



**HAL**  
open science

# A Posteriori Local Subcell Correction of High-Order Discontinuous Galerkin Scheme for Conservation Laws on Two-Dimensional Unstructured Grids

François Vilar, Rémi Abgrall

► **To cite this version:**

François Vilar, Rémi Abgrall. A Posteriori Local Subcell Correction of High-Order Discontinuous Galerkin Scheme for Conservation Laws on Two-Dimensional Unstructured Grids. *SIAM Journal on Scientific Computing*, 2024, 46 (2), pp.A851-A883. 10.1137/22M1542696 . hal-04556790

**HAL Id: hal-04556790**

**<https://hal.science/hal-04556790>**

Submitted on 23 Apr 2024

**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.

# A *POSTERIORI* LOCAL SUBCELL CORRECTION OF HIGH-ORDER DISCONTINUOUS GALERKIN SCHEME FOR CONSERVATION LAWS ON TWO-DIMENSIONAL UNSTRUCTURED GRIDS

FRANÇOIS VILAR \* AND RÉMI ABGRALL†

**Abstract.** In this paper, we present the two-dimensional unstructured grids extension of the *a posteriori* local subcell correction (APLSC) of discontinuous Galerkin (DG) schemes introduced in [42]. The technique is based on the reformulation of DG scheme as a finite volume (FV) like method through the definition of some specific numerical fluxes referred to as reconstructed fluxes. High-order DG numerical scheme combined with this new local subcell correction technique ensures positivity preservation of the solution, along with a low oscillatory and sharp shocks representation.

The main idea of this correction procedure is to retain as much as possible the high accuracy and the very precise subcell resolution of DG schemes, while ensuring the robustness and stability of the numerical method, as well as preserving physical admissibility of the solution. Consequently, an *a posteriori* correction will only be applied locally at the subcell scale where it is needed, but still ensuring the scheme conservativity. Practically, at each time step, we compute a DG candidate solution and check if this solution is admissible (for instance positive, non-oscillating, ...). If it is the case, we go further in time. Otherwise, we return to the previous time step and correct locally, at the subcell scale, the numerical solution. To this end, each cell is subdivided into subcells. Then, if the solution is locally detected as bad, we substitute the DG reconstructed flux on the subcell boundaries by a robust first-order numerical flux. For subcell detected as admissible, we keep the high-order DG reconstructed flux which allows us to retain the very high accurate resolution and conservation of the DG scheme. As a consequence, only the solution inside troubled subcells and its first neighbors will have to be recomputed, elsewhere the solution remains unchanged. Another technique blending in a convex combination fashion DG reconstructed fluxes and first-order FV fluxes for admissible subcells in the vicinity of troubled areas will also be presented and prove to improve results in comparison to the original algorithm introduced in [42]. Numerical results on various type of problems and test cases will be presented to assess the very good performance of the designed correction algorithm.

**Key word.** *a posteriori* correction, subcell correction, arbitrary high-order, DG subcell FV formulation, positivity-preserving scheme, hyperbolic conservation laws, subcell conservative scheme

**1. Introduction.** This paper is devoted to the extension of the *a posteriori* local subcell correction (APLSC) method introduced in [42] to the two-dimensional unstructured case. It is well known that hyperbolic partial differential equations generally lead, in finite time, to discontinuous weak solutions. In numerical simulations, this aspect has to be dealt with, and this issue has been one of the main questions to address in the design of numerical methods for hyperbolic problems. Another core problem is the one of accuracy. The last one is that the problem is well set only if the solution belongs to an admissibility set called invariant domain. For example, for the gas dynamics Euler equations, the density and the internal energy must stay positive. It is particularly difficult to address, simultaneously, these three questions together, because those constraints are antagonistic. These questions have been at the center of algorithmic developments since decades, one may mention [3, 41, 18, 23, 52] and the reference herein.

The Discontinuous Galerkin (DG) method is one of the most widely used numerical scheme, especially in the context of computational fluid dynamics. Cockburn, Shu *et al.* in a series papers (see [7] and the reference therein) have paved the road to efficient methods for fluid dynamics. DG methods allow to reach any arbitrary order of accuracy, while keeping the stencil compact, along with other good properties such as built-in entropy stability and *hp*-adaptivity. However, though recent progress have been done [37], designing methods that are oscillation free and compliant with the invariant domain is not a trivial matter. Controlling spurious oscillations has been studied in many papers, among which [2, 4, 25, 50, 22, 28, 33]. Staying in the invariant domain has also been made possible, see for example [52, 49, 15]. However, this is often achieved to the cost of enlarging the width of the discontinuous patterns, and some time a loss of accuracy.

Other methods have recently gained in popularity, the so-called subcell techniques. Here the idea is to subdivide the bad cells, and to adopt a special procedure with the hope of curing the negative aspects of the original scheme. Some examples of this strategy can be found in [20, 40]. For example, [20], the authors use a convex combination between high-order DG schemes and first-order Finite-Volumes (FV) on a subgrid, allowing them to retain the very high accurate resolution of DG in smooth areas and ensuring the scheme's robustness in the presence of shocks. Similarly, in [40], after having detected the troubled zones, cells are

---

\*IMAG, Univ Montpellier, CNRS, Montpellier, France ([francois.vilar@umontpellier.fr](mailto:francois.vilar@umontpellier.fr)).

†Institut für Mathematik, Universität Zürich, CH-8057 Zürich, Switzerland ([remi.abgrall@uzh.ch](mailto:remi.abgrall@uzh.ch)).

then subdivided into subcells, and a robust first-order finite volume scheme is performed on the subgrid in troubled cells. Alternatively, some robust high-order scheme as MUSCL or WENO could either be used to avoid too much accuracy discrepancy. Note that in general these methods are tuned for Cartesian meshes.

Another approach which is worth to be mentioned is the so-called MOOD technique, [6, 8, 9]. Through this procedure, the order of approximation of the numerical scheme is locally reduced in an *a posteriori* sequence until the solution becomes admissible, *i.e.* oscillation free and the solution lives in the invariant domain. In [11, 10], a subcell FV technique similar to the one presented in [40] has been applied to the *a posteriori* paradigm. In practice, if the numerical solution in a cell is detected as bad, the cell is then subdivided into subcells and a first-order finite volume, or alternatively other robust scheme (second-order TVD FV scheme, WENO scheme, ...), is applied on each subcell. Then, through these new subcell mean values, a high-order polynomial is reconstructed on the primal cell. This correction procedure has the benefit to be very simple and robust, and is able to preserve the high accuracy of DG schemes in smooth areas.

In all these aforementioned limitation techniques, *a priori* and *a posteriori*, the high-order DG polynomial is either globally modified in the cell, or even discarded as it is in the (H)WENO limiter or any *a posteriori* correction technique in troubled cells. Since one of the main advantage of high-order scheme is to be able to use coarse grids while still being very precise, one can see that there is a waste of information here, as well as unnecessary computational effort made. This problem was addressed for the one-dimensional case in [42]. This new technique relies on the reformulation of DG schemes as a FV-like scheme defined on a subgrid. First, as the number of subcells matches the dimension of the polynomial solution space, the numerical solution inside a cell can be uniquely defined by either a high degree polynomial, or as a piecewise constant solution through its different subcell mean values; the connection between the two being done via a projector. Second, by means of particular basis functions introduced in [42], which are nothing but the  $L_2$  projection over the polynomial space of the subcells indicator function, the DG volume and boundary terms can be rewritten as flux differences. Such theoretical analysis is relatively simple in one dimension in space, but much more challenging in two dimensions.

The format of the paper is as follows. Extending the theoretical analysis introduced in [42] and using ideas from [36], we first reinterpret unstructured grid DG scheme as a subgrid FV-like scheme, through the definition of particular fluxes that we referred to as reconstructed fluxes. These DG reconstructed fluxes are analytically computed, and the analysis shows how those fluxes are connected to the interior polynomial flux and the jump at the cell interface between the interior flux and the DG numerical flux. Let us emphasize that reformulation of DG scheme as a subcell finite-volume method can be performed regardless the form of the element: this is thus not limited to quadrilateral nor triangular cells, but can be done on general polygonal elements. This theoretical part is done in section 2, where a discussion on the type of subcells is also provided. Using this equivalent formulation, we can proceed by means of an *a posteriori* paradigm as follows: at each time step, we compute a DG candidate solution and check if this solution is admissible. If it is the case, we go further in time. Otherwise, we return to the previous time step and correct locally, at the subcell scale, the numerical solution. In the subcells where the solution was detected as bad, we substitute the DG reconstructed flux on the subcell boundaries by a robust first-order numerical flux. For subcells detected as admissible, we keep the high-order reconstructed flux which allows us to retain the very high accurate resolution and conservation property of DG schemes. Consequently, only the solution inside troubled subcells and their first neighbors will have to be recomputed. Elsewhere, the solution remains unchanged. This correction procedure is then extremely local. Another technique blending in a convex combination fashion DG reconstructed fluxes and first-order FV fluxes for admissible subcells in the vicinity of troubled areas will also be presented and prove to improve results in comparison to the original algorithm. This procedure and all the technical details are described in section 3. Section 4 provides numerical results, and a conclusion follows.

**2. DG method as a subcell FV scheme.** This section is devoted to the demonstration of the equivalency between DG schemes and a FV-like method on a subgrid provided the definition of particular fluxes. To remain as simple as possible, two-dimensional Scalar Conservation Laws (SCL) will be considered. Let  $u = u(\mathbf{x}, t)$ , for  $\mathbf{x} \in \omega \subset \mathbb{R}^2$  and  $t \in [0, T]$ , be the solution of the following system,

$$\begin{aligned} (2.1a) \quad & \begin{cases} \partial_t u(\mathbf{x}, t) + \nabla_{\mathbf{x}} \cdot \mathbf{F}(u(\mathbf{x}, t)) = 0, \\ u(\mathbf{x}, 0) = u_0(\mathbf{x}), \end{cases} & (\mathbf{x}, t) \in \omega \times [0, T], \\ (2.1b) \quad & & \mathbf{x} \in \omega, \end{aligned}$$

where  $u_0$  is the initial data and  $\mathbf{F}(u)$  the 2D flux function. For the subsequent discretization, let us introduce the following notation. Let  $\{\omega_c\}_c$  be a generic partition of the domain  $\omega$  into non-overlapping cells, with  $|\omega_c|$  being the size of  $\omega_c$ . We also partition the time domain in intermediate times  $(t^n)_n$  with  $\Delta t^n = t^{n+1} - t^n$  the  $n^{\text{th}}$  time step. In the DG frame, the numerical solution is considered piecewise polynomial over the domain, and hence developed on each cell onto  $\mathbb{P}^k(\omega_c)$ , the set of polynomials of degree up to  $k$  defined on cell  $\omega_c$ . This space approximation theoretically leads to a  $(k+1)^{\text{th}}$  space order accurate scheme. Let  $u_h^c$  be the restriction of  $u_h$ , the piecewise polynomial approximation of the solution  $u$ , over the cell  $\omega_c$

$$(2.2) \quad u_h^c(\mathbf{x}, t) = \sum_{m=1}^{N_k} u_m^c(t) \sigma_m^c(\mathbf{x}),$$

where the  $u_m^c$  are the  $N_k$  successive components of  $u_h^c$  over the polynomial basis  $\{\sigma_m^c\}_m$ . We recall that in the two-dimensional case  $N_k = \frac{(k+1)(k+2)}{2}$ . The coefficients  $u_m^c$  present in (2.2) are the solution's moments to be computed through a local variational formulation on  $\omega_c$ . To this end, one has to multiply equation (2.1a) by  $\psi \in \mathbb{P}^k(\omega_c)$ , a polynomial test function, and then integrate then on  $\omega_c$ . By means of an integration by parts and substituting the solution  $u$  by its approximated polynomial counterpart  $u_h^c$ , one gets

$$(2.3) \quad \int_{\omega_c} \frac{\partial u_h^c}{\partial t} \psi \, dV = \int_{\omega_c} \mathbf{F}(u_h^c) \cdot \nabla_x \psi \, dV - \int_{\partial\omega_c} \psi \mathcal{F}_n \, dS, \quad \forall \psi \in \mathbb{P}^k(\omega_c).$$

The numerical solution  $u_h^c$  is then the unique polynomial function defined in  $\mathbb{P}^k(\omega_c)$  satisfying equation (2.3) for all function  $\psi \in \mathbb{P}^k(\omega_c)$ . In (2.3), the numerical flux function  $\mathcal{F}_n$ , in addition to ensure the scheme conservation, is the cornerstone of any finite volume or DG scheme regarding fundamental considerations as stability, positivity and entropy among others. In this context, this numerical flux is defined as a function of the two states on the left and right of each interface

$$(2.4) \quad \mathcal{F}_n = \mathcal{F}(u^-, u^+, \mathbf{n}),$$

with  $u^- = \lim_{\epsilon \rightarrow 0^+} u_h^c(\mathbf{x}_i - \epsilon \mathbf{n}, t)$  and  $u^+ = \lim_{\epsilon \rightarrow 0^+} u_h^c(\mathbf{x}_i + \epsilon \mathbf{n}, t)$ , where  $\omega_v$  is a face neighboring cell of  $\omega_c$ , while  $\mathbf{x}_i$  and  $\mathbf{n}$  respectively stand for a point and the unit outward normal of the separating interface. From now on, we refer  $\mathcal{V}_c$  to as the set containing the face neighboring cells of  $\omega_c$ . Function  $\mathcal{F}$  is generally obtained through the resolution of an exact or approximated Riemann problem. In the remainder of this paper, for sake of simplicity, we make use of the very well-known global Lax-Friedrichs numerical flux which reads

$$(2.5) \quad \mathcal{F}(u, v, \mathbf{n}) = \frac{(\mathbf{F}(u) + \mathbf{F}(v))}{2} \cdot \mathbf{n} - \frac{\gamma}{2} (v - u),$$

where  $\gamma = \sup_w (||d_w \mathbf{F}(w)||_2)$ .

Now, taking in (2.3) the test function  $\psi$  among the polynomial basis functions leads to the following linear system allowing the calculation of the solution moments  $u_m^c$

$$(2.6) \quad \sum_{m=1}^{N_k} \frac{d u_m^c}{dt} \int_{\omega_c} \sigma_m^c \sigma_p^c \, dV = \int_{\omega_c} \mathbf{F}(u_h^c) \cdot \nabla_x \sigma_p^c \, dV - \int_{\partial\omega_c} \sigma_p^c \mathcal{F}_n \, dS, \quad \forall p \in [1, N_k].$$

The terms  $\int_{\omega_c} \mathbf{F}(u_h^c) \cdot \nabla_x \sigma_p^c \, dV$  and  $\int_{\partial\omega_c} \sigma_p^c \mathcal{F}_n \, dS$  are respectively referred to as volume and surface integrals. In (2.6), we identify  $\int_{\omega_c} \sigma_m^c \sigma_p^c \, dV = (M_c)_{mp}$  as the generic coefficient of the symmetric mass matrix  $M_c$ . The scheme (2.6) can then be reformulated in a compact matrix-vector form as follows

$$(2.7) \quad M_c \frac{d \mathbf{U}_c}{dt} = \Phi_c,$$

with  $(\mathbf{U}_c)_m = u_m^c$  the solution vector filled with the polynomial moments, and where the so-called DG residuals  $\Phi_c$  write

$$(2.8) \quad (\Phi_c)_m = \int_{\omega_c} \mathbf{F}(u_h^c) \cdot \nabla_x \sigma_m^c \, dV - \int_{\partial\omega_c} \sigma_m^c \mathcal{F}_n \, dS.$$

Similarly to what has been done in the one-dimensional case in [42], let us now demonstrate the equivalency between discontinuous Galerkin schemes and a finite volume like method on a subgrid, and exhibit the corresponding subcell numerical fluxes that will be referred to as high-order reconstructed fluxes. To do so, we first need to subdivide the mesh cells into subcells. Let us emphasize that to obtain a relation of equivalency, one would need the same number of Degrees of Freedom (DoF) as number of subcells. Even if the choice of the cells subdivision may have an effect of the DG correction technique, for the following theoretical part it has no influence what so ever. The only constrain is that the projection matrix  $P_c$  further introduced in (2.11) has to be non-singular. Many subdivisions can be found in the literature for other methods relying on subgrid, as spectral volume methods for instance [47, 16] or subcell shock capturing technique as [20]. In Figure 1, three types of subdivision are displayed for both a triangular cell and a polygonal cell, in the case of a 4<sup>th</sup>-order DG scheme. Let us note that the one depicted in Figure 1(b) is the most widely used in triangle mesh subgrid techniques, and has the advantages to be invariant by rotation and can be generated for any order of accuracy. For numerical applications, this subdivision will be compared to a way simpler one, Figure 1(a), which has the benefits to be extremely simple to implement and where the subcells normals are nothing but the primal triangular cell ones.

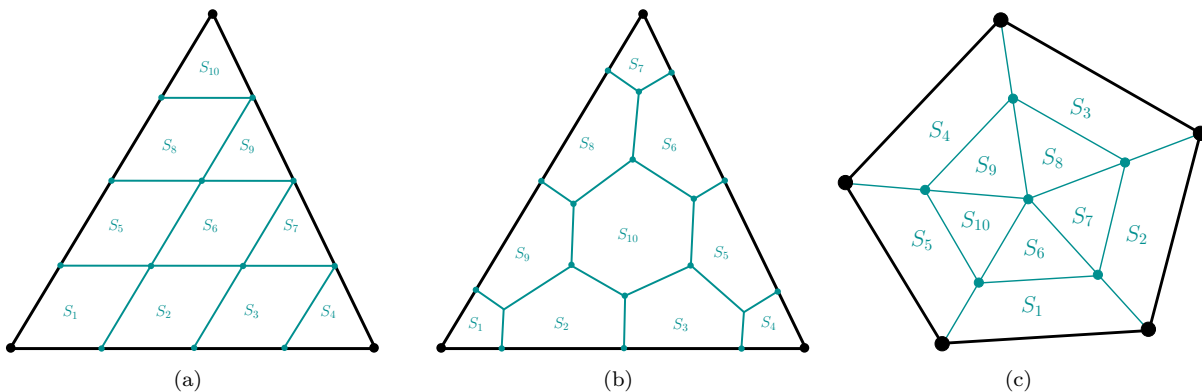


FIG. 1. Examples of subdivision for a  $\mathbb{P}^3$  DG scheme.

Even if only triangular grids are considered for numerical applications, let us emphasize that the following demonstration as well as the subcell correction technique presented in this paper are not limited to this case. Any grid made of generic polygonal cells can be considered. An example of a possible subdivision for  $\mathbb{P}^3$ -DG scheme is displayed in Figure 1(c). Curvilinear meshes with curvilinear cell subdivision could also be used, and will be the topic of a near future paper.

That being said, let us consider a cell  $\omega_c$  and its subdivision into  $N_k$  subcells  $S_m^c$ , for  $m \in [1, N_k]$ . For any function  $\psi \in L^2(\omega_c)$ , we define the corresponding subcell mean values, also referred to as submean values

$$(2.9) \quad \bar{\psi}_m^c = \frac{1}{|S_m^c|} \int_{S_m^c} \psi \, dV.$$

Let us now reformulated DG as a FV-like scheme provided the definition of the so-called reconstructed fluxes.

**2.1. Reconstructed flux through residuals.** Applying definition (2.9) to  $\partial_t u_h^c$  and by means of (2.2)

$$(2.10) \quad \frac{d\bar{u}_m^c}{dt} = \frac{1}{|S_m^c|} \sum_{q=1}^{N_k} \frac{d u_q^c}{dt} \int_{S_m^c} \sigma_q^c \, dV,$$

which can be put into into a matrix-vector form as  $\frac{d\bar{U}_c}{dt} = P_c \frac{dU_c}{dt}$  where the projection matrix  $P_c$ , which has to be invertible for the subdivision to be admissible, is defined as

$$(2.11) \quad (P_c)_{mp} = \frac{1}{|S_m^c|} \int_{S_m^c} \sigma_p^c \, dV.$$

The vector  $\bar{U}_c$  contains the cell submean values, *i.e.*  $(\bar{U}_c)_m = \bar{u}_m^c$ . Now, by means of DG scheme definition (2.7), it follows that  $\frac{d\bar{U}_c}{dt} = P_c M_c^{-1} \Phi_c$ . To express this relation as a FV-like scheme, we now introduce the DG reconstructed flux  $\widehat{F}_n$  such that

$$(2.12) \quad \frac{d\bar{u}_m^c}{dt} = -\frac{1}{|S_m^c|} \int_{\partial S_m^c} \widehat{F}_n \, dS.$$

By introducing  $\mathcal{V}_m^c$ , the set of face neighboring subcells of  $S_m^c$ , this last expression rewrites

$$(2.13) \quad \frac{d\bar{u}_m^c}{dt} = -\frac{1}{|S_m^c|} \sum_{S_p^v \in \mathcal{V}_m^c} \int_{f_{mp}^c} \widehat{F}_n \, dS,$$

where  $f_{mp}^c$  denotes the face between subcells  $S_m^c$  and  $S_p^v$ . Let us emphasize that  $S_p^v \in \mathcal{V}_m^c$  can either be inside cell  $\omega_c$  or in one of its neighbors  $\omega_v$  with  $v \neq c$ . This situation is displayed in Figure 2 where  $S_m^c$  would be colored red, while its faces neighboring subcells would be colored green. Now, similarly to what has been done in the 1D case, we impose that on the boundary of cell  $\omega_c$  the reconstructed flux coincides with the DG numerical flux

$$(2.14) \quad \widehat{F}_n|_{\partial\omega_c} = \mathcal{F}_n.$$

Expression (2.13) rewrites as

$$(2.15) \quad \frac{d\bar{u}_m^c}{dt} = -\frac{1}{|S_m^c|} \left( \sum_{S_p^v \in \widetilde{\mathcal{V}}_m^c} \int_{f_{mp}^c} \widehat{F}_n \, dS + \int_{\partial S_m^c \cap \partial\omega_c} \mathcal{F}_n \, dS \right),$$

where  $\widetilde{\mathcal{V}}_m^c$  stands for the set containing only the face neighboring subcells of  $S_m^c$  inside  $\omega_c$ . For now on, an orientation will be assigned to each face. Then, taking two subcells  $S_m^c$  and  $S_p^v$ , we introduce the sign function  $\varepsilon_{mp}^c$  defining the orientation of face  $f_{mp}^c$

$$(2.16) \quad \varepsilon_{mp}^c = \begin{cases} 1 & \text{if } S_p^v \in \mathcal{V}_m^c \text{ and } v \neq c, \\ 1 & \text{if } S_p^v \in \mathcal{V}_m^c \text{ with } v = c \text{ and } m < p, \\ -1 & \text{if } S_p^v \in \mathcal{V}_m^c \text{ with } v = c \text{ and } m > p, \\ 0 & \text{if } S_p^v \notin \mathcal{V}_m^c. \end{cases}$$

Obviously,  $\forall S_p^v \in \widetilde{\mathcal{V}}_m^c$ , we have  $\varepsilon_{pm}^c = -\varepsilon_{mp}^c$ . Now, unlike the 1D case where a pointwise definition of the reconstructed flux was given inside the cell, we make use here of a face integrated version of the high-order DG reconstructed flux. Indeed, for a face  $f_{mp}^c$ , let  $\widehat{F}_{mp}^c$  be defined as follows

$$(2.17) \quad \int_{f_{mp}^c} \widehat{F}_n \, dS = \varepsilon_{mp}^c \widehat{F}_{mp}^c.$$

Since the face orientation has been carried through  $\varepsilon_{mp}^c$ , the face integrated quantity  $\widehat{F}_{mp}^c$  is then continuous, *i.e.*  $\forall S_p^v \in \widetilde{\mathcal{V}}_m^c$ ,  $\widehat{F}_{pm}^c = \widehat{F}_{mp}^c$ . Now, denoting by  $N_f^c$  the number of subcells' faces inside  $\omega_c$ , meaning not belonging to  $\partial\omega_c$  (see Fig. 1 where those interior faces are colored green), let us introduce  $\widehat{F}_c \in \mathbb{R}^{N_f^c}$  the vector containing all the interior faces reconstructed fluxes. The subcell mean values governing equations (2.15) then yield the following system

$$(2.18) \quad -A_c \widehat{F}_c = D_c \frac{d\bar{U}_c}{dt} + B_c,$$

where  $A_c \in \mathcal{M}_{N_k \times N_f^c}$ , defined as  $(A_c)_{mp} = \varepsilon_{mp}^c$ , stands for the adjacency matrix, the subcells volume matrix  $D_c = \text{diag}(|S_1^c|, \dots, |S_{N_k}^c|)$  and  $(B_c)_m = \int_{\partial S_m^c \cap \partial \omega_c} \mathcal{F}_n \, dS$  carries the cell boundary contribution. Finally, we get

$$(2.19) \quad -A_c \widehat{\mathbf{F}}_c = D_c P_c M_c^{-1} \Phi_c + B_c.$$

To solve such system, and obtain an explicit expression of the reconstructed fluxes  $\widehat{\mathbf{F}}_c$ , we make use of the same graph Laplacian technique employed in [36] in a similar context. To do so, let us introduce  $L_c = A_c A_c^t$ , the graph Laplacian matrix which its generic coefficient writes as follows

$$(2.20) \quad (L_c)_{mp} = \begin{cases} |\widetilde{\mathcal{V}}_m^c| & \text{if } m = p, \\ -1 & \text{if } S_p^v \in \widetilde{\mathcal{V}}_m^c, \\ 0 & \text{otherwise.} \end{cases}$$

Such matrix is symmetric, and its rank and kernel are respectively  $N_k - 1$  and  $\text{span}\{\mathbf{1}\}$ . Due to these properties, and by means of  $L_c$  diagonalization through an orthogonal matrix, it follows that

$$(2.21) \quad L_c \mathcal{L}_c^{-1} = \mathcal{L}_c^{-1} L_c = I_{N_k} - \Pi.$$

Here,  $I_{N_k}$  is the identity matrix,  $\Pi = \frac{1}{N_k} (\mathbf{1} \otimes \mathbf{1})$  and  $\mathcal{L}_c^{-1}$  is the pseudo-inverse of  $L_c$ . Now, in order to define this latter matrix, let us introduce  $T = L_c + \lambda \Pi$ , for any constant  $\lambda \neq 0$ . Matrix  $T$  has the same eigendecomposition as  $L_c$ , apart from the null eigenvalue which has been substituted by  $\lambda$ . Such matrix is then invertible and it directly follows that  $\mathcal{L}_c^{-1} = (L_c + \lambda \Pi)^{-1} - \Pi/\lambda$ . In order to solve (2.19), let us first focus on the following problem

$$(2.22) \quad \left\{ \text{For a given } C \in \mathbb{R}^{N_k}, \text{ find } Y \in \mathbb{R}^{N_k} \text{ s.t. } L_c Y = C \right\}.$$

This problem admits solutions if and only if  $C \in \text{Im}(L_c)$ , condition we can recast into  $L_c \mathcal{L}_c^{-1} C = C$ . By means of relation (2.21), this previous condition can be reformulated as  $C \cdot \mathbf{1} = 0$ . Then, under this particular condition, problem (2.22) admits the following general solutions

$$(2.23) \quad Y = \mathcal{L}_c^{-1} C + (I_{N_k} - \mathcal{L}_c^{-1} L_c) Z, \quad \forall Z \in \mathbb{R}^{N_k}.$$

Seeking solutions of (2.19) we are not directly concerned with (2.22) but with a  $A_c X = C$  type of problem, with  $X \in \mathbb{R}^{N_f^c}$ . But by setting  $X = A_c^t Y$ , it follows again that the latter problem admits solutions if and only if  $C \cdot \mathbf{1} = 0$ , and it appears that this solution is in fact unique, as

$$(2.24) \quad \begin{aligned} X &= A_c^t \mathcal{L}_c^{-1} C + A_c^t \underbrace{(I_{N_k} - \mathcal{L}_c^{-1} L_c)}_{\Pi} Z, \quad \forall Z \in \mathbb{R}^{N_k}, \\ &= A_c^t \mathcal{L}_c^{-1} C, \end{aligned}$$

since  $A_c^t \Pi = \mathbf{0}_{N_k}$ . Following this procedure, we are now able to exhibit the following unique definition of the DG reconstructed flux

$$(2.25) \quad \widehat{\mathbf{F}}_c = -A_c^t \mathcal{L}_c^{-1} (D_c P_c M_c^{-1} \Phi_c + B_c).$$

This unique solution does exist since  $(D_c P_c M_c^{-1} \Phi_c + B_c) \cdot \mathbf{1} = 0$  (see Appendix A.1).

**REMARK 1.** *In the definition of the reconstructed flux through DG residual, (2.25), the only time dependent terms are  $\Phi_c$ , the residual which is directly available in any DG code, and the boundary contribution  $B_c$ . All the other terms can be evaluated initially, once and for all. Furthermore, if all mesh cells have the same structure, as triangles for example, then by means of a mapping to a referential element, the projection matrix  $P_c$ , the adjacency matrix  $A_c$  and the generalized inverse of the graph Laplacian matrix  $\mathcal{L}_c^{-1}$  do not depend on the cell under consideration, but only on the order of approximation and on the choice of the subdivision.*

Once we have computed the reconstructed flux  $\widehat{F}_c$ , we can simply recover the polynomial solution governing equation as follows

$$(2.26) \quad \frac{dU_c}{dt} = -P_c^{-1} D_c^{-1} (A_c \widehat{F}_c + B_c).$$

Now, for a deeper understanding of the reconstructed flux  $\widehat{F}_c$  defined in (2.25), let us seek its relation with the interior flux  $\mathbf{F}(u_h^c)$  and the DG numerical flux  $\mathcal{F}_n$ .

**2.2. Reconstructed flux through fluxes.** In this section, we aim at expressing the reconstructing fluxes only through the interior flux  $\mathbf{F}(u_h^c)$  and a correction term taking into account the jump at the cell boundaries, similarly to what is done in SBP operator with SAT boundary treatment [13, 14] or CPR schemes [21, 45]. To do so, we will not make use of the definition of DG scheme through residual (2.7) but through fluxes (2.3). The first step is to substitute in (2.3) the interior flux  $\mathbf{F}(u_h^c)$  by  $\mathbf{F}_h^c$  its  $L_2$  projection onto  $(\mathbb{P}^{k+1}(\omega_c))^2$ . If one uses nodal DG or any collocation of the interior flux, this step can obviously be skipped. Performing a second integration by parts leads to the so-called strong form of DG scheme

$$(2.27) \quad \int_{\omega_c} \frac{\partial u_h^c}{\partial t} \psi \, dV = - \int_{\omega_c} \psi \nabla_x \cdot \mathbf{F}_h^c \, dV + \int_{\partial\omega_c} \psi (\mathbf{F}_h^c \cdot \mathbf{n} - \mathcal{F}_n) \, dS, \quad \forall \psi \in \mathbb{P}^k(\omega_c).$$

Similarly to what has been done in [42] for the one-dimensional case, let us introduce the  $N_k$  sub-resolution basis functions  $\{\phi_m\}_m$ . These particular basis functions of  $\mathbb{P}^k(\omega_c)$ , which can be seen as the  $L_2$  projection of the subcell indicator functions  $\mathbb{1}_{S_m^c}(\mathbf{x})$  onto  $\mathbb{P}^k(\omega_c)$ , are defined such that  $\forall \psi \in \mathbb{P}^k(\omega_c)$

$$(2.28) \quad \int_{\omega_c} \phi_m \psi \, dV = \int_{S_m^c} \psi \, dV, \quad \forall m = 1, \dots, N_k.$$

Because equation (2.27) holds for any polynomial  $\psi$  of degree  $k$ , let us substitute  $\phi_m$  for  $\psi$  in DG schemes. Then, through the sub-resolution property (2.28), one can recast equation (2.27) into

$$(2.29) \quad |S_m^c| \frac{d\bar{u}_m^c}{dt} = - \int_{\partial S_m^c} \mathbf{F}_h^c \cdot \mathbf{n} \, dS + \int_{\partial\omega_c} \phi_m (\mathbf{F}_h^c \cdot \mathbf{n} - \mathcal{F}_n) \, dS.$$

The use of reconstructed flux definition (2.12) directly leads to the following relation

$$(2.30) \quad \int_{\partial S_m^c} \widehat{F}_n \, dS = \int_{\partial S_m^c} \mathbf{F}_h^c \cdot \mathbf{n} \, dS - \int_{\partial\omega_c} \phi_m (\mathbf{F}_h^c \cdot \mathbf{n} - \mathcal{F}_n) \, dS.$$

One can see how the reconstructed flux is connected to the interior polynomial flux and the jump at the cell interface between the interior flux and the DG numerical flux. Similarly to (2.17), let  $F_{mp}$  be the face integrated value of the polynomial interior flux

$$(2.31) \quad \int_{f_{mp}^c} \mathbf{F}_h^c \cdot \mathbf{n} \, dS = \varepsilon_{mp}^c F_{mp}.$$

Then, if  $F_c$  is the vector containing all the interior faces fluxes, one gets  $A_c \widehat{F}_c = A_c F_c - G_c$ , where  $G_c$  contains the boundary contribution as

$$(2.32) \quad (G_c)_m = \int_{\partial\omega_c} (\phi_m - \mathbb{1}_{\partial S_m^c}) (\mathbf{F}_h^c \cdot \mathbf{n} - \mathcal{F}_n) \, dS.$$

The term  $\mathbb{1}_{\partial S_m^c}$  comes from assumption (2.14) where the reconstructed flux is set to be the DG numerical flux on the primal cell boundary. Finally, by means of the same graph Laplacian technique used previously, we are able to express the reconstructed flux through the interior flux and a boundary correction term

$$(2.33) \quad \widehat{F}_c = F_c - A_c^t \mathcal{L}_c^{-1} G_c.$$



REMARK 2. We can rewrite (2.33) as  $\widehat{F}_c = F_c - E(\mathbf{F}_h^c \cdot \mathbf{n} - \mathcal{F}_n)$ , where  $E(\cdot)$  would be a correction function taking into account the jump between the polynomial flux and the numerical flux on the cell boundary. This permits to demonstrate once more that in DG schemes the numerical diffusion deriving from the jump in term of flux across the cell interfaces is distributing elsewhere inside the cell, here at the subcells faces. The sub-resolution basis functions act as weighted functions in the diffusion distribution.

Let us note that another choice in the correction term function  $E(\cdot)$  would lead to a different scheme. For instance, setting  $E(\cdot) = 0$  leads to spectral volume scheme of Z.J. Wang [46].

Both definitions (2.25) and (2.33) are perfectly equivalent. While the definition of DG reconstructed fluxes (2.33) enables a better comprehension of how those fluxes are related to the interior flux  $\mathbf{F}(u_h^c)$  and how the numerical diffusion deriving from the jump at the interface is distributed inside the cell, the definition of reconstructed fluxes through residual (2.25) is a lot easier to implement and do not require the definition of sub-resolution basis functions.

**3. A posteriori local subcell correction.** The previous reformulations of DG scheme into subcell FV-like scheme through the definition of reconstructed fluxes enable us to construct our *a posteriori* local subcell correction (APLSC). In few words, the reconstructed fluxes  $\widehat{F}_{mp}$  will be modified in a robust way in subcells where the original DG scheme has failed. Let us emphasize that the popularity and number of subcell correction techniques have extensively grown these past years, see [20, 40, 10, 42, 29, 34, 19]. Those shock capturing and property preserving methods generally rely on a low-order scheme combined, at the subcell level, with a high-order one. However, let us note that in most cases, all the subcells contained in a cell will be impacted if something bad happened somewhere in the cell. In [42], we have introduced, for the one-dimensional case, a new technique permitting to correct the solution in a subcell without modifying the solution elsewhere. This particular feature allow to retain the very precise subcell resolution of high-order DG schemes. The present paper aims at presenting the two-dimensional version of this correction. Let us emphasize that, up to our knowledge, this is the only technique working on totally unstructured grids and permitting the modification of the scheme, locally at the subcell level, without impacting the solution everywhere in the cell under consideration.

Let us mention that until now, only the semi-discrete version of schemes and their corresponding analysis were presented. To achieve high-accuracy in time, we make use of SSP Runge-Kutta time integration method [38]. But, in the light of the fact that these multistage time integration methods write as convex combinations of first-order forward Euler scheme, the correction DG procedure will be presented for the simple case of this latter time numerical scheme, for sake of simplicity.

Let us now introduce the correction procedure. First, we assume that at time  $t^n$  the numerical solution  $u_h^n$  is satisfactory in the sens that, on any cell  $\omega_c$ , the subcell mean values are admissible regarding some criteria yet to be defined. Then, we compute  $u_h^{n+1}$  a candidate solution through the uncorrected DG scheme. The third step is then crucial. Indeed, we then have to check if the new uncorrected solution is admissible. If it is the case, we can go further in time without any special treatment. Otherwise we have to return to time  $t^n$  and recompute the solution locally by means of a more robust scheme. This step is crucial in the sens that it will tell us if and where a new computation would be required.

**3.1. Troubled zone detector.** Regarding troubled zone detectors, we simply extend to the two-dimensional unstructured case the ones used in [11, 42]. In those work, two detection criteria were mainly used, namely one ensuring the physical admissibility of the numerical solution (PAD) and another addressing the apparition of spurious oscillations (NAD). Let us then recall these two criteria.

*Physical admissibility detection (PAD).*

- Check if the different submean values  $\bar{u}_m^{c, n+1}$  lie in a chosen convex physical admissible set (maximum principle for SCL, positivity of the pressure and density for Euler, ...). Entropy stability can be added to this admissible set.
- Check if there is any NaN values

Those are the minimum requirements if one wants to enforce code robustness. Now, in order to tackle the issue of spurious oscillations, we make use of a local maximum principle. Indeed, through the respect of the CFL, the solution in cell  $\omega_c$  at time  $t^{n+1}$  has to remain in the bounds of the solution at the previous time step

$t^n$  wherein  $\omega_c$  and its first face neighbors. This condition is reformulated in the following detection criterion.

*Numerical admissibility detection (NAD).*

- Check if the following Discrete Maximum Principle (DMP) on submean values is ensured:

$$\min_{v \in \mathcal{N}(S_m^c)} (\bar{u}_v^n) \leq \bar{u}_m^{c,n+1} \leq \max_{v \in \mathcal{N}(S_m^c)} (\bar{u}_v^n),$$

where  $\mathcal{N}(S_m^c)$  is some set of  $S_m^c$  neighboring subcells, including subcell  $S_m^c$ .

REMARK 3. *The smaller the set  $\mathcal{N}(S_m^c)$  is, the more constraining the NAD criterion will be, meaning more subcells will be considered as problematic. Different sets will be considered for linear and non-linear problems.*

Let us enlighten that because the NAD criterion relies on a maximum principle based on subcell mean values, one has to relax it to preserve scheme accuracy in the presence of smooth extrema. To do so, we make use of the two-dimensional version of the one introduced in the 1D case, [42].

*Detection of smooth extrema.* This smooth extrema detection criterion is based on an idea present in different limitations, as the hierarchical slope limiter [28]. In this work, the numerical solution is supposed to exhibit a smooth extrema if at least the linearized version of the numerical solution spatial derivatives, *i.e.*

$$(3.1a) \quad \begin{cases} v_x^c(\mathbf{x}) = \overline{\partial_x u_h^{c,n+1}} + \overline{\nabla_x (\partial_x u_h^{c,n+1})} \cdot (\mathbf{x} - \mathbf{x}_c), \\ v_y^c(\mathbf{x}) = \overline{\partial_y u_h^{c,n+1}} + \overline{\nabla_x (\partial_y u_h^{c,n+1})} \cdot (\mathbf{x} - \mathbf{x}_c), \end{cases}$$

$$(3.1b) \quad \begin{cases} v_x^c(\mathbf{x}) = \overline{\partial_x u_h^{c,n+1}} + \overline{\nabla_x (\partial_x u_h^{c,n+1})} \cdot (\mathbf{x} - \mathbf{x}_c), \\ v_y^c(\mathbf{x}) = \overline{\partial_y u_h^{c,n+1}} + \overline{\nabla_x (\partial_y u_h^{c,n+1})} \cdot (\mathbf{x} - \mathbf{x}_c), \end{cases}$$

present a monotonous profile. In (3.1),  $\mathbf{x}_c$  denotes the centroid of cell  $\omega_c$ , while  $\overline{\partial_{x \setminus y} u_h^{c,n+1}}$  and  $\overline{\nabla_x (\partial_{x \setminus y} u_h^{c,n+1})}$  are nothing but the averaged values on  $\omega_c$  of the successive partial derivatives of  $u_h^c$ . In practice, the NAD relaxation used here works as a vertex-based limiter on  $v_{x \setminus y}^c$ . Due to their linearity, functions  $v_{x \setminus y}^c$  attain their extrema at the vertices  $\mathbf{x}_p \in \mathcal{P}_c$ , where  $\mathcal{P}_c$  stands for the set of vertices of cell  $\omega_c$ . Then, we consider that the exact weak solution underlying the numerical solution  $u_h$  presents a smooth profile in cell  $\omega_c$  if, for any vertex  $\mathbf{x}_p \in \mathcal{P}_c$ , the linearized spatial derivative functions ensure the following constraints

$$(3.2) \quad v_{x,p}^{\min} \leq v_x^c(\mathbf{x}_p) \leq v_{x,p}^{\max} \quad \text{and} \quad v_{y,p}^{\min} \leq v_y^c(\mathbf{x}_p) \leq v_{y,p}^{\max},$$

where  $v_{x \setminus y,p}^{\min} = \min_{v \in \mathcal{V}_p} v_{x \setminus y}^c(\mathbf{x}_p)$  and  $v_{x \setminus y,p}^{\max} = \max_{v \in \mathcal{V}_p} v_{x \setminus y}^c(\mathbf{x}_p)$ . Here,  $\mathcal{V}_p$  represents the set of cells that share  $\mathbf{x}_p$  as a vertex, *i.e.*  $\omega_v \in \mathcal{V}_p \implies \mathbf{x}_p \in \mathcal{P}_v$ . Practically, if for any vertex  $\mathbf{x}_p \in \mathcal{P}_c$ , conditions (3.2) are ensured, we then consider that the numerical solution presents a smooth profile on cell  $\omega_c$ . In this particular case, the NAD criterion is relaxed allowing the preservation of smooth extrema along with the order of accuracy for smooth problems, see Section 4.

REMARK 4. *We have presented here the detection based on the linearized first-order derivatives of the solution. This would work for any higher order derivative. Furthermore, such relaxation procedure can also be applied at the subcell level to also preserve smooth extrema even within a cell at a smaller length scale. This would be useful in the context of coarse grids. Actually, because the subcell smooth extrema relaxation technique works well for both coarse and fine meshes, this will be the procedure used for the numerical applications, Section 4.*

**3.2. Correction.** Now that we have detailed the troubled subcell detector, the correction procedure will be presented. The very simple idea that forms the basis of the original correction procedure introduced in [42] is the following: if the uncorrected DG scheme has produced a numerical solution  $u_h^{c,n+1}$  on cell  $\omega_c$ , which is not admissible in subcell  $S_m^c$  in regards to the detection criteria presented previously, the subcell mean value  $\bar{u}_m^{c,n+1}$  will be recomputed by means of a more robust scheme. To do so, and because uncorrected DG scheme is equivalent to subcell finite volume scheme with the appropriate high-order reconstructed fluxes, we substitute on the boundaries of subcell  $S_m^c$  the high-order reconstructed fluxes with some first-order finite volume numerical fluxes. Submean value  $\bar{u}_m^{c,n+1}$  will then be recomputed by means of a simple and robust

first-order finite volume scheme. To this end, we introduce  $\widetilde{F}_{mp}$  the corrected reconstructed fluxes so that

$$\begin{cases} \widetilde{F}_{mp} = \varepsilon_{mp}^c l_{mp}^c \mathcal{F}(\bar{u}_m^{c,n}, \bar{u}_p^{v,n}, \mathbf{n}_{mp}) & \text{if } S_m^c \text{ or } S_p^v \in \mathcal{V}_m^c \text{ is either marked,} \\ \widetilde{F}_{mp} = \widehat{F}_{mp} & \text{otherwise,} \end{cases}$$

where  $l_{mp}^c$  is the length of face  $f_{mp}^c$ . Through those corrected reconstructed fluxes, we recompute the submean values for tagged subcells and their first neighboring subcells as

$$(3.3) \quad \bar{u}_m^{c,n+1} = \bar{u}_m^{c,n} - \frac{\Delta t}{|S_m^c|} \sum_{S_p^v \in \mathcal{V}_m^c} \varepsilon_{mp}^c \widetilde{F}_{mp}.$$

This concept is depicted in Figure 2, where the troubled subcells are colored red.

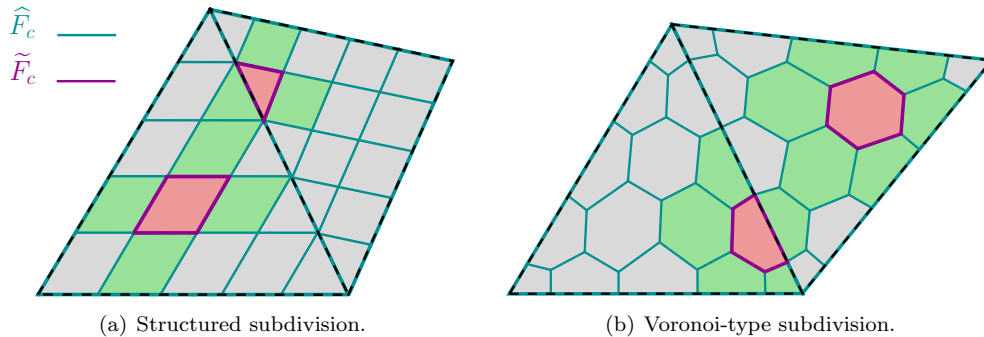


FIG. 2. Original correction of the DG reconstructed flux.

Because it is of fundamental importance to preserve scheme conservation, the first face neighboring subcells, colored green in Figure 2, have to be also recomputed since we have modify the reconstructed fluxes on the boundary of the troubled subcell. The submean values of the neighboring subcells are then computed through a FV-like scheme with first-order numerical flux on one or more faces and high-order reconstructed fluxes on the remaining interfaces. For the remaining subcells, colored gray in Figure 2, because the corresponding reconstructed fluxes have not been modified, there is no need to recompute them. The corresponding submean values are hence the values obtained through the uncorrected DG scheme. It is clear that through this technique, the DG solution will only be affected at the subcell scale. Furthermore, the corrected scheme is conservative at the subcell level by construction.

REMARK 5. *Let us emphasize that since this a posteriori correction is based on first-order finite volume scheme, maximum principle or positivity preservation in the case of systems are enforced by construction. However, it is crucial to note that to start from a discrete representation of the initial datum that respects its bounds (or ensures positivity), the initialization has to be carried out through the integration of the initial solution on the subcells to obtain the different submean values, and then through the projection matrix  $P_c$  compute the corresponding polynomial moments of the solution. Traditionally, in DG schemes the initialization is handled by either assigning at the solution points the initial datum value to the numerical solution, or by a  $L_2$  projection onto  $\mathbb{P}^k$ . In either of these procedures, nothing ensures that the submean values would respect the bounds of the initial datum.*

REMARK 6. *Let us highlight the fact that the admissible properties are enforced on the subcell mean values of the solution, and not its polynomial representation. Consequently, the numerical solution may present some non-admissible values, for instance at quadrature or boundary points. The "Check if there is any NaN values" present in the PAD criterion is here for that matter. Indeed, if for instance the speed of sound becomes negative in the hydrodynamics case or some square root of negative quantity is computed, then the whole solution inside the cell will exhibit extreme values and possibly some NaN. However, no additional slope limiter is required here, as if such situation presents itself, the a posteriori correction loop will automatically end up by correcting all the subcells inside the problematic cell.*

While this procedure has already proved its efficiency in one-dimensional configuration, see [42], for non-linear problems using very high-order schemes on coarse grids, the numerical solution has showed to remain slightly oscillatory at the subcell level. To overcome this issue, we were artificially enlarging the stencil to correct by also marking, for  $k \geq 3$ , the first face neighboring subcells of a troubled subcell. Now, in this 2D frame, a modified algorithm making use of convex combination of DG reconstructed fluxes and first-order FV fluxes for admissible subcells in the vicinity of troubled areas will be presented and has proven to produce better results than the original correction, in this context of very high-orders and coarse meshes. Indeed, to avoid to locally jump, in two subcells scale, from a very precise approximation to a robust but very low accurate first-order representation, we now introduce a new definition for the corrected reconstructed fluxes by means of convex combination between DG reconstructed fluxes and first-order FV fluxes as follows

$$(3.4) \quad \widetilde{F}_{mp} = \theta_{mp} \varepsilon_{mp}^c I_{mp}^c \mathcal{F}(\bar{u}_m^{c,n}, \bar{u}_p^{v,n}, \mathbf{n}_{mp}) + (1 - \theta_{mp}) \widehat{F}_{mp},$$

where  $\theta_{mp}$  is a function of the distance to a non-admissible subcell. Obviously, the farther from the troubled subcell we are, the less of first-order FV we want. Many smoothness indicators exist in the literature, see for instance [41, 17, 4, 26, 35, 1], and could be used in this context to determine a relevant bending coefficient. Another way around could be to adopt flux limiting or Flux-Corrected transport (FCT) approach, see [3, 51, 31, 30], to find the proper blending ensuring a discrete maximum principle. To remain as simple as possible, and because this paper is mainly concerned with the introduction of the DG-FV equivalency and the basic correction principle on 2D unstructured grids, we make use here of a very naive procedure. The principle is the following, for marked subcells detected through the troubled subcell indicator, first-order FV numerical flux is used on its boundaries, *i.e.*  $\theta_{mp} = 1$  in (3.4). Then, for its face neighboring subcells  $S_p^v \in \mathcal{V}_m^c$ , convex combination (3.4) with  $\theta_{mp} = \frac{3}{4}$  would be chosen for their remaining boundaries. Now, introducing  $\widetilde{\mathcal{V}}_m^c$  the set containing both the face and node neighboring subcells, fluxes on faces of  $S_p^v \in \widetilde{\mathcal{V}}_m^c \setminus \mathcal{V}_m^c$  which yet have not been corrected will be defined through at the blending coefficient  $\theta_{mp} = \frac{1}{2}$ . Finally, the remaining face neighbors of  $S_p^v \in \widetilde{\mathcal{V}}_m^c$  will see their remaining boundaries associated with a numerical flux calculated through (3.4) with  $\theta_{mp} = \frac{1}{4}$ . This naive technique is displayed in Figure 3.

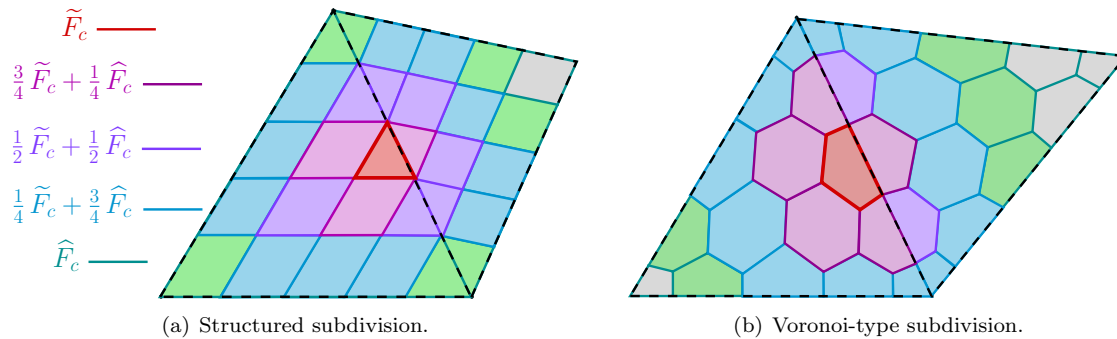


FIG. 3. New correction of the DG reconstructed flux.

One could think that doing so, we have substantially enlarged the stencil of subcells to be corrected, and thus reduced the simulation code efficiency. However, in practice it is generally not the case, as one can see in Figure 4. Moreover, the computational time will also be generally slightly reduced. Emphasizing that the correction is done in an *a posteriori* fashion and thus has to be potentially iterated multiple times at a time step to reach an admissible solution, if the stiff original correction introduces small oscillations at the subcell scale, the stencil will be automatically enlarged along with an increase of the computational time, through the correction iteration. With the new correction principle, only one iteration is generally needed. So even if at first, the set containing marked subcells is smaller than in the original correction, it will end up with approximately the same size as this new approach, and required more iterations and thus more computational effort. To assess this matter and compare both corrections, we make use of the Burgers equation, defined through (2.1a) and flux function  $\mathbf{F}(u) = \frac{1}{2} (u^2, u^2)^t$ , with the smooth initial

solution  $u_0(\mathbf{x}) = \sin(2\pi(x+y))$ . The domain is chosen as the unit square  $[0, 1]^2$  with periodic boundary condition. Through time, the exact solution will exhibit two stationary shocks along the lines defined by  $(\mathbf{x} \in [0, 1]^2, x+y=0.5)$  and  $(\mathbf{x} \in [0, 1]^2, x+y=1.5)$ . We run this test case until  $t = 0.5$  with a sixth-order DG scheme, on a quite coarse unstructured grid made of 576 cells, with both corrections. To compare the two approaches, let us first display the marked subcells to be corrected. In Figures 4, we color all the subcells that have been corrected in the different Runge-Kutta steps during the last time step.

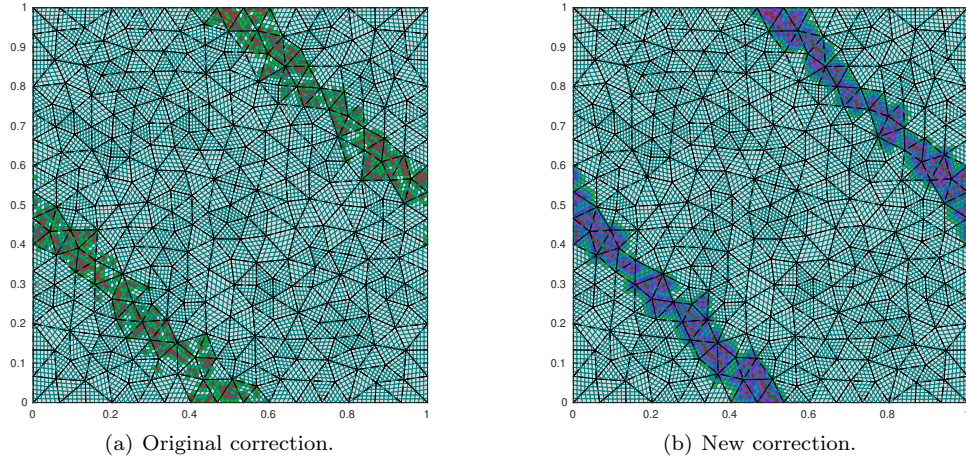


FIG. 4. Comparison between original and new correction procedure: corrected subcells.

Firstly, one can note that both corrections, through the NAD criterion, have accurately capture the discontinuities, as the subcells to be corrected remain in a small vicinity of the shocks. Secondly, we can also observe how these corrections operate locally at the subcell scale. Finally, as depicted in Figure 4, the number of subcells to be corrected remains approximately the same with both approaches. Actually, with the original correction on average 10% of the total number of subcells have to be corrected through the computation, while 14% of the subcells with the new approach. But, even if indeed the number of subcells to be corrected has slightly increased through this new procedure, if we compare computational efficiency, it took 2 minutes and 22 seconds to the code with the original correction, and 2 minutes 15 seconds with the new one. As said previously, it comes from the fact that more correction iterations are required with the original procedure. For this calculation, on average 2.86 iterations are needed with the original approach when the correction has been triggered, for a maximum of 6 iterations, while only 1.46 iterations with the new approach, with a maximum of 3 iterations during the whole calculation. One can also observe that a lot less subcells are corrected through a purely first-order FV with the new correction than with the original one, which enables even more the preservation of the high accurate subcell resolution of DG schemes. In Figure 4(a), one can see the two types of subcells to be corrected, meaning the troubled subcells colored red and their face neighbors colored green which have to be recalculated to preserve scheme conservation. In Figure 4(b), there are five types of subcells to be corrected, from the troubled subcells to be recomputed through a first-order FV scheme to the magenta, purple and blue ones which are recalculated through a convex combination between first-order FV numerical flux and the high-order DG reconstructed flux, with different weights. The green ones are the face neighbors of the previous subcells to be recomputed to preserve scheme conservation.

While the previous results showed how those two approaches work and do have a quasi-equivalent computation cost, let us now compare the approximated solutions obtained and assess the benefit of the new correction. In Figure 5, the subcell mean values obtained by means of a sixth-order DG scheme corrected through the original and the new procedures are displayed. Comparing Figures 5(a) and 5(b), we can see that the new correction has improved the accuracy of the scheme by a sharper and less oscillatory representation of the shock, despite the very coarseness of the mesh used. Finally, we plot in Figure 6 for both corrections, as well as for the exact entropic solution, the subcell mean values  $\bar{u}_m^c$  versus  $r_m^c = x_c + y_c - 1$ , where  $(x_c, y_c)$  stands for the barycenter of subcell  $S_m^c$ . These results demonstrates once more how the new correction leads to a

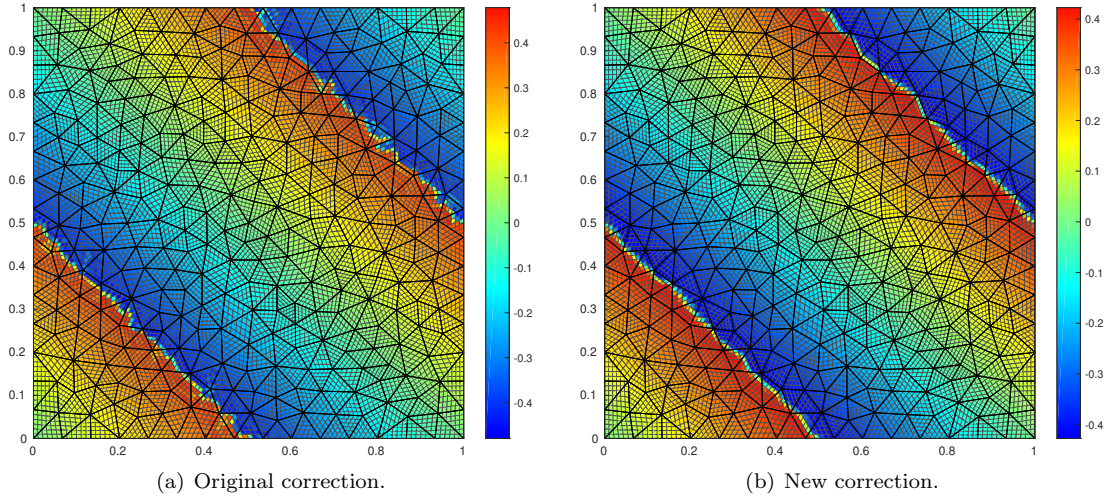


FIG. 5. Comparison between original and new correction procedure: subcell mean values.

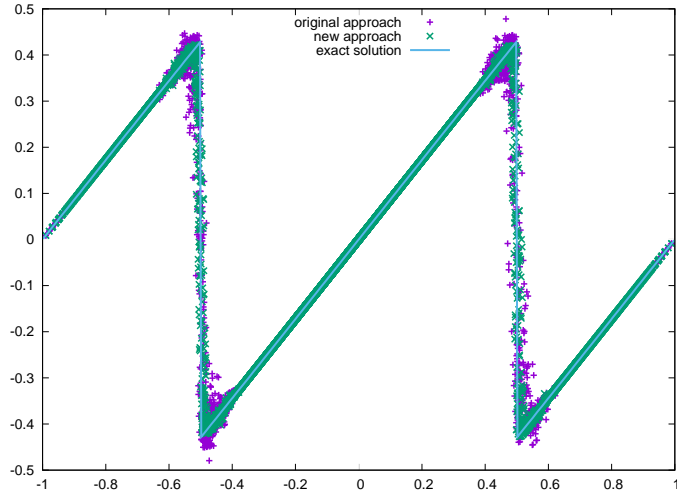


FIG. 6. Comparison between original and new correction procedure: submean values versus  $(x + y - 1)$  coordinate.

better representation of the solution. For that reason, the new correction based on a convex combination of first-order FV numerical flux and high-order DG reconstructed flux, with decreasing weight in the vicinity of a troubled subcell, will be adopted for the numerical applications presented in the next section.

**4. Numerical results.** In this numerical results section, we make use of several widely addressed and challenging test cases to demonstrate the performance and robustness of this *a posteriori* local subcell correction of discontinuous Galerkin schemes. In all following test cases, the simple case of global Lax-Friedrichs numerical flux will be used for both the DG scheme and the first-order finite volume reconstructed flux correction. Regarding the cell decomposition into subcells, while this has no impact on the reformulation of DG schemes into subcell finite volume method, see Section 2, it may have a slight impact on the results obtained by means of the present correction. In the one-dimensional case, see [42], it has been demonstrated that the use of a non-uniform subdivision, for instance by means of the Gauss-Lobatto points, leads to better results compared to a uniform subdivision. In this two-dimensional framework the two types of subdivisions introduced respectively in Figure 1(a) and 1(b) will be experimented.

Regarding the time integration, we make use of the classical third-order SSP Runge-Kutta scheme, see for instance [38]. As the correction described earlier combines both DG scheme on the primal cells  $\omega_c$  and FV

scheme on the subcells  $S_m^c$ , the time step is computed adaptively using the following CFL condition

$$(4.1) \quad \Delta t = \frac{1}{\gamma} \min \left( \frac{d_c}{2k+1}, \min_m d_m^c \right),$$

where  $\gamma = \max_{c,m} (\|\mathbf{F}'(\bar{u}_m^c)\|_2)$ , and where the cell and subcell characteristic lengths  $d_c$  and  $d_m^c$  are defined as

$$(4.2) \quad d_c = \frac{|\omega_c|}{\sum_{\omega_v \in \mathcal{V}_c} l_{cv}} \quad \text{and} \quad d_m^c = \frac{|S_m^c|}{\sum_{S_p^v \in \mathcal{V}_m^c} l_{mp}^c}.$$

We recall that  $l_{mp}^c$  stands for the length of face  $f_{mp}^c$  separating subcell  $S_m^c$  and its neighbor  $S_p^v$ , while  $l_{cv}$  is the length of the interface between cell  $\omega_c$  and its neighbor  $\omega_v$ . Let us note that in cases where we compute rates of convergence, a time step  $\Delta t \leq \min_c d_c^{\frac{k+1}{3}}$  is used in order to make the time error negligible in comparison to the spatial discretization error. For each result, the solution subcell mean values are displayed. Consequently, for a mesh made of  $N_c$  cells, the numerical solution will be represented on  $N_k N_c$  subcells, where  $N_k = (k+1)(k+2)/2$ .

**4.1. Linear case.** Let us first assess the performance and accuracy of the present *A Posteriori* Local Subcell Corrected DG (APLSC-DG) scheme in the case of 2D linear conservation laws. In Section 3.1, when the NAD criterion based on a discrete maximum principle has been introduced, we did not specify the set  $\mathcal{N}(S_m^c)$  of subcell  $S_m^c$  characterizing the DMP. For linear problems, we make use of a cell-wise DMP, meaning  $\mathcal{N}(S_m^c)$  will be constituted by all the subcells of cell  $\omega_c$ , as well as the subcells of  $\omega_v$  the face neighboring cells of  $\omega_c$ . By means of notations previously introduced, this definition can be rewritten as

$$(4.3) \quad \mathcal{N}(S_m^c) = \left\{ S_q^v; \forall \omega_v \in \mathcal{V}_c \cup \{\omega_c\}, \forall q \in \llbracket 1, N_k \rrbracket \right\}.$$

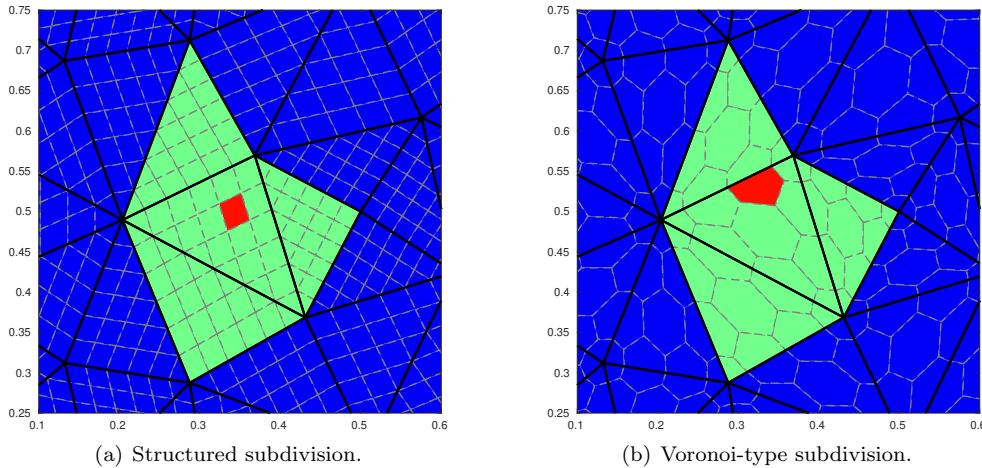


FIG. 7. Neighboring subcells set  $\mathcal{N}(S_m^c)$  for the NAD criterion in the linear case: subcell  $S_m^c$  is colored red while the subcells in  $\mathcal{N}(S_m^c)$  are colored green.

This particular set is depicted in Figure 7, for both the simple structured subdivision as well as the polygonal Voronoi-type one. In this figures, the subcell  $S_m^c$  under consideration would be colored red, while the subcells constituting  $\mathcal{N}(S_m^c)$  would be colored green. Let us emphasize that subcell  $S_m^c$  is also part of  $\mathcal{N}(S_m^c)$ .

**4.1.1. Linear advection.** To display the efficiency of DG schemes plus correction, let us first assess how the APLSC behaves in the linear advection case. To this end, we consider the following equation

$$\begin{cases} \partial_t u(\mathbf{x}, t) + \mathbf{A}(\mathbf{x}) \cdot \nabla_{\mathbf{x}} u(\mathbf{x}, t) = 0, & (\mathbf{x}, t) \in [0, 1]^2 \times [0, T], \\ u(\mathbf{x}, 0) = u_0(\mathbf{x}), & \mathbf{x} \in [0, 1]^2, \end{cases}$$

**Linear advection of a smooth signal.** We start from a smooth initial datum  $u_0(x, y) = \sin(2\pi(x+y))$ , and consider periodic boundary conditions. We assess the scheme accuracy after one period, namely at time  $t = 1$ . In Figure 8, the numerical solution of the APLSC of sixth-order DG scheme, obtained on grid made of only 100 cells, is depicted. Let us note that the correction procedure does not activate in this case, which

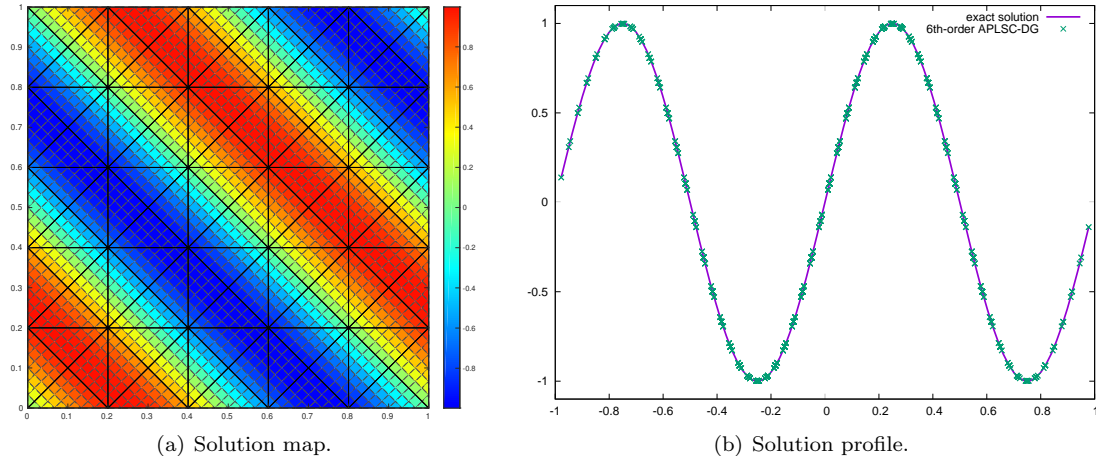


FIG. 8. Linear advection a smooth signal with a 6th corrected DG scheme, on a  $5 \times 5 \times 4$  grid after 1 period.

proves that the relaxation criterion on smooth extrema works properly. The rates of convergence are gathered in Table 1 and do exhibit a convergence to six.

	$L_1$		$L_2$		$L_\infty$	
$h$	$E_{L_1}^h$	$q_{L_1}^h$	$E_{L_2}^h$	$q_{L_2}^h$	$E_{L_\infty}^h$	$q_{L_\infty}^h$
$\frac{1}{10}$	1.62E-7	6.00	1.81E-7	6.00	3.98E-7	5.96
$\frac{1}{20}$	2.53E-9	5.97	2.82E-9	5.96	6.38E-9	5.41
$\frac{1}{40}$	4.03E-11	-	4.52E-11	-	1.50E-10	-

TABLE 1

Convergence rates for the linear advection case for a 6th-order APLSC-DG scheme

**Linear advection of a crenel signal.** To assess the efficiency of the correction presented in the presence of discontinuity, let us start with the simple case of the advection of a crenel signal, where the initial solution  $u_0$  is defined as follows

$$(4.5) \quad u_0(\mathbf{x}) = \begin{cases} 1 & \text{if } (x+y) \in [1/4, 1/2] \cup [5/4, 3/2], \\ 0 & \text{if } (x+y) \in [3/4, 1] \cup [7/4, 2], \\ 1/2 & \text{otherwise.} \end{cases}$$

In Figure 9, we compare the uncorrected and corrected versions of the 6th-order DG scheme on an unstructured grid made of 576 cells, after one period. One can see that in both cases, the numerical solution obtained is very accurate, but the one obtained through the APLSC-DG method respects the maximum principle as the final solution remains in the bounds of the initial one. In Figure 10, the submean values versus  $(x+y-1)$  coordinate are compared for both the uncorrected and the APLSC-DG schemes. We note how the spurious oscillations in the vicinity of the discontinuities have been removed, while preserving the very precise resolution of the 6th-order DG scheme. In the previous examples, the simple uniform structured cell subdivision has been used. Let us emphasize that the cell subdivision does not theoretically impact the equivalency between DG scheme and subcell FV-like scheme defined through the reconstructed fluxes introduced previously. This choice, if it does, may only influence the corrected scheme. To illustrate such statement, let us run the linear advection of the crenel signal on five periods for different types of cell subdivision, see Figures 11.



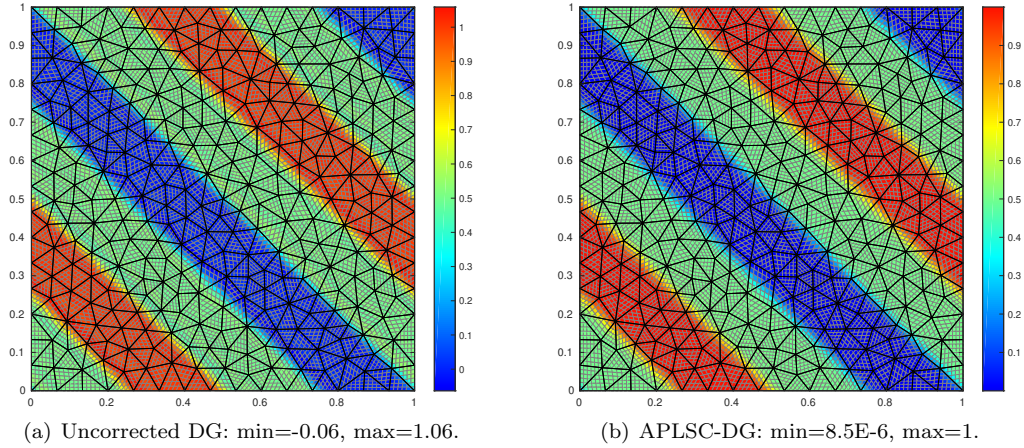


FIG. 9. 6th-order DG solutions for the linear advection of a crenel signal on 576 cells after one period.

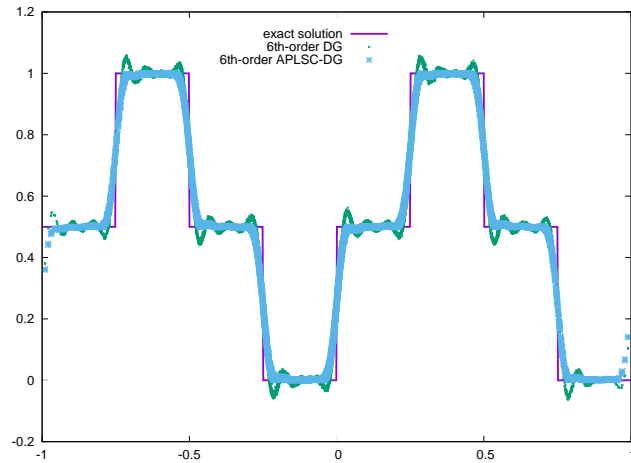


FIG. 10. 6th-order solutions for the crenel advection case on 576 cells: submean values versus  $(x + y - 1)$  coordinate.

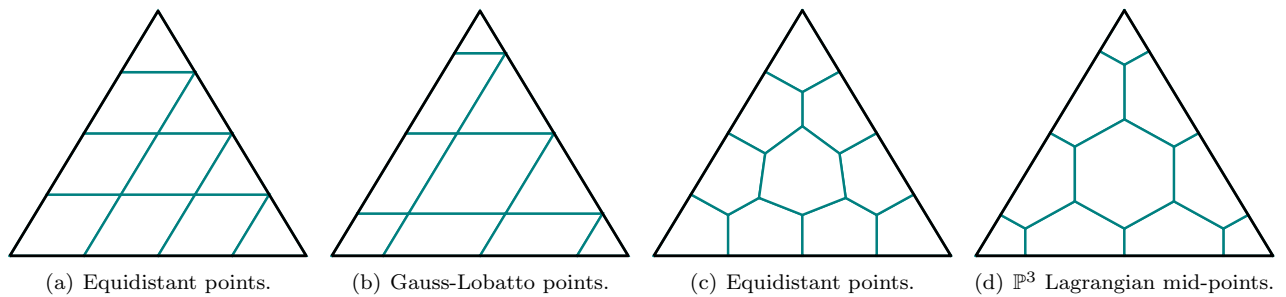


FIG. 11. Examples of structured and Voronoi-type subdivisions for a triangular cell and a  $\mathbb{P}^3$  DG scheme.

In Figures 11(a) and 11(b), the simple case of structured subdivision is depicted where the cell is partitioned as a quadrilateral cell would be and then split into two to fit the triangular shape. In this case, one can note that the subdivision is not rotation invariant and that a choice has to be made as one corner subcell is a quadrilateral while the other two are triangles. In this work, we made the choice to start the structured cell subdivision from the wider angle corner, which induces the quadrilateral subcell to stand as this particular

corner. In Figure 11(a), the cell boundary points defining the subdivision are distributed in an uniform manner, while in Figure 11(b) there are defined as Gauss-Lobatto quadrature points. Let us emphasize that in both cases, those subdivisions are very simple to implement and to generalize to any order of accuracy. Furthermore, the subcells normals are nothing but the ones to the original triangular cell, and the subcells are either triangles or quadrilaterals. In Figures 11(c) and 11(d), polygonal subdivisions are displayed. Those Voronoi type subdivisions are widely used in subcell techniques and can be found for instance in [5, 20] and references within. In Figure 11(d), the subdivision is obtained as follows: first, we define a sub-triangulation of the element by joining the  $\mathbb{P}^k$  Lagrangian nodes. Then, the centroid of each sub-triangle and the midpoint of the edges form a set of  $N_k$  polygonal subcells. In Figure 11(c), the previous procedure is modified to yield a more uniform subdivision by setting equidistant cell boundary points.

Now, to make sure that cell subdivision does not have any impact on the numerical solution, we display in Figures 12 and 13 the solutions obtained through subcell finite volume scheme with high-order reconstructed fluxes as numerical fluxes, with the four different cell subdivisions previously introduced. In Figure 13, where submean values versus  $(x+y-1)$  coordinate are displayed, it is clear that the four calculations lead to exactly the same numerical solution.

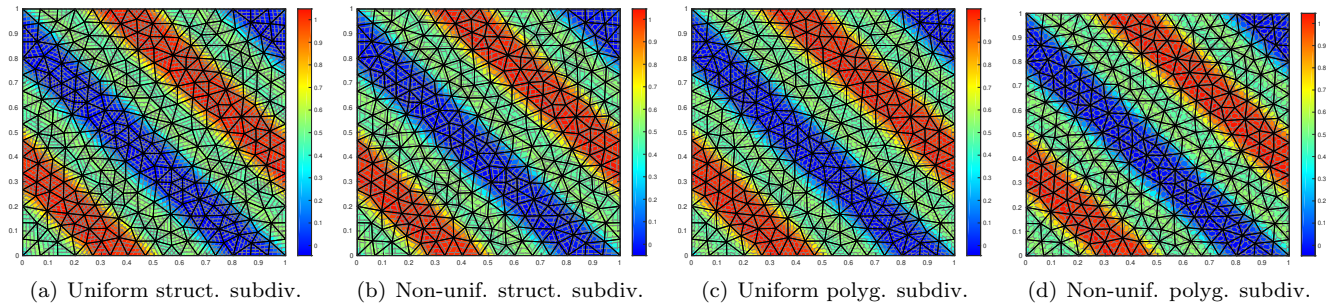


FIG. 12. 4th-order DG solutions for the crenel signal advection on 576 cells after five periods.

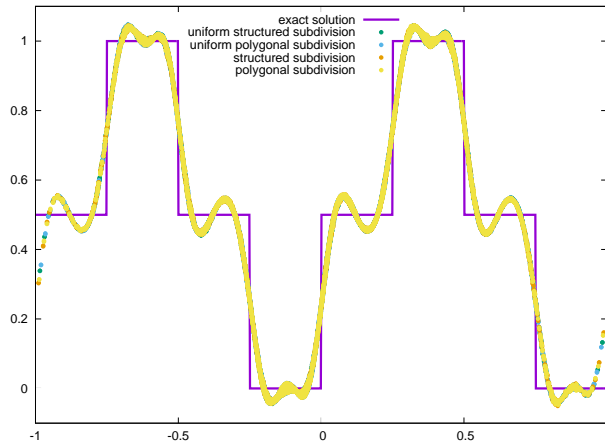


FIG. 13. 4th-order DG solutions for the crenel signal advection on 576 cells using different cell subdivisions: submean values versus  $(x+y-1)$  coordinate.

While the cell subdivision does not affect the uncorrected DG numerical solution, see Figure 13, it has been proved in the 1D case [42] that it does have an impact on the quality of the results when the correction is used, especially in the linear advection case. To assess if any similar phenomenon exists in the 2D unstructured case, we use the same setup as before but this time with the full APLSC-DG method, see Figures 14 and 15. In the light of those results, it appears that uniform subdivisions lead to a way better resolution of the problem under consideration. This huge difference mainly derives from the application of the NAD troubled

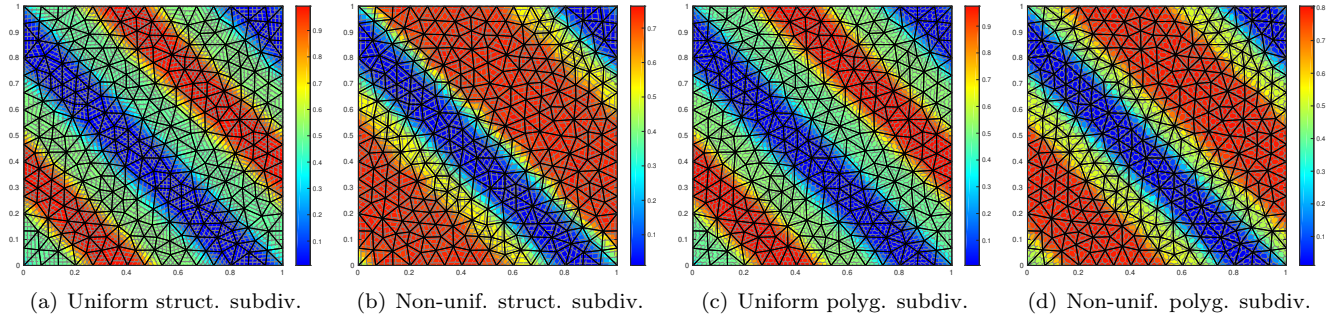


FIG. 14. 4th-order APLSC-DG solutions for crenel advection on 576 cells after five periods

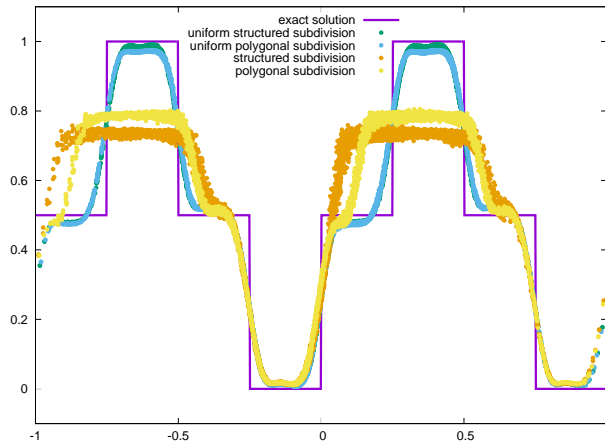


FIG. 15. 4th-order APLSC-DG solutions for the crenel advection case on 576 cells using different cell subdivisions: submean values versus  $(x + y - 1)$  coordinate.

subcell detector, as if we make just use of the PAD criterion to only enforce the global maximum principle, the difference in the results is a lot more slight. In the linear advection case, NAD criterion based on a discrete maximum principle is highly dependent of the subcell aspect ratio. However, considering other type of problem, the difference in the quality of the results is a lot more slight.

**4.1.2. Solid body rotation.** We make use of the classical test case taken from [32]. Let us then consider (4.4) with a divergence-free velocity field corresponding to a rigid rotation, defined by  $\mathbf{A}(\mathbf{x}) = (\frac{1}{2} - y, x - \frac{1}{2})^t$ . We apply this solid body rotation to the initial data displayed in Figure 16(a), which includes both a plotted disk, a cone and a smooth hump. In Figure 16, the submean values obtained through the 6th-order APLSC-DG scheme on a 576 cells coarse grid are displayed. The uniform structured cell subdivision has been used. One can see on Figure 16(b) how the corrected DG scheme produces a very accurate solution, even using a quite coarse mesh, while still ensuring a global maximum principle as well as a mainly non-oscillatory behavior. In Figure 17, cross-sections of the solution along lines  $y = 0.25$  and  $y = 0.75$  have been plotted. Those results further demonstrates the very high capability of the correction procedure presented.

Now, similarly to the linear advection equation, we want to understand how the choice in the cell subdivision impacts the quality of the results now considering a solid body rotation problem. To do so, we compare the solutions obtained by means of the 4th-order APLSC-DG scheme and the four different subdivisions introduced before, see Figure 11, after five full rotations. In the light of Figures 18 and 19, we note how more uniform cell subdivision lead to less diffused solutions. However, it is worth mentioning that the difference in the result is not as critical as it was in the linear advection case. Furthermore, even if the uniform structured subdivision is not rotation invariant, it led to comparable results to the uniform Voronoi-type subdivision, while being a lot more simpler to implement and to generalize to any order of accuracy.

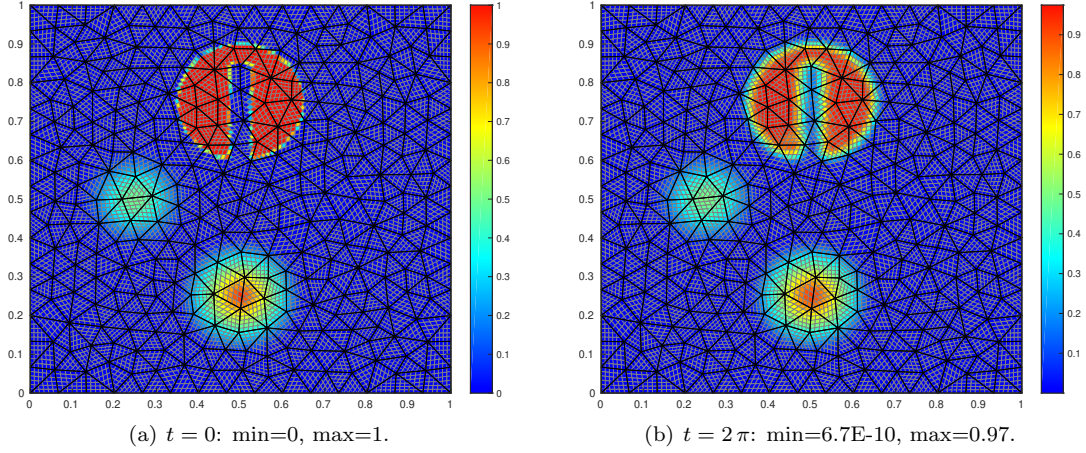


FIG. 16. 6th-order APLSC-DG solution for the rigid rotation case on 576 cells.

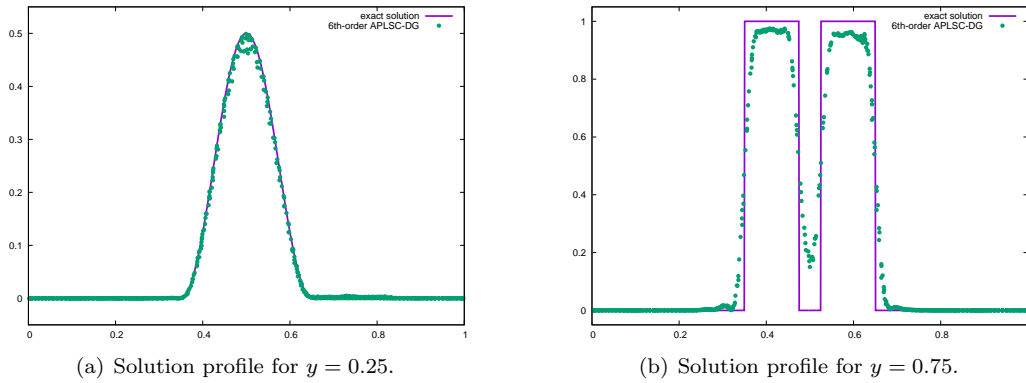


FIG. 17. 6th-order APLSC-DG solution for rigid rotation on 576 cells after one full rotation: solution profiles.

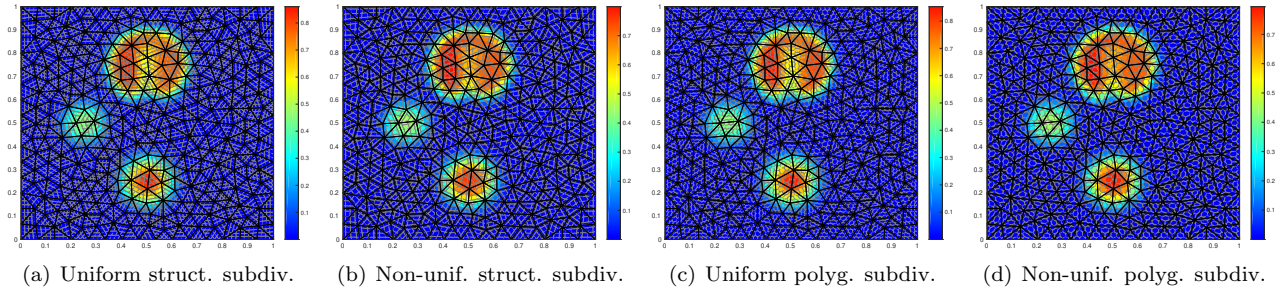


FIG. 18. 4th-order APLSC-DG solutions for rigid rotation on 576 cells after five full rotations.

**4.2. Non-linear case.** Let us now assess the performance and accuracy of our APLSC-DG technique in the 2D non-linear case. Both cases of scalar conservation laws as well as system of conservation laws will be addressed. Similarly to the subsection devoted to the linear case, let us defined set  $\mathcal{N}(S_m^c)$  when considering the NAD criterion on subcell  $S_m^c$ . Unlike the linear case, we make use here of a subcell-wise DMP, meaning  $\mathcal{N}(S_m^c)$  will be constituted by subcell  $S_m^c$ , as well as all its face and node neighboring subcells  $S_q^v$ , either they belong to the same cell or not. By introducing  $\mathcal{P}_m^c$  the set of vertices of subcell  $S_m^c$  as well as  $\mathcal{N}_p$  the set of subcells that share  $\mathbf{x}_p$  as a vertex, this definition can be rewritten as  $\mathcal{N}(S_m^c) = \bigcup_{\mathbf{x}_p \in \mathcal{P}_m^c} \mathcal{N}_p$ . This particular

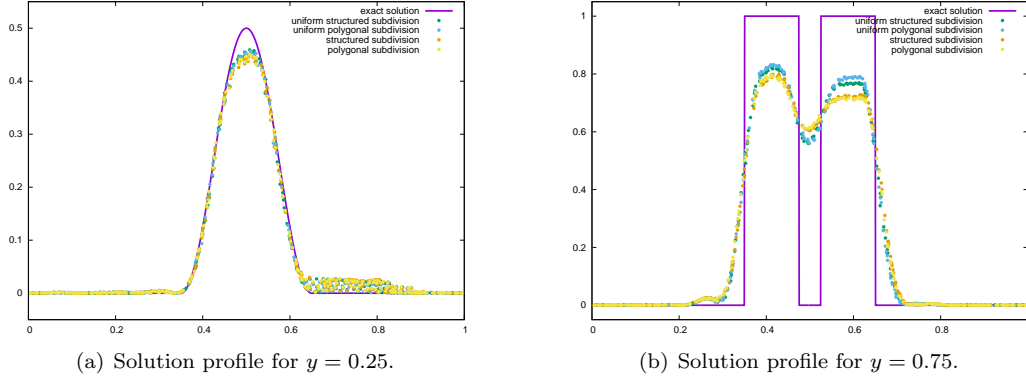


FIG. 19. 4th-order APLSC-DG solutions for rigid rotation on 576 cells after five full rotations: solution profiles.

set is depicted in Figure 7, for both the simple structured subdivision as well as the polygonal Voronoi-type one. In this figures, the subcell  $S_m^c$  under consideration would be colored red, while the subcells constituting  $\mathcal{N}(S_m^c)$  would be colored green. Let us emphasize that subcell  $S_m^c$  is also part of  $\mathcal{N}(S_m^c)$ .

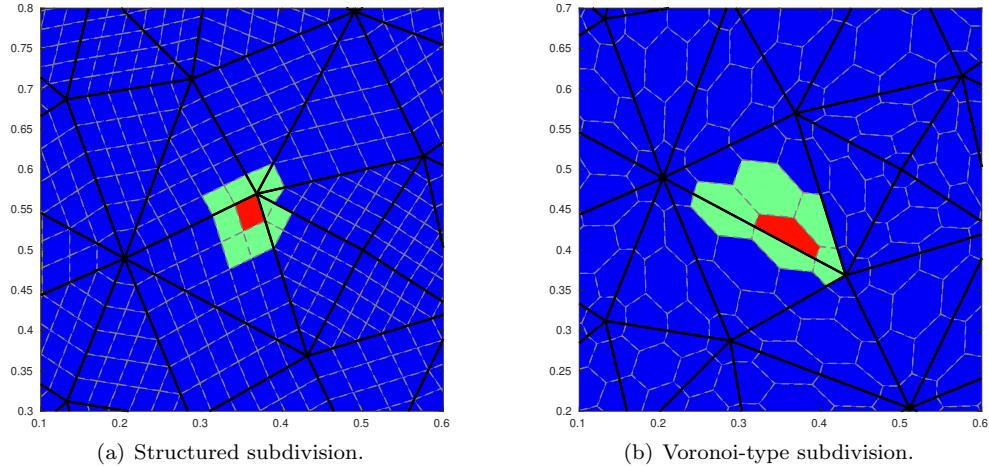


FIG. 20. Neighboring subcells set  $\mathcal{N}(S_m^c)$  for the NAD criterion in the non-linear case: subcell  $S_m^c$  is colored red while the subcells in  $\mathcal{N}(S_m^c)$  are colored green.

**4.2.1. Burgers equation with a smooth initial solution.** To highlight the efficiency of the developed APLSC-DG scheme in the non-linear case, let us first consider the 2D Burgers equation, (2.1), where the convex flux function writes  $\mathbf{F}(u) = \frac{1}{2} (u^2, u^2)^t$ . As seen previously, starting from the smooth initial condition  $u_0(x) = \sin(2\pi(x, y))$  on  $[0, 1]^2$ , two stationary discontinuities form along the lines  $\{(x, y) \in [0, 1]^2, x+y = \frac{1}{2}\}$  and  $\{(x, y) \in [0, 1]^2, x+y = \frac{3}{2}\}$ . To emphasize how important a limiter or a correction technique is needed in this non-linear context, we first represent the numerical solution obtained by means of the 6th-order uncorrected DG on a very coarse grid made of 242 cells, see Figure 21. One can see how oscillating the numerical solution is. Furthermore, the two shocks are absolutely not well resolved. In Figure 22, the numerical solution obtained with the 6th-order APLSC-DG scheme is illustrated at time  $t = 0.5$ . We can see in this totally anisotropic triangles coarse grid, the corrected scheme quite accurately recovers the two straight line shocks, while ensuring a robust low oscillatory behavior. Now, to investigate once more if the cell subdivision has any influence on the quality of the results, we simulate this two-shocks Burgers test case with an even coarser grid with the four subdivisions depicted in Figures 11.

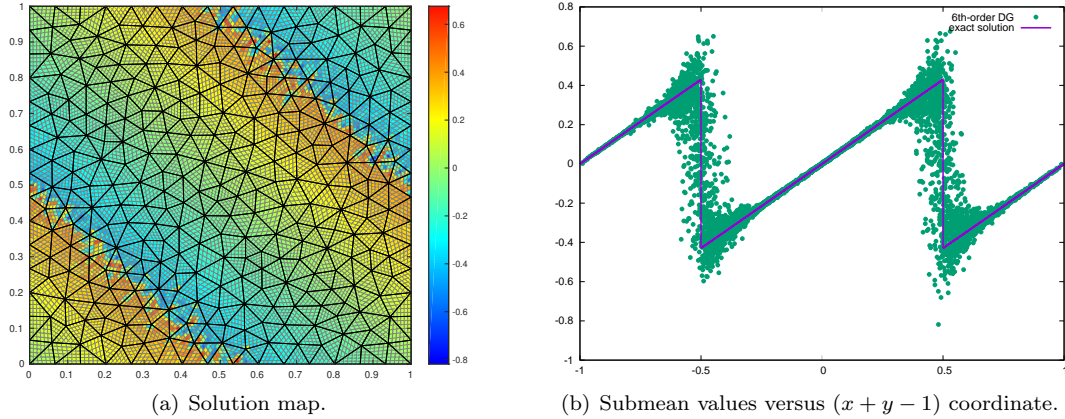


FIG. 21. 6th-order uncorrected DG solution for 2D Burgers equation on a 576 cells mesh at  $t = 0.5$ .

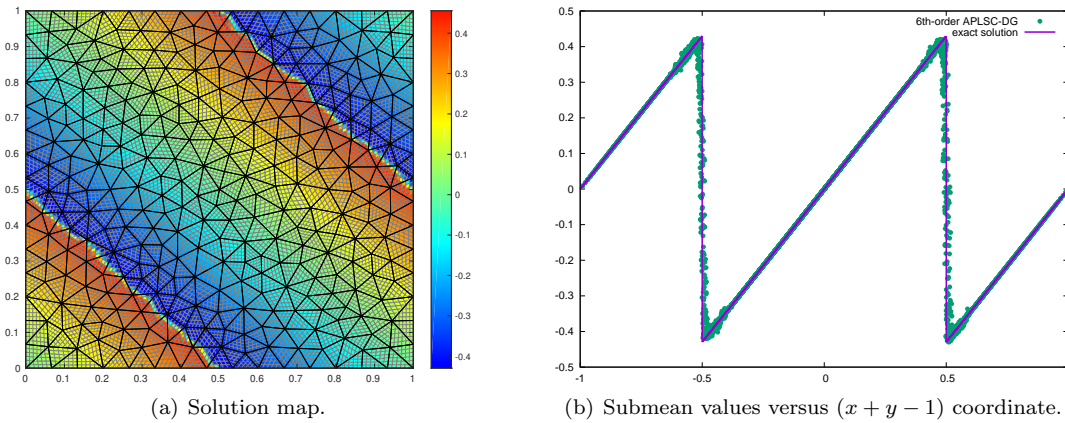


FIG. 22. 6th-order APLSC-DG solution for 2D Burgers equation on a 576 cells mesh at  $t = 0.5$ .

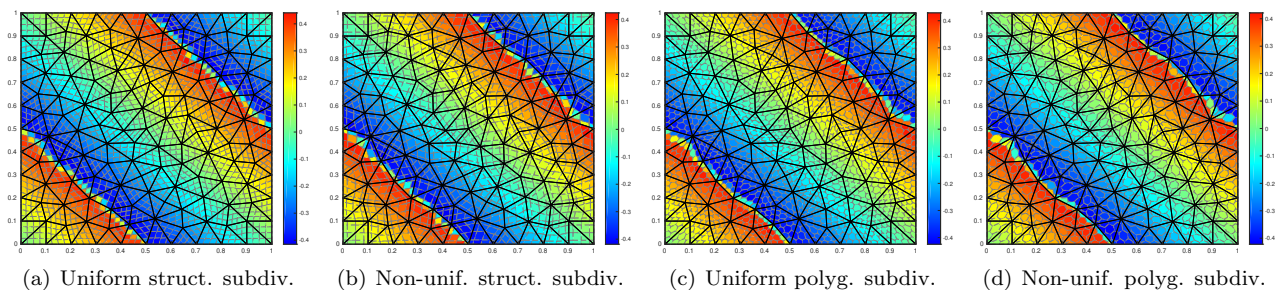


FIG. 23. 4th-order APLSC-DG solutions for 2D Burgers equation on 242 cells at  $t = 0.5$

In the light of Figure 24, the four different subdivisions seem to produce equivalent results, which are further quite satisfactory considering the extremely coarse grid used. However, the use of the uniform structured cell partition, Figure 23(a), appears to capture in a sharper fashion the two straight line shocks. As we have shown that uniform subdivision, structured or Voronoi-type, lead to better results when APLSC-DG scheme is used, only those two subdivisions will be used in the remainder of the article.

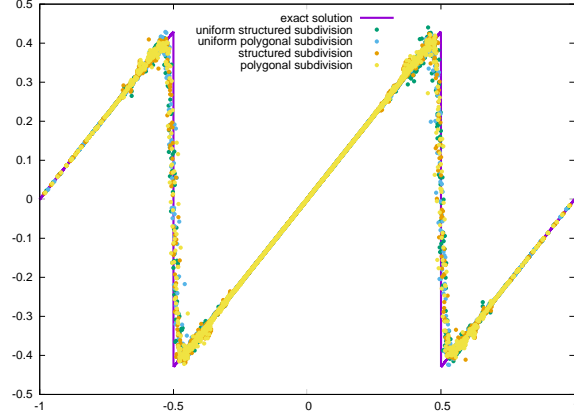


FIG. 24. 4th-order APLSC-DG solutions for 2D Burgers equation on 242 cells: submean values versus  $(x + y - 1)$  coordinate.

**4.2.2. KPP problem.** Before investigating the non-linear system case, we now turn our attention to non-linear conservation laws with non-convex fluxes. To this end, we consider the KPP problem proposed by Kurganov, Petrova, Popov (KPP) in [27] to test the convergence properties of some WENO schemes in the context of non-convex fluxes. For this particular problem, we study the non-linear problem (2.1) where the flux function is given by  $F(u) = (\sin(u), \cos(u))^t$ . Considering the computational domain  $[-2, 2] \times [-2.5, 1.5]$ , the initial condition reads as follows

$$u_0(x) = \begin{cases} 7\pi/2 & \text{if } x < \frac{1}{2}, \\ \pi/4 & \text{if } x > \frac{1}{2}. \end{cases}$$

This test is very challenging to many high-order schemes as the solution has a two-dimensional composite wave structure, and as generally numerical methods fail to converge to the unique entropic exact solution. In most cases, to be able to capture such rotation composite structure, very fine grids must be used. Here, by means of 6th-order uncorrected DG and then APLSC-DG scheme, we make use of an unstructured mesh made of 1054 triangular cells, which is very coarse in this quite complex situation. Results are displayed in Figure 25. Let us emphasize that DG scheme, without any additional correction or treatment, would produce

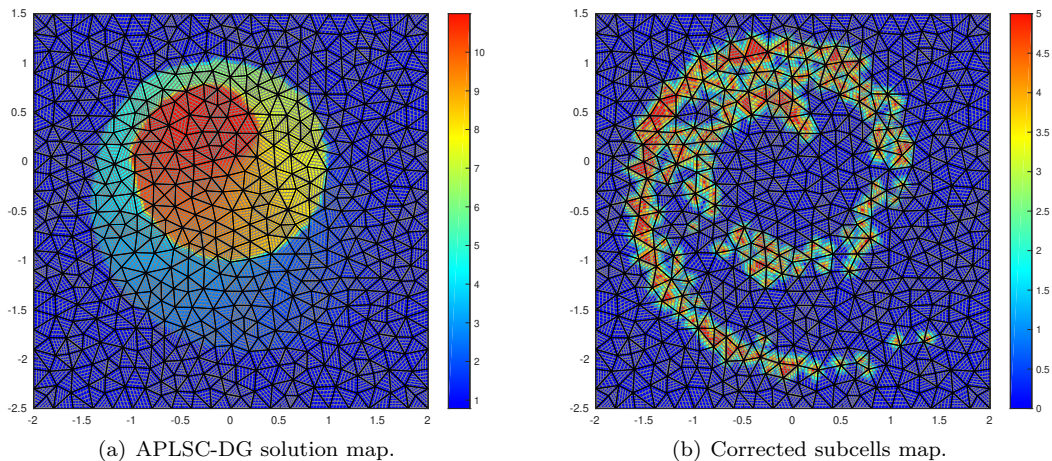


FIG. 25. 6th-order APLSC-DG solution for the KPP problem on a 1054 cells mesh at time  $t = 1$ .

a non-entropic solution which will further be extremely oscillatory. The application of our APLSC technique permits to capture to correct entropic solution, while avoiding the apparition of spurious oscillations, as

displayed in Figure 25(a). Furthermore, we can observe that although the coarseness of the grid used, the APLSC-DG scheme allowed to recover the two-dimensional vortex-like wave structure of the solution. In Figure 25(b), subcells corrected during the different Runge-Kutta stages of the last time iteration are displayed, with different colors accordingly the amount of first-order correction applied. One can see how the NAD and PAD troubled subcell detection criteria accurately track the spiral discontinuity of the entropic exact solution. Here, a different colormap compared to Figure 4 has been used for a better readability of the results. In Figure 26, we once more assess the impact of the subdivision on the quality of the results obtained through the APLSC-DG method. As the uniform ones have proved to yield better results, we only compare here the uniform structured subdivision with the uniform Voronoi-type one. Anew, the 4th-order scheme is used here. As one can see in Figure 26, comparable results have been obtained regardless the type of cell subdivision. This is the reason why we choose to utilize the simple uniform structure subdivision in the remainder for sake of simplicity and computational cost.

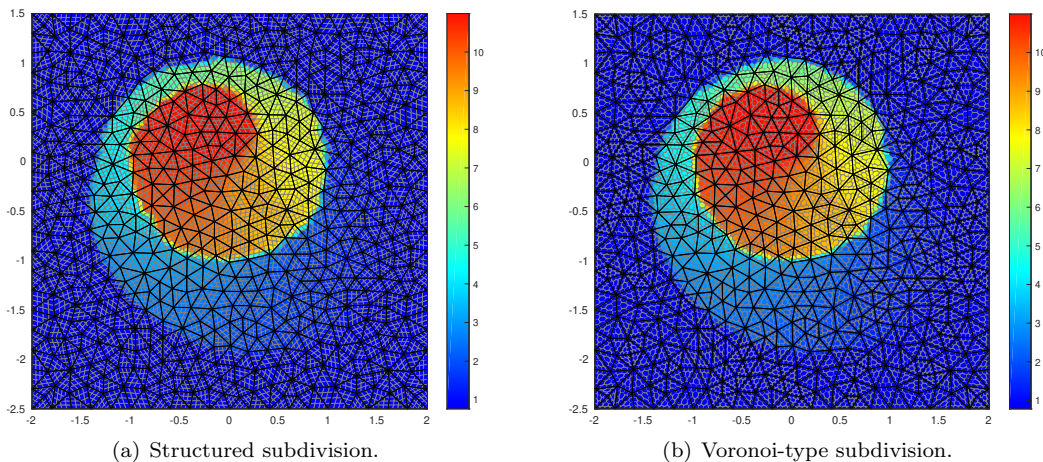


FIG. 26. 4th-order APLSC-DG solution for the KPP problem on a 1054 cells mesh at time  $t = 1$ : subdivisions comparison.

**4.3. 2D Euler system.** To close this numerical application section, and assess once again the high capability of the *a posteriori* local subcell correction technique presented here, the non-linear system case will be now addressed. To this end, let us consider the 2D Euler compressible gas dynamics system

$$\begin{cases} \partial_t \rho + \nabla_x \cdot \mathbf{q} = 0, \\ \partial_t \mathbf{q} + \nabla_x \cdot (\rho \mathbf{v} \otimes \mathbf{v} + p I_d) = 0, \\ \partial_t E + \nabla_x \cdot ((E + p) \mathbf{v}) = 0, \end{cases}$$

where the conserved variables  $\rho$ ,  $\mathbf{q} = \rho \mathbf{v}$  and  $E$  respectively stand for the density, momentum and total energy, while  $\mathbf{v}$  characterizes the fluid velocity. The thermodynamic closure is given by the equation of state  $p = p(\rho, \varepsilon)$  where  $\varepsilon = E - \frac{1}{2} \rho \|\mathbf{v}\|^2$  denotes the internal energy. In this paper, we make use of a gamma gas law, *i.e.*  $p = (\gamma - 1) \varepsilon$ , where  $\gamma$  is the polytropic index of the gas.

Although the whole theory presented here has been introduced in the simple case of scalar conservation laws, the extension to the system case is perfectly straightforward. The only part which may differ is the troubled detection part. For the PAD, we consider that a solution is physically admissible if the density and the internal energy are strictly positive. The use of other equations of state may lead to a different convex set of admissibility, see [44] for instance. For the NAD, the natural system counterpart would be to apply the previously introduced detection criteria to the Riemann invariants. However, in the non-linear system case, those quantities are not easy to get nor to manipulate. We could have use a linearized version of the Riemann invariants, as in [43] for instance, but for sake of simplicity we naively apply the NAD criterion to one of the conserved variable. Here, we choose to work with the energy, as this physical quantity would be sensitive to any type of wave. Once again, the simple global Lax-Friedrichs numerical flux will be used in the remainder.



**4.3.1. Sod shock tube problem.** We consider the extension of the classical Sod shock tube [39] to the case of the cylindrical geometry. This problem consists of a cylindrical shock tube of unity radius. The interface is located at  $r = 0.5$ . At the initial time, the states on the left and on the right sides of the interface are constant. The left state is a high pressure fluid characterized by  $(\rho_0^L, p_0^L, \mathbf{v}_0^L) = (1, 1, \mathbf{0})$ , the right state is a low pressure fluid defined by  $(\rho_0^R, p_0^R, \mathbf{v}_0^R) = (0.125, 0.1, \mathbf{0})$ . The gamma gas law is defined by  $\gamma = \frac{7}{5}$ . The computational domain is defined in polar coordinates by  $(r, \theta) \in [0, 1] \times [0, \frac{\pi}{4}]$ . We prescribe symmetry boundary conditions at the boundaries  $\theta = 0$  and  $\theta = \frac{\pi}{4}$ , and an outflow condition at  $r = 1$ . The exact solution consists of three circular waves, a shock followed by a contact discontinuity and rarefaction wave. The aim of this test case is then to assess the APLSC-DG scheme accuracy while ensuring a non-oscillatory behavior, and its ability to preserve the radial symmetry.

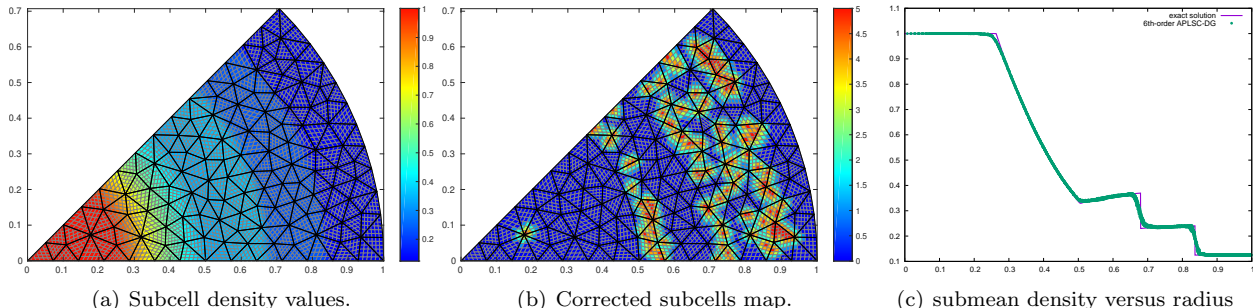


FIG. 27. 6th-order APLSC-DG solution for the polar Sod shock tube problem on 230 cells.

In Figure 27, the 6th-order APLSC-DG scheme has been used on a very coarse anisotropic mesh made of 230 triangular cells. In the light of Figure 27(a), one can see how the radial wave structure has been accurately captured, even in this coarse mesh case. Figure 27(c), where the subcell mean values versus the subcell centroid radii  $\sqrt{x^2 + y^2}$  are displayed, confirms this statement as the different points for a given radius do coincide. In Figure 27(b), subcells corrected during the different Runge-Kutta stages of the last time iteration are colored accordingly the amount of first-order FV correction used. It illustrates how this *a posteriori* correction procedure has been activated on zones corresponding to the solution loss of smoothness, meaning the left and right ends of the expansion fan, the contact discontinuity and the shock. One can also observe how the correction does operate locally inside the cell at a subcell scale, allowing the preservation of DG subcell high accurate resolution. Let us emphasize that this *a posteriori* correction procedure is not limited to the case of very high-order of accuracy on coarse grids. It also performs very well at second or third order.

**4.3.2. Sedov point blast problem.** We consider the Sedov problem for a point-blast in a uniform medium. An exact solution based on self-similarity arguments is available, see for instance [24]. The initial conditions are characterized by  $(\rho_0, p_0, \mathbf{v}_0) = (1, 10^{-14}, \mathbf{0})$ , and the polytropic index is equal to  $\frac{7}{5}$ . We set an initial delta-function energy source at the origin prescribing the pressure in a control volume, yet to be defined, containing the origin as follows,  $p_{or} = (\gamma - 1) \frac{\varepsilon_0}{v_{or}}$ , where  $v_{or}$  denotes the measure of the chosen control volume and  $\varepsilon_0$  the total amount of release energy. By choosing  $\varepsilon_0 = 0.244816$ , as suggested in [24], the solution consists of a diverging infinite strength shock wave whose front is located at radius  $r = 1$  at  $t = 1$ , with a peak density reaching 6. The computational domain is defined in polar coordinates by  $(r, \theta) \in [0, 1.2] \times [0, \frac{\pi}{4}]$ . Similarly to the polar Sod shock tube problem, we prescribe symmetry boundary conditions at the boundaries  $\theta = 0$  and  $\theta = \frac{\pi}{4}$ , and an outflow condition at  $r = 1.2$ .

Regarding the control volume in which the delta-function energy will be dropped off, generally the cell containing the origin is considered. Here, to make this test case even more challenging, we choose to restrict the energy source only to the one subcell containing the origin. This means that initially, in one grid element the pressure in one subcell will be set to  $p_{or}$ , while in the remainder of the cell the pressure will be  $10^{-14}$ . Let us further emphasize that generally in this test case, because one cannot simulate vacuum, the initial pressure is set to  $10^{-6}$  over the domain, except at the origin. Here, to make it once again more challenging,

we set the initial pressure to  $10^{-14}$ .

We run this modified Sedov point blast problem with the sixth-order APLSC-DG scheme on a very coarse grid made of 271 triangular cells. In this particular case, the amount of total energy contained in the subcell located at the origin reaches 1947.5, while in the rest of the cell as well as in the remainder of the domain the total energy is set to  $2.5E-14$ . Any scheme lacking positivity-preserving property or a rigorous stabilization technique would fail solving this test problem.

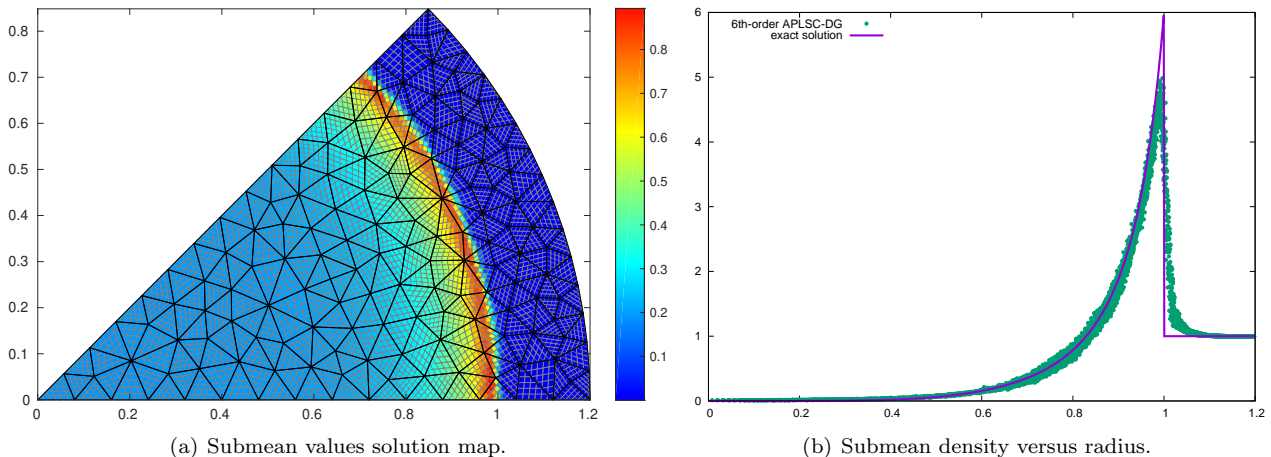


FIG. 28. 6th-order APLSC-DG solution for the Sedov problem on 271 cells: subcell mean total energy values.

In Figure 28(a), one can see how the circular aspect of the shock has been accurately captured by the scheme, and the shock wave front is correctly located. This latter further goes through and inside different cells, enlightening the very precise subcell resolution of the APLSC-DG method. The numerical solution produced remains quite close to the one-dimensional self-similar exact solution, see Figure 28(b).

**4.3.3. The forward-facing step problem.** We now consider the forward facing step problem, which has been initially introduced by A. Emery in [12], and further studied by P. Woodward and P. Colella in [48]. This challenging test case consists in a Mach 3 flow in a 3 units in length and 1 unit in width wind tunnel. Initially, the tunnel is filled with a gamma gas law with  $\gamma = \frac{7}{5}$ , which everywhere has density  $\rho_0 = 1.4$ , pressure  $p_0 = 1$  and velocity  $\mathbf{v}_0 = (3, 0)^t$ . The 0.2 high step being located at  $x = 0.6$ , the computational domain is then  $([0, 3] \times [0, 1]) \setminus ([0.6, 3] \times [0.2, 1])$ . Gas with this density, pressure and velocity is continually fed in from the left-hand boundary. Let us emphasize that unlike as it is generally done, we did not refine the mesh near the corner, see Figure 29, nor modify in any way our APLSC-DG scheme.

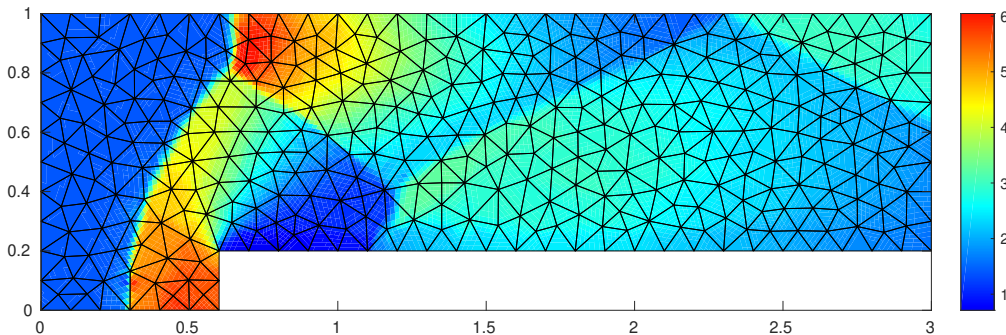


FIG. 29. 6th-order APLSC-DG solution for the facing step problem on 680 cells at  $t = 4$ : submean density map.

In Figure 29, the numerical solution obtained by means of 6th-order APLSC-DG scheme on an unstructured grid made of 680 cells is displayed. Let us note that despite the coarseness of the mesh used, the shocks

and the rarefaction fan created around the corner are quite well resolved, while ensuring a low oscillatory behavior. This result demonstrates once again the high capability of the presented *a posteriori* local subcell correction combined with high-order discontinuous Galerkin methods.

**5. Conclusion.** The paper aims at presenting the two-dimensional unstructured extension of the new correction technique of DG schemes introduced in [42]. This *a posteriori* procedure relies on the expression of DG methods as a FV-like scheme on a subgrid. By means of this theoretical part, we modify at the subcell level the so-called reconstructed fluxes only where the uncorrected DG scheme has failed. Consequently, only very few subcells require this particular treatment. In this paper, a new version of the correction procedure is also introduced, where a convex blending of high-order DG reconstructed fluxes and first-order FV fluxes is applied in the vicinity of troubled zones. For the remaining subcells, the submean values obtained through the uncorrected DG method are kept, as they have been detected as admissible by troubled zone criteria. This correction procedure allows us to retain the very precise subcell resolution of DG schemes, along with addressing the issues of spurious oscillations or non-entropic behavior. A wide number of test cases on different problems have been used to depict the good performance and robustness of the presented correction technique. Different types of cell subdivision have also been compared.

In the future, we intend to extend this *a posteriori* correction technique to moving grid configurations, with both ALE and Lagrangian formalisms, as well as the case of curvilinear meshes. We also plan to adapt this local subcell reconstructed flux correction framework to the *a priori* paradigm, by means of the FCT methodology, in order to obtain an automatic very high-order and property preserving scheme.

**Acknowledgment.** R.A. was partially funded by SNF project "Structure preserving and fast methods for hyperbolic systems of conservation laws" number 200020\_204917.

### Appendix A. Graph Laplacian existence and uniqueness condition.

This appendix aims at giving further details on the existence and uniqueness condition of the solution in the graph Laplacian technique.

**A.1. Through residuals.** Let us check that  $(D_c P_c M_c^{-1} \Phi_c + B_c) \cdot \mathbf{1} = 0$ . First, from  $B_c$  definition it follows that

$$B_c \cdot \mathbf{1} = \sum_{m=1}^{N_k} (B_c)_m = \sum_{m=1}^{N_k} \int_{\partial S_m^c \cap \partial \omega_c} \mathcal{F}_n \, dS = \int_{\partial \omega_c} \mathcal{F}_n \, dS.$$

Now, evaluating the first term, and by means of (2.7), one gets

$$(D_c P_c M_c^{-1} \Phi_c) \cdot \mathbf{1} = M_c^{-1} \Phi_c \cdot P_c^t D_c \mathbf{1} = \frac{d}{dt} \begin{pmatrix} u_1^c \\ \vdots \\ u_{N_k}^c \end{pmatrix} \cdot \begin{pmatrix} \int_{\omega_c} \sigma_1^c \, dV \\ \vdots \\ \int_{\omega_c} \sigma_{N_k}^c \, dV \end{pmatrix} = \frac{d}{dt} \int_{\omega_c} u_h^c \, dV.$$

Finally, making use of (2.3) with  $\psi = 1$ , the previous relation reduces to  $(D_c P_c M_c^{-1} \Phi_c) \cdot \mathbf{1} = - \int_{\partial \omega_c} \mathcal{F}_n \, dS$ .

**A.2. Through fluxes.** Let us check that  $G_c \cdot \mathbf{1} = 0$ . Let us emphasize that due to (2.28) the sub-resolution  $\{\phi_m\}_m$  sum to one. Furthermore, for any  $\mathbf{x} \in \partial \omega_c$ , we have that  $\sum_m \mathbb{1}_{\partial S_m^c} = 1$ . Consequently, it directly follows that

$$G_c \cdot \mathbf{1} = \int_{\omega_c} \left( \underbrace{\sum_{m=1}^{N_k} \phi_m^c}_1 - \underbrace{\sum_{m=1}^{N_k} \mathbb{1}_{\partial S_m^c}}_1 \right) (\mathbf{F}_h^c \cdot \mathbf{n} - \mathcal{F}_n) \, dS = 0.$$

## REFERENCES

- [1] D. Balsara, C. Altmann, C.D. Munz, and M. Dumbser. A sub-cell based indicator for troubled zones in RKDG schemes and a novel class of hybrid RKDG+HWENO schemes. *J. Comp. Phys.*, 226:586–620, 2007.
- [2] R. Biswas, K. Devine, and J.E. Flaherty. Parallel adaptive finite element methods for conservation laws. *Applied Numerical Mathematics*, 14:255–284, 1994.
- [3] J.-P. Boris and D.-L. Book. Flux-corrected transport. I. SHASTA, a fluid transport algorithm that works. *J. Comp. Phys.*, 11(1):38–69, 1973.
- [4] A. Burbeau, P. Sagaut, and C.-H. Bruneau. A problem-independent limiter for high-order Runge Kutta discontinuous Galerkin methods. *J. Comp. Phys.*, 169:111–150, 2001.
- [5] Qian-Yong Chen. Partitions of a simplex leading to accurate spectral (finite) volume reconstruction. *SIAM J. Sci. Comput.*, 27:1458–1470, 2006.
- [6] S. Clain, S. Diot, and R. Loubère. A high-order finite volume method for hyperbolic systems: Multi-dimensional Optimal Order Detection (MOOD). *J. Comp. Phys.*, 230:4028–4050, 2011.
- [7] B. Cockburn, S. Hou, and C.-W. Shu. The Runge-Kutta Discontinuous Galerkin Method for Conservation Laws V: Multidimensional Systems. *J. Comp. Phys.*, 141:199–224, 1998.
- [8] S. Diot, S. Clain, and R. Loubère. Improved detection criteria for the multi-dimensional optimal order detection (MOOD) on unstructured meshes with very high-order polynomials. *Computers and Fluids*, 64:43–63, 2012.
- [9] S. Diot, R. Loubère, and S. Clain. The MOOD method in the three-dimensional case: very-high-order finite volume method for hyperbolic systems. *Int. J. Numer. Meth. Fluids*, 73:362–392, 2013.
- [10] M. Dumbser and R. Loubère. A simple robust and accurate a posteriori sub-cell finite volume limiter for the discontinuous Galerkin method on unstructured meshes. *J. Comp. Phys.*, 319:163–199, 2016.
- [11] M. Dumbser, O. Zanotti, R. Loubère, and S. Diot. A posteriori subcell limiting of the discontinuous Galerkin finite element method for hyperbolic conservation laws. *J. Comp. Phys.*, 278:47–75, 2014.
- [12] A. Emery. An evaluation of several differencing methods for inviscid flow problems. *J. Comp. Phys.*, 2(3):306–331, 1968.
- [13] T. C. Fisher and M. H. Carpenter. High-order entropy stable finite difference schemes for nonlinear conservation laws: Finite domains. *J. Comp. Phys.*, 252:518–557, 2013.
- [14] G. Gassner. A Skew-Symmetric Discontinuous Galerkin Spectral Element Discretization and Its Relation to SBP-SAT Finite Difference Methods. *SIAM J. Sci. Comput.*, 35:1233–1253, 2013.
- [15] J. L. Guermond and R. Pasquetti. Entropy viscosity method for high-order approximations of conservation laws. *Lecture Notes in computational Science and Engineering*, 2009.
- [16] R. Harris and Z. J. Wang. Partition design and optimization for high-order spectral volume schemes. In *47th AIAA Aerospace Sciences Meeting, Orlando*, 2009.
- [17] A. Harten. ENO schemes with subcell resolution. *J. Comp. Phys.*, 83:148–184, 1989.
- [18] A. Harten, B. Engquist, S. Osher, and S. Chakravarthy. Uniformly high order essentially non-oscillatory schemes, III. *J. Comp. Phys.*, 71:231–303, 1987.
- [19] S. Hennemann, A.M. Rueda-Ramirez, F.J. Hindenlang, and G.J. Gassner. A provably entropy stable subcell shock capturing approach for high order split form dg for the compressible euler equations. *J. Comp. Phys.*, 426:109935, 2021.
- [20] A. Huerta, E. Casoni, and J. Peraire. A simple shock-capturing technique for high-order discontinuous Galerkin methods. *Int. J. Numer. Meth. Fluids*, 69:1614–1632, 2012.
- [21] H. T. Huynh. A Flux Reconstruction Approach to High-Order Schemes Including Discontinuous Galerkin Methods. In *18th AIAA Computational Fluid Dynamics Conference, Miami*, 2007.
- [22] J. S. Park and S.-H. Yoon and C. Kim. Multi-dimensional limiting process for hyperbolic conservation laws on unstructured grids. *J. Comp. Phys.*, 229:788–812, 2010.
- [23] G.-S. Jiang and C.-W. Shu. Efficient implementation of weighted eno schemes. *J. Comp. Phys.*, 126:202–228, 1996.
- [24] J.R. Kamm and F.X. Timmes. On efficient generation of numerically robust Sedov solutions. Technical Report LA-UR-07-2849, Los Alamos National Laboratory, 2007.
- [25] L. Krivodonova. Limiters for high-order discontinuous Galerkin methods. *J. Comp. Phys.*, 226:879–896, 2007.
- [26] L. Krivodonova, J. Xin, J.-F. Remacle, N. Chevaugeon, and J.E. Flaherty. Shock detection and limiting with discontinuous Galerkin methods for hyperbolic conservation laws. *Appl. Numer. Math.*, 48:323–338, 2004.
- [27] A. Kurganov, G. Petrova, and B. Popov. Adaptive semi-discrete central-upwind schemes for non convex hyperbolic conservation laws. *SIAM J. Sci. Comput.*, 29:2381–2401, 2007.
- [28] D. Kuzmin. A vertex-based hierarchical slope limiter for p-adaptive discontinuous Galerkin methods. *J. Comp. Appl. Math.*, 233:3077–3085, 2009.
- [29] D. Kuzmin and M.Q. de Luna. Subcell flux limiting for high-order bernstein finite element discretizations of scalar hyperbolic conservation laws. *J. Comp. Phys.*, 411:109411, 2020.
- [30] D. Kuzmin, M.Q. de Luna, D.I. Ketcheson, and J. Gröll. Bound-preserving flux limiting for high-order explicit Runge Kutta time discretizations of hyperbolic conservation laws. *J. Sci. Comput.*, 91, 2022.
- [31] D. Kuzmin, R. Löhner, and S. Turek. *Flux-Corrected Transport: Principles, Algorithms and Applications*. Scientific Computation. Springer, 2012.
- [32] R. J. LeVeque. High-resolution conservative algorithms for advection in compressible flow. *SIAM J. Numer. Anal.*, 33:627–665, 1996.
- [33] L. Li and Q. Zhang. A new vertex-based limiting approach for nodal discontinuous Galerkin methods on arbitrary unstructured meshes. *Computers and Fluids*, 159:316–326, 2017.
- [34] Will Pazner. Sparse invariant domain preserving discontinuous galerkin methods with subcell convex limiting. *Computer Methods in Applied Mechanics and Engineering*, 382:113876, 2021.

- [35] P.-O. Persson and J. Peraire. Sub-cell shock capturing for discontinuous Galerkin methods. *AIAA paper*, 2006.
- [36] R. Abgrall. Some Remarks About Conservation for Residual Distribution Schemes. *Comput. Methods Appl. Math.*, 18:327–351, 2018.
- [37] C.-W. Shu. TVB uniformly high-order schemes for conservation laws. *Math. Comp.*, 49:105–121, 1987.
- [38] C.-W. Shu and S. Osher. Efficient implementation of essentially non-oscillatory shock-capturing schemes. *J. Comp. Phys.*, 77:439–471, 1988.
- [39] G. A. Sod. A survey of several finite difference methods for systems of non-linear hyperbolic conservation laws. *J. Comp. Phys.*, 27:1–31, 1978.
- [40] M. Sonntag and C. D. Munz. Shock capturing for discontinuous Galerkin methods using finite volume subcells. In *Finite Volumes for Complex Applications VII*, pages 945–953. Springer, 2014.
- [41] B. van Leer. Towards the ultimate conservative difference scheme. V-A second-order sequel to Godunov’s method. *J. Comput. Phys.*, 32:101–136, 1979.
- [42] F. Vilar. A posteriori correction of high-order discontinuous Galerkin scheme through subcell finite volume formulation and flux reconstruction. *J. Comp. Phys.*, 387:245–279, 2018.
- [43] F. Vilar, P.-H. Maire, and R. Abgrall. A discontinuous Galerkin discretization for solving the two-dimensional gas dynamics equations written under total Lagrangian formulation on general unstructured grids. *J. Comp. Phys.*, 276:188–234, 2014.
- [44] F. Vilar, C.-W. Shu, and P.-H. Maire. Positivity-preserving cell-centered Lagrangian schemes for multi-material compressible flows: From first-order to high-orders. Part I: The one-dimensional case. *J. Comp. Phys.*, 312:385–415, 2016.
- [45] Z. J. Wang and H. Gao. A unifying lifting collocation penalty formulation including the discontinuous Galerkin, spectral volume/difference methods for conservation laws on mixed grids. *J. Comp. Phys.*, 228:8161–8186, 2009.
- [46] Z.J. Wang. Spectral (Finite) Volume Method for Conservation Laws on Unstructured Grids: Basic Formulation. *J. Comp. Phys.*, 178:210–251, 2002.
- [47] Z.J. Wang and Y. Liu. Spectral (Finite) Volume Method for Conservation Laws on Unstructured Grids: II. Extension to two-dimensional scalar equation. *J. Comp. Phys.*, 178:210–251, 2002.
- [48] P. Woodward and P. Collela. The numerical-simulation of two-dimensional fluid-flow with strong shocks. *J. Comp. Phys.*, 54(1):115–173, 1984.
- [49] K. Wu and C.-W. Shu. Geometric quasilinearization framework for analysis and design of bound-preserving schemes, 2021.
- [50] M. Yang and Z.J. Wang. A parameter-free generalized moment limiter for high-order methods on unstructured grids. *Adv. Appl. Math. Mech.*, 4:451–480, 2009.
- [51] S.T. Zalesak. Fully multidimensional flux-corrected transport algorithms for fluids. *J. Comp. Phys.*, 31(3):335–362, 1979.
- [52] X. Zhang and C.-W. Shu. Maximum-principle-satisfying and positivity-preserving high-order schemes for conservation laws: survey and new developments. *Proc. R. Soc. Lond., Ser. A, Math. Phys. Eng. Sci.*, 467(2134):2752–2776, 2011.