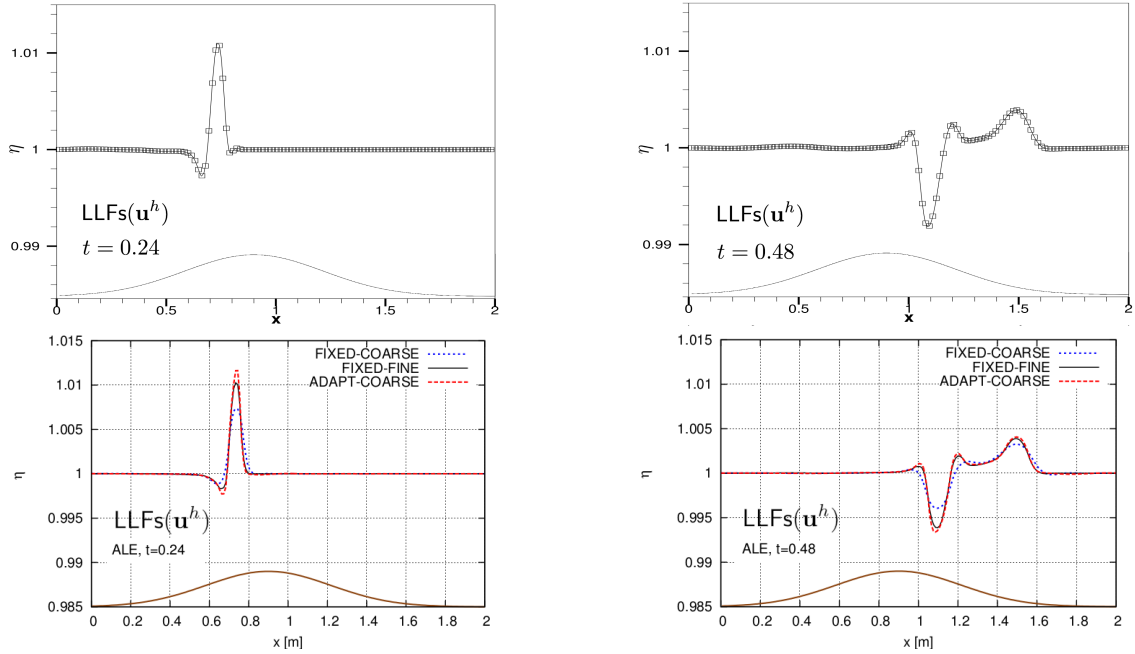


Figure 6: Perturbation of the lake at rest over a smooth topography: centerline free surface levels. Top: implicit scheme of [57]. Bottom : explicit adaptive moving mesh approach by [18]

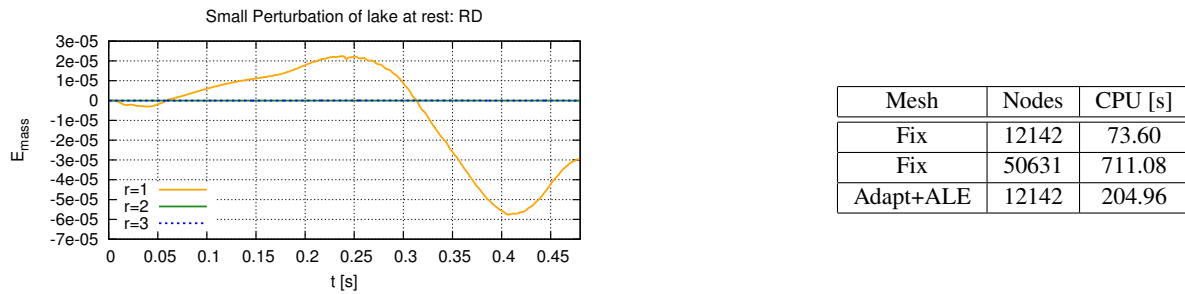


Figure 7: Mass conservation in function of the quadrature order for the topography remap, and CPU times for the explicit schemes

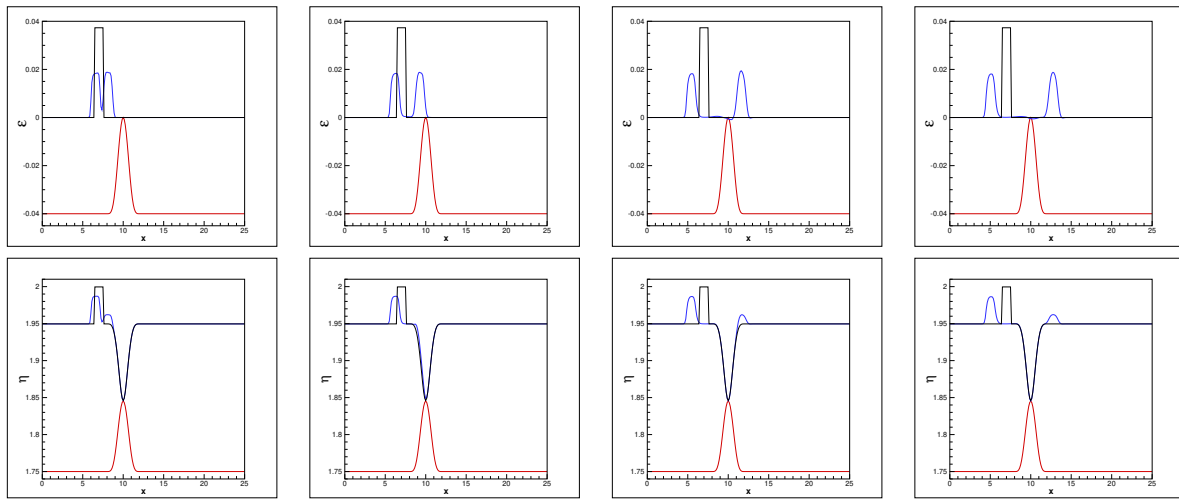


Figure 8: Example of a 1D perturbation in a constant energy steady flow.s

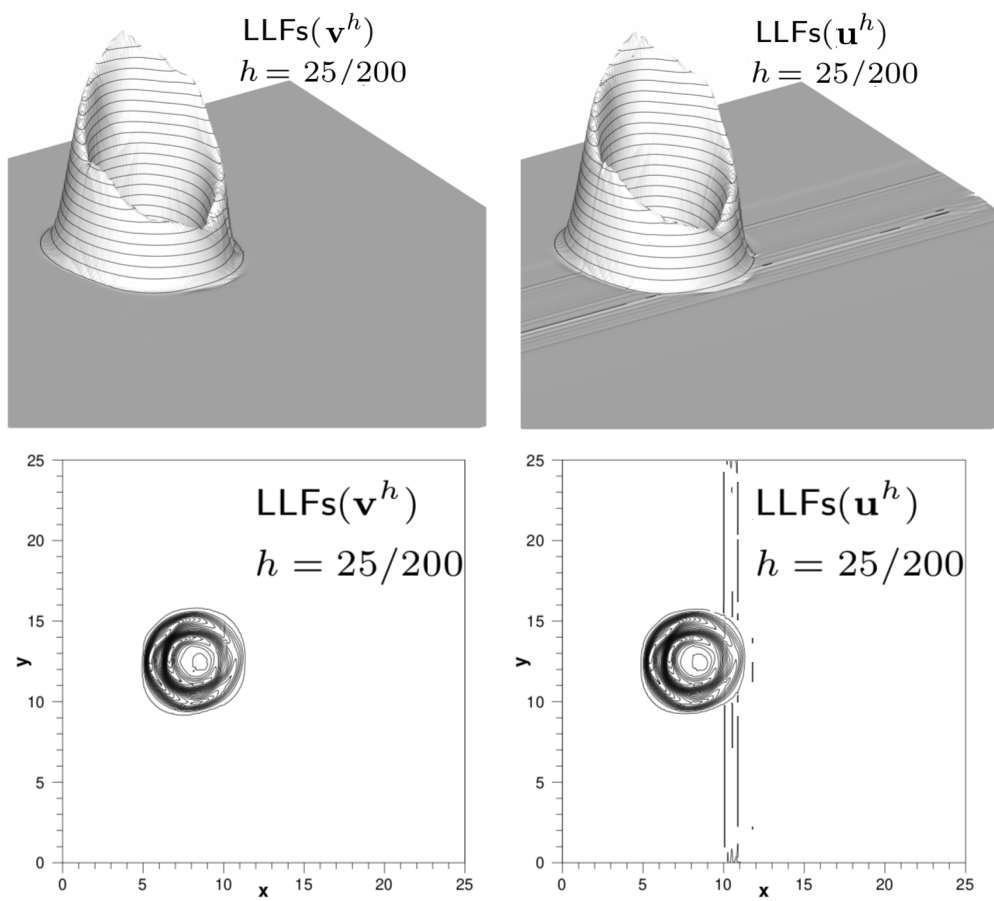


Figure 9: Example of a 2D perturbation in a constant energy steady flow.

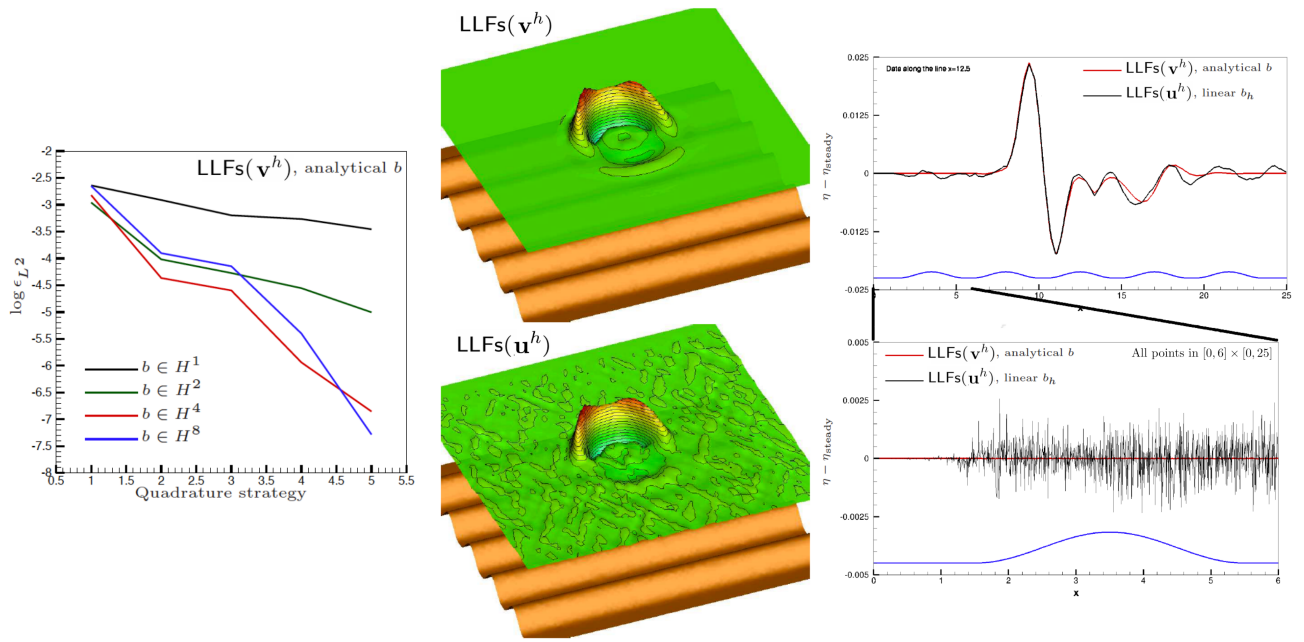


Figure 10: Example of a 2D perturbation in a constant energy steady flow.

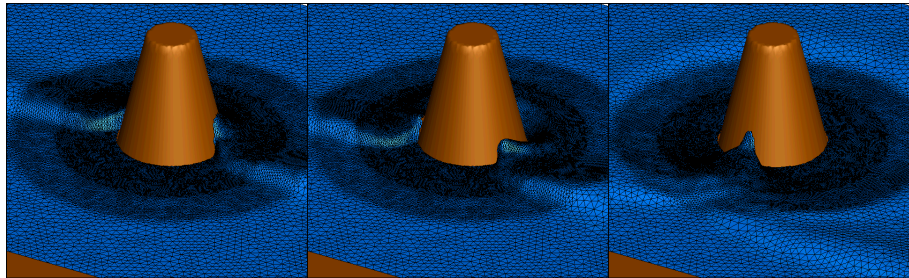


Figure 11: Conical Island. 3D view of the wave runup process, and gauge locations.

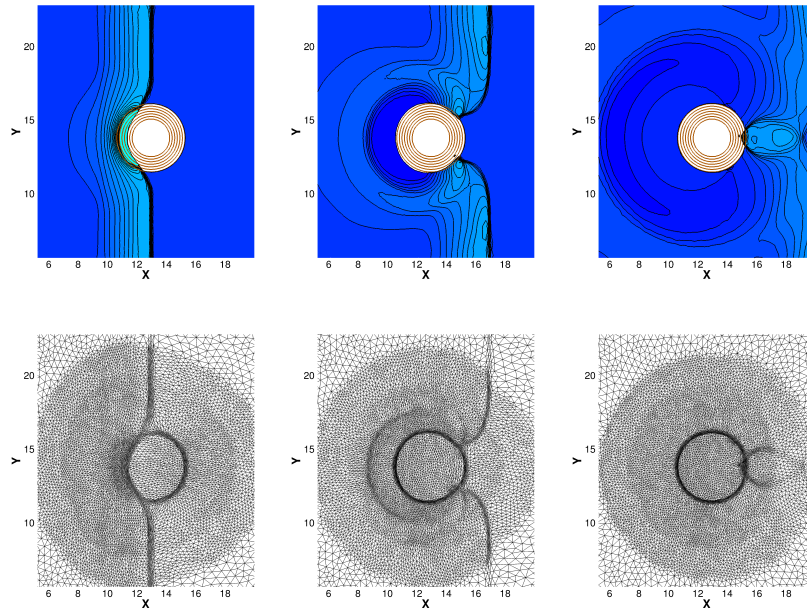
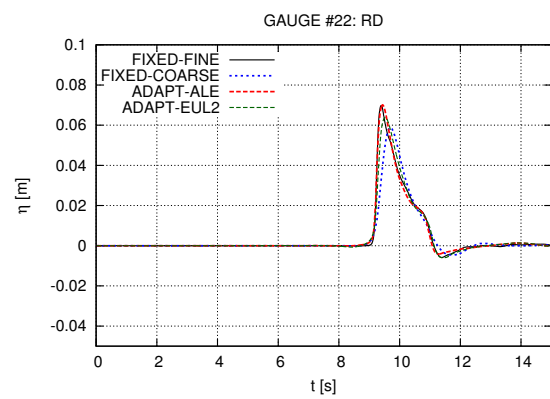
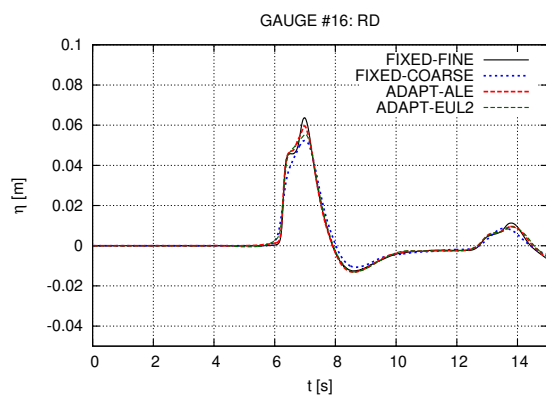


Figure 12: Conical Island; free surface contours at  $t = 6.0, 8.0, 10.0$  [s], and adapted meshes



Mesh	Nodes	CPU [s] (%MMPDE)
Fix	10401	171.30
Fix	37982	1785.96
Adapt+ALE	10401	510.65 (38.8%)

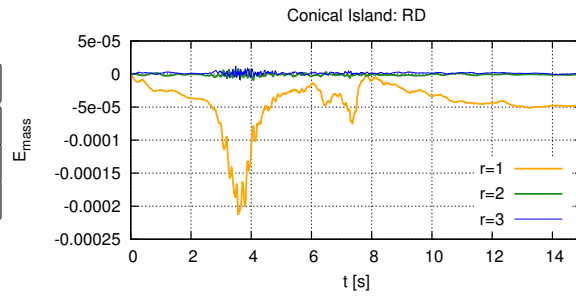


Figure 13: Conical Island. Top: free surface time series at the front and rear runup gauges. Bottom: mass conservation (left) and CPU times (right).