# New two-dimensional slope limiters for discontinuous Galerkin methods on arbitrary meshes

H. Hoteit[1,2,‡], Ph. Ackerer[1,*,†], R. Mosé[1,§,**], J. Erhel[2,¶] and B. Philippe[2,‖]

[1]*Institut de Mécanique des Fluides, Univ. Louis Pasteur de Strasbourg, CNRS/UMR 7507, 2 rue Boussingault, F-67000 Strasbourg, France*
[2]*IRISA-INRIA, Campus de Beaulieu, 35042 Rennes cedex, France*

## SUMMARY

In this paper, we introduce an extension of Van Leer's slope limiter for two-dimensional discontinuous Galerkin (DG) methods on arbitrary unstructured quadrangular or triangular grids. The aim is to construct a non-oscillatory shock capturing DG method for the approximation of hyperbolic conservative laws without adding excessive numerical dispersion. Unlike some splitting techniques that are limited to linear approximations on rectangular grids, in this work, the solution is approximated by means of piecewise quadratic functions. The main idea of this new reconstructing and limiting technique follows a well-known approach where local maximum principle regions are defined by enforcing some constraints on the reconstruction of the solution. Numerical comparisons with some existing slope limiters on structured as well as on unstructured meshes show a superior accuracy of our proposed slope limiters. Copyright © 2004 John Wiley & Sons, Ltd.

KEY WORDS: hyperbolic conservative laws; discontinuous Galerkin methods; slope limiters; upwind schemes

## 1. INTRODUCTION

The success of discontinuous Galerkin (DG) methods in approximating various physical problems notably hyperbolic systems of conservative laws has attracted the attention to explore the benefits of this approach. One favourable property of DG methods is that they conserve mass at the element level in a finite element framework. Consequently, they inherit the flexibility

of finite elements in handling complicated geometries. Furthermore, the particular approximation space of these methods, where continuity across inter-element boundaries is not enforced, allows a simple treatment of non-homogeneous finite element geometries as well as different degree of approximating polynomials. It is known that when using constant cell approximations the numerical diffusion due to upwinding is big enough to keep the scheme stable. However, by using higher order approximation spaces the scheme produces non-physical oscillations near shocks. In such a case, the use of an appropriate slope limiter is crucial to ensure the stability of the method.

In one dimension, discontinuous finite-elements can be interpreted as a generalization of high order Godunov finite differences [1–4]. Such high resolution schemes are usually stabilized using some form of total variation diminishing (TVD) limiters (see, e.g. References [5, 6]) so that spurious oscillations can be avoided without destroying the high-order accuracy of the schemes. One commonly used technique is the Van Leer's Monotonic Upstream Centred Schemes for Conservative Laws (MUSCL) slope limiter [4]. In the works of Cockburn and Shu [7–9], this slope limiter is extended to the so-called generalized slope limiter where a $(k+1)$st order of accuracy is achieved in smooth regions by using DG method with polynomials of degree $k$ for the spatial discretization and a special $(k+1)$st order explicit Runge–Kutta method for temporal discretization. The generalized slope limiter does not totally suppress oscillations near shocks so that the scheme accuracy is preserved in smooth regions. Thus, the resulting scheme is no more TVD, however it satisfies a total variation bounded (TVB) property.

In multi-dimensional spaces, DG methods are still facing difficulties to attain the same degree of accuracy as in the one-dimensional case, specially on unstructured meshes. The troublesome part is the construction of appropriate multi-dimensional slope limiters that preserve the accuracy of the scheme. Nevertheless, it is proved that any scheme combined with a slope limiting operator that enforces a TVD condition is at most first-order accurate [10]. Consequently, a great deal of effort has been oriented for the construction of genuinely multi-dimensional slope limiters that can eliminate non-physical oscillations without adding excessive numerical viscosity. One simple approach in the case of rectangular grids is to use the DG method with linear polynomials ($P^1$) for the space discretization instead of quadratic ones ($Q^1$) [11]. This approach can be considered as a dimensional splitting technique [5]. In this case, the slope limiting process can be carried out by applying a one-dimensional slope limiter sequentially in the $x$ and $y$ directions.

In our work, we concentrate on a genuinely multi-dimensional slope limiter in the sense that it does not require any operator splitting. This slope limiting operator was introduced by Chavent and Jaffré [12] as a generalization of Van Leer's MUSCL limiter [4]. It can be applied in a geometric manner so that slopes are limited in such a way that each sub-reconstruction lies between the cell averages of its neighbours. In the one-dimensional case, Gowda and Jaffré have analysed this limiter and proved the stability of the DG method with the TVD property [11, 13]. Nevertheless, we have found that the proposed extension of this limiter to the multi-dimensional case does not give satisfactory results. We have detected some cases for both triangular and rectangular discretizations where the limiting operator fails to completely eliminate undershots and overshots. The origin of this drawback is due to the fact that limiting slopes by using the nodal values of the solution does not prohibit non-physical values at the midpoints of the cell edges. As a result, this approach does not satisfy a local maximum principle. This paper proposes a remedy whereby this limiting technique can be improved by giving more weight to the averages of the cell edges.

For triangular elements, piecewise linear approximations are used with degrees of freedom at the grid vertices. Our limiting process intends to reconstruct the solution first at the midpoints of the cell edges by preventing local extrema then at the cell vertices by using the midpoint reconstructions. This requires less restrictive constraints than reconstructing the function at the cell vertices. On the other hand, we have found that by taking the degrees of freedom at the midpoints of the grid edges, the scheme leads to excessive smearing.

For rectangular elements, the solution is approximated by using piecewise quadratic function where the degrees of freedom are indeed at the grid vertices. We use similar techniques for the reconstructions at the midpoints of the edges within each cell. Unfortunately, the information at the midpoints of the edges is not sufficient to give a unique reconstruction at the cell vertices. Consequently, we append supplementary constraints in order to overcome the singularity of the system.

The DG finite element method for scalar, linear conservation laws is reviewed in the next section. In Section 3, we present the slope limiter introduced by Chavent and Jaffré in one- and higher dimensional spaces. We give a simple numerical test where this limiter fails to eliminate all oscillations. Section 4 is devoted to describe our modified slope limiter for unstructured rectangular and triangular grids. In Section 5, we briefly review some existing slope limiters, in particular those introduced by Cockburn and Shu. Finally, before ending with a conclusion in Section 7, we give, in Section 6, some critical comparisons between the described reconstruction techniques by using several numerical experiments.

## 2. DG FINITE ELEMENT METHOD

The ultimate goal of this work is to check out some reconstruction techniques for DG methods. Thus, for the sake of brevity, we restrict our attention to two-dimensional, linear, scalar advection equations. The extension to three-dimensional general conservation laws is an ongoing work.

Hence, we consider the hyperbolic-type equation of the form

$$\frac{\partial u}{\partial t} + \nabla . f(u) = 0 \quad \text{in } \Omega \times (0, T) \tag{1}$$

with the initial conditions

$$u(x, 0) = u^0(x) \quad \text{in } \Omega \tag{2}$$

and appropriate boundary conditions. Here $f(u) = u\beta$ where $u = u(x, t)$ is a scalar unknown representing a concentration for example, $\beta = (\beta_1(x), \ldots, \beta_d(x))$, $(d = 1, 2)$ is a given vector field, $\Omega \subset \mathbb{R}^d$ and $(0, T)$ is a given time interval.

### 2.1. Space integration

In this presentation of the DG method, some materials are drawn from these works [11, 12, 14–16]. The polygonal domain $\Omega$ is discretized into a mesh $\mathscr{T}_h$ consisting of quadrilaterals or triangles where $h$ refers to the maximal element diameter. We also denote by $\mathscr{N}_K$ the number of vertices of the discretized element $K$.

The DG method is based on using the following discontinuous finite element space:

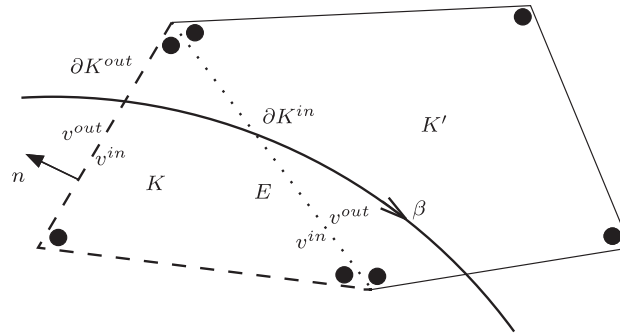$$V_h = \{v \in L^\infty(\Omega) : v_h|_K \in V(K) \ \forall K \in \mathscr{T}_h\}$$

Figure 1. Inflow and outflow boundaries with the local degrees of freedom.

where $V(K)$ is the space of linear $P^1$ (resp. quadratic $Q^1$) polynomials if $K$ is a triangle (resp. quadrilateral). In this work, we are restricted to polynomials of degree one.

In order to define the upwind technique [6], we need to split the boundary $\partial K$ of a discretized element $K$ into an inflow part $\partial K^{\text{in}}$ and an outflow part $\partial K^{\text{out}}$ defined by:

$$\partial K^{\text{in}} = \{x \in \partial K : n(x).\beta < 0\}$$

$$\partial K^{\text{out}} = \{x \in \partial K : n(x).\beta \geqslant 0\}$$

where $n(x)$ denotes the unit outward normal to $\partial K$ (see Figure 1).

Let $E$ be a common edge between any two adjacent elements $K$ and $K'$. Since discontinuity for any function $v \in V_h$ is allowed across interelement boundaries, we need to define the jump discontinuity of $v$ across $E$. We introduce the notations $v^{\text{in}}$ and $v^{\text{out}}$ to denote, respectively, the inner and the outer values of $v$ over $E$ with respect to $K$ (see Figure 1).

The formulation obtained by using the DG method is formulated by multiplying Equation (1) by a sufficiently smooth test function $v$ and by integrating by parts over an element $K \in \mathscr{T}_h$

$$\int_K \frac{\partial u}{\partial t} v \, dx - \int_K f(u).\nabla v \, dx + \int_{\partial K} v f(u).n \, d\ell = 0 \qquad (3)$$

Then, we replace $u$ by the approximate solution $u_h$ which can be expressed as follows:

$$u_h(x, t) \equiv u_h(x, t)|_K = \sum_{j=1}^{\mathscr{N}_K} u_{K,j}(t)$$

where $\varphi_{K,j}$ are some test functions in $V_h$ that form a basis for the local approximation space $V(K)$ and where $u_{K,j}(t)$ are associated degrees of freedom. The standard finite element shape functions may be chosen.

Due to the discontinuity of $u$ across $\partial K$, the flux function $f(u)$ is approximated by solving a one-dimensional Riemann problem. In our case $f(u) = u\beta$ is a linear function, consequently

the Riemann solver is evident (see, e.g. Reference [5]), that is,

$$f(u) = \begin{cases} f(u^{\text{in}}) & \text{over } \partial K^{\text{out}} \\ f(u^{\text{out}}) & \text{over } \partial K^{\text{in}} \end{cases}$$

By replacing $v$ successively by the test functions $\varphi_{K,i}$, $i = 1, \ldots, \mathcal{N}_K$, the weak formulation (3) takes the following form:

$\forall K \in \mathcal{T}_h$, we seek the approximation solution $u_h \equiv u_h|_K \in V_h$ with the initial data (2) such that,

$$\sum_{j=1}^{\mathcal{N}_K} \frac{\mathrm{d}u_{K,j}}{\mathrm{d}t} \int_K \varphi_{K,i}\varphi_{K,j} \, \mathrm{d}x = \sum_{j=1}^{\mathcal{N}_K} \left( u_{K,j} \int_K \varphi_{K,j}\beta.\nabla\varphi_{K,i} \, \mathrm{d}x - u_{K,j}^{\text{in}} \int_{\partial K^{\text{out}}} \varphi_{K,i}\varphi_{K,j}\beta.n \, \mathrm{d}\ell \right.$$

$$\left. - u_{K,j}^{\text{out}} \int_{\partial K^{\text{in}}} \varphi_{K,i}\varphi_{K,j}\beta.n \, \mathrm{d}\ell \right) \tag{4}$$

Note that the inner values of the functions $\varphi_{K,i}$ are taken in the integrals across the boundaries of $K$ in Equation (4).

## 2.2. Time integration

The DG approximation leads to a system of $\mathcal{N}_K$ ordinary differential equations over each element $K \in \mathcal{T}_h$. After inverting the local mass matrix, which corresponds to the integrals on the left-hand side of Equation (4), this system can be rewritten in matrix form as follows:

$$\frac{\mathrm{d}U_K}{\mathrm{d}t} = \mathscr{A}(U_K^{\text{in}}, U_K^{\text{out}}) \tag{5}$$

where $U_K^{\text{in}}$ is a vector of dimension $\mathcal{N}_K$ containing the cell unknowns $u_{K,j}$ and $\mathscr{A}$ represents the components of the right-hand side of Equation (4) multiplied by the inverse of the mass matrix. In order to approximate system (5), we subdivide the time interval $[0, T]$ into a finite number of sub-interval $[t^n, t^{n+1}]$. Let $\Delta t = t^{n+1} - t^n$ denote the time step. We specify the following schemes:

*2.2.1. Forward Euler method.* A simple approach is to use Euler forward time discretization scheme. However, Chavent and Cockburn [17] showed that without using a suitable slope limiter this scheme is unconditionally unstable. Thus, the reconstruction process is crucial in order to stabilize the scheme. Therefore a stable DG computation procedure consists of the following two steps:

1. Calculation of $\widetilde{U}_K^{n+1}$ for a given $u_h^n$ as follows:

$$\widetilde{U}_K^{n+1} = U_K^n + \Delta t \mathscr{A}(U_K^{\text{in},n}, U_K^{\text{out},n}) \quad \forall K \in \mathcal{T}_h$$

2. Reconstruction of the updated solution $\widetilde{U}_K^{n+1}$ by applying

$$U_K^{n+1} = \mathscr{L}(\widetilde{U}_L^{n+1}, L \in \mathscr{V}(K)) \quad \forall K \in \mathcal{T}_h$$

where $\mathscr{V}(K)$ is the set of adjacent elements of $K$ and $\mathscr{L}$ denotes a slope limiting operator to be discussed in the next section.

*2.2.2. Explicit Runge–Kutta method.* A second-order accuracy in time may be obtained by using an explicit Runge–Kutta method. The time-stepping algorithm reads in four steps as follows:

1. Compute an intermediate function $\widetilde{U}_K^{n+1/2}$ for given $u_h^n$,

$$\widetilde{U}_K^{n+1/2} = U_K^n + \frac{\Delta t}{2} \mathscr{A}(U_K^{\mathrm{in},n}, U_K^{\mathrm{out},n}) \quad \forall K \in \mathscr{T}_h$$

2. Apply the slope limiter operator, $U_K^{n+1/2} = \mathscr{L}(\widetilde{U}_L^{n+1/2}, L \in \mathscr{V}(K)) \quad \forall K \in \mathscr{T}_h$.
3. Compute $\widetilde{U}_K^{n+1}$ for given $u_h^n$ and $u_h^{n+1/2}$,

$$\widetilde{U}_K^{n+1} = U_K^n + \Delta t \, \mathscr{A}(U_K^{\mathrm{in},n+1/2}, U_K^{\mathrm{out},n+1/2}) \quad \forall K \in \mathscr{T}_h$$

4. Apply the slope limiter operator, $U_K^{n+1} = \mathscr{L}(\widetilde{U}_L^{n+1}, L \in \mathscr{V}(K)) \quad \forall K \in \mathscr{T}_h$.

*2.2.3. Simplified Runge–Kutta.* Due to the expensive computing cost for Riemann solvers as well as slope limiting process, a simplified version of the above Runge–Kutta method was introduced. This schema is used in these works [18–21]. Thus, the following three-steps algorithm can be used instead:

1. Compute an intermediate function $\widetilde{U}_K^{n+1/2}$ for given $u_h^n$,

$$\widetilde{U}_K^{n+1/2} = U_K^n + \frac{\Delta t}{2} \mathscr{A}(U_K^{\mathrm{in},n}, U_K^{\mathrm{in},n}) \; \forall K \in \mathscr{T}_h$$

   Note that in this step, the intermediate functions are calculated by means of local interior values of $u_h^n$ which makes this computation local on $K$.
2. Compute $\widetilde{U}_K^{n+1}$ for given $u_h^n$ and $\widetilde{u}_h^{n+1/2}$,

$$\widetilde{U}_K^{n+1} = U_K^n + \Delta t \, \mathscr{A}(\widetilde{U}_K^{\mathrm{in},n+1/2}, \widetilde{U}_K^{\mathrm{out},n+1/2}) \quad \forall K \in \mathscr{T}_h$$

3. Apply the slope limiter operator, $U_K^{n+1} = \mathscr{L}(\widetilde{U}_L^{n+1}, L \in \mathscr{V}(K)) \quad \forall K \in \mathscr{T}_h$.

# 3. DATA RECONSTRUCTION

In this section, we focus on the slope limiter introduced by Chavent and Jaffré [12]. This limiter can be interpreted as a generalization of Van Leer's MUSCL limiter [4]. The essential idea of this technique is to impose some local constraints in a geometric manner so that the reconstructed solution satisfies an appropriate maximum principle. Before starting with the multi-dimensional case, we present the slope limiter with one variable in space.

## 3.1. One-dimensional slope limiter

Let us now denote by $K_i = ]x_{i-1/2}, x_{i+1/2}[$ the sub-intervals of the one-dimensional space discretization. The sought function $u_h$ is approximated by means of piecewise linear functions. We denote by $\overline{u}_i$ the average of $u_h$ over $K_i$ which is indeed the midpoint of the two boundary

nodal values, that is,

$$\overline{u}_i = \frac{1}{|K_i|} \int_{K_i} u_h \, \mathrm{d}x = \frac{1}{2}(u_{i-1/2} + u_{i+1/2})$$

Given a function $\widetilde{u}_h \in V_h$, we want to define its slope limited version $u_h = \mathscr{L}(\widetilde{u}_h) \in V_h$ in such a way that, over an element $K_i$, $u_h$ depends only on $\widetilde{u}_h$ over the elements $K_{i-1}$, $K_i$ and $K_{i+1}$. Thus, the slope limiter is defined as the solution $U_i = (u_{i-1/2}, u_{i+1/2})$ of the following conditions:

1. Conservation of mass:

$$\overline{u}_i = \tfrac{1}{2}(u_{i-1/2} + u_{i+1/2}) = \tfrac{1}{2}(\widetilde{u}_{i-1/2} + \widetilde{u}_{i+1/2})$$

2. Avoid creating local extremum for some $\alpha \in [0, 1]$, one requires that:

$$(1-\alpha)\overline{u}_i + \alpha \min(\overline{u}_{i-1}, \overline{u}_i) \leqslant u_{i-1/2} \leqslant (1-\alpha)\overline{u}_i + \alpha \max(\overline{u}_{i-1}, \overline{u}_i)$$

$$(1-\alpha)\overline{u}_i + \alpha \min(\overline{u}_i, \overline{u}_{i+1}) \leqslant u_{i+1/2} \leqslant (1-\alpha)\overline{u}_i + \alpha \max(\overline{u}_i, \overline{u}_{i+1})$$

   The parameter $\alpha$ controls the degree of constraints on the slopes, that is, the added numerical viscosity (see Figure 2).

3. Minimum modification of $\widetilde{u}_h$ : $U_i$ is chosen as close as possible to $\widetilde{U}_i$ with respect to the $L^2$ norm, that is,

$$\|U_i - \widetilde{U}_i\|_2 \text{ is minimal}$$

The above problem can also be rewritten using another more familiar form:

$$u_{i-1/2} = \overline{u}_i - \mathscr{M}(\overline{u}_i - \widetilde{u}_{i-1/2}, \alpha(\overline{u}_i - \overline{u}_{i-1}), \alpha(\overline{u}_{i+1} - \overline{u}_i))$$

$$u_{i+1/2} = \overline{u}_i + \mathscr{M}(\widetilde{u}_{i+1/2} - \overline{u}_i, \alpha(\overline{u}_i - \overline{u}_{i-1}), \alpha(\overline{u}_{i+1} - \overline{u}_i))$$

where $\mathscr{M}$ is the well-known minmod function [22],

$$\mathscr{M}(a_1, a_2, a_3) = \begin{cases} s \min_{1 \leqslant i \leqslant 3} |a_i| & \text{if } s = \mathrm{sign}(a_1) = \mathrm{sign}(a_2) = \mathrm{sign}(a_3) \\ 0 & \text{otherwise} \end{cases} \tag{6}$$

Note that $\mathscr{M}$ needs to be applied only once since

$$\mathscr{M}(\overline{u}_i - \widetilde{u}_{i-1/2}, \alpha(\overline{u}_i - \overline{u}_{i-1}), \alpha(\overline{u}_{i+1} - \overline{u}_i))$$

$$= \mathscr{M}(\widetilde{u}_{i+1/2} - \overline{u}_i, \alpha(\overline{u}_i - \overline{u}_{i-1}), \alpha(\overline{u}_{i+1} - \overline{u}_i))$$

By choosing specific values for $\alpha$, some well-known slope limiters are obtained.

- For $\alpha = 0$, the slope limiter enforces constant piecewise approximations. Therefore, the scheme boils down to first-order Godunov finite difference method [1].
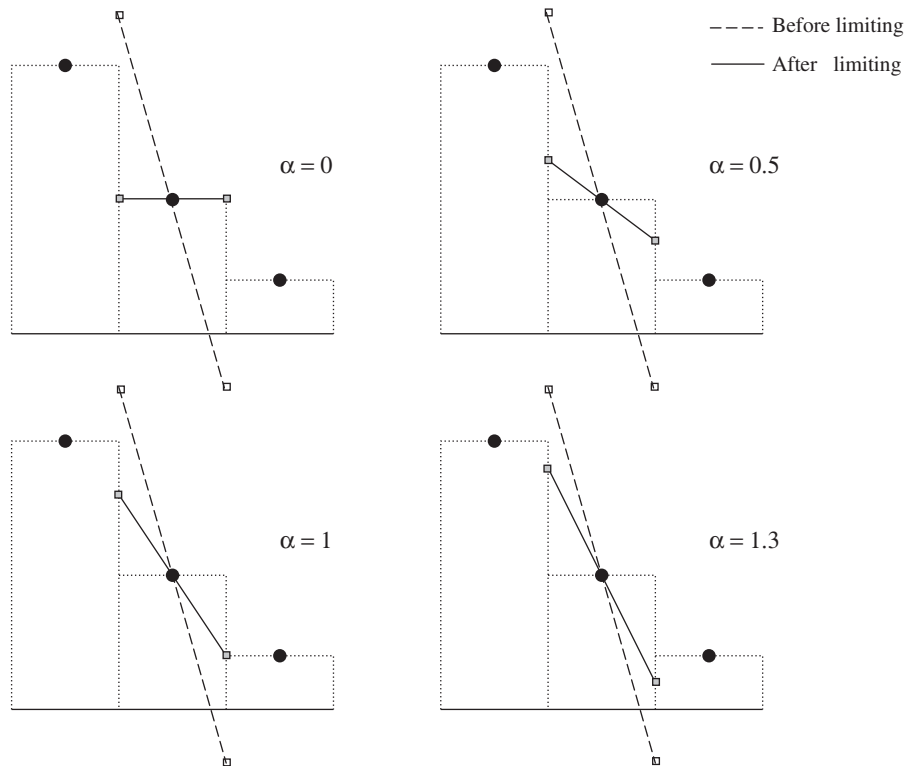
Figure 2. Slope limiting with different values for α.

• For $\alpha = 1/2$, we obtain the slope limiter of the MUSCL schemes of Van Leer [4].
• For $\alpha = 1$, we find a less restrictive limiter found in References [12, 23].

By using Harten's TVD conditions [24], Gowda and Jaffré [11] proved the stability of the DG method combined with this slope limiter for $\alpha \in [0, 1]$. However, by taking $\alpha$ slightly greater than one, it is found that this slope limiter for smooth initial conditions behaves in a similar manner as the TVB generalized slope limiter introduced by Cockburn and Shu [8, 9].

### 3.2. Multi-dimensional slope limiter

The extension of the slope limiter to the multi-dimensional case is formulated in such a way that in each cell $K$ each state variable at a vertex $A_i$ lies between the cell averages of all neighbouring elements containing $A_i$ as a vertex. For any $K \in \mathcal{T}_h$, we introduce the following notations:

$$T(A) = \{K \in \mathcal{T}_h | \; A \text{ is a vertex of } K\}$$

$$U_K = (u_{K,i})_{i=1,\ldots,\mathcal{N}_K}$$
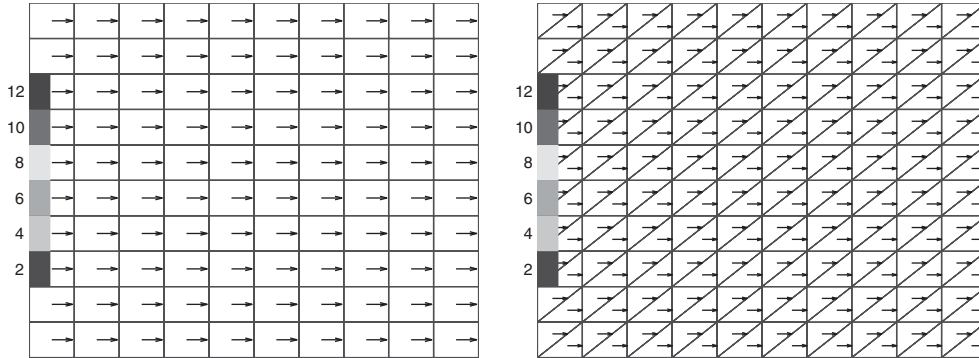
$$\overline{u}_K = \frac{1}{|K|} \int_K u_h \, dx$$

Figure 3. Uniform grids with Dirichlet boundary conditions at the left-hand side of the domain.

$$\overline{u}_{\min,i} = \min_{K \in T(A_i)} \overline{u}_K$$

$$\overline{u}_{\max,i} = \max_{K \in T(A_i)} \overline{u}_K$$

The slope limiting process seeks $U_K \in V(K)$, $\forall K \in \mathscr{T}_h$, as the solution of the following least squares problem:

$$\min_W \| W - \widetilde{U}_K \|_2, \quad \text{subject to the linear constraints}$$

$$\overline{w} = \frac{1}{\mathscr{N}_K} \sum_{j=1}^{\mathscr{N}_K} w_i = \overline{u}_K \tag{7}$$

$$(1-\alpha)\overline{u}_K + \alpha\overline{u}_{\min,i} \leqslant w_i \leqslant (1-\alpha)\overline{u}_K + \alpha\overline{u}_{\max,i}, \quad i = 1, \dots, \mathscr{N}_K$$

where $\alpha \in [0, 1]$. It is easy to check that this minimization problem has a unique solution [12]. See Appendix A for a robust algorithm.

We have found that this slope limiter sometimes fails to smear completely the spurious oscillations. Its weak point is that it does not prevent creating new extrema at the midpoints of the grid edges. In other words, it is possible to obtain a value of the average over an edge $E$ which is beyond the cell averages of the two adjacent grid elements having $E$ as a common edge. As a result, we could have regions where a local maximum principle is violated.

### 3.3. Numerical test

In order to illustrate the drawback of this slope limiter, we consider a very simple numerical test. Let $\Omega = (0, 10) \times (0, 10)$ be the computational domain and $\beta = (1, 0)$ be the velocity field. Two uniform grids of rectangles and triangles are considered with space steps $\Delta x = \Delta y = 1$ (see Figure 3). The scalar convection equation (1) is considered with zero initial conditions and a decreasing piecewise constant Dirichlet boundary conditions on the left-hand side of the domain (Figure 3). Even though, this problem is physically one-dimensional, the above-described multi-dimensional slope limiter is used. In Figures 4 and 5, the cell average values of
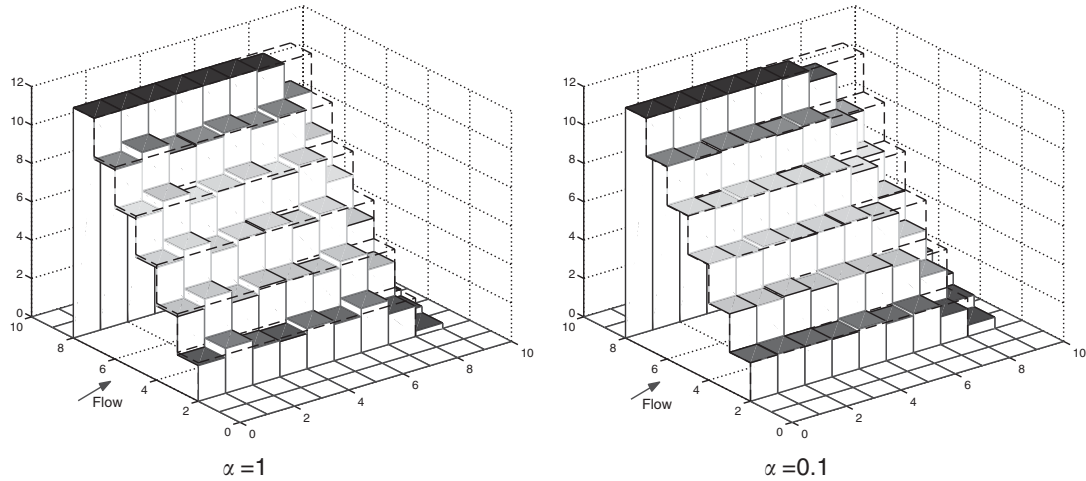
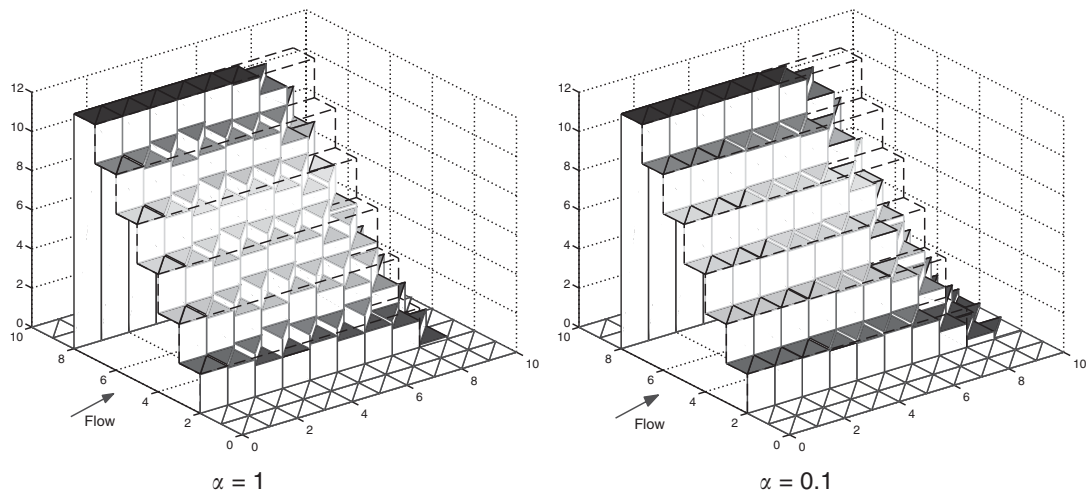Figure 4. DG solutions on a rectangular grid at $T = 8$ with different values of $\alpha$.



Figure 5. DG solutions on a triangular grid at $T = 8$ with different values of $\alpha$.

the solutions obtained by the DG method are presented without any graphical smoothing. The dashed lines represent the profile of the exact solution. It is clear that the slope limiter for both triangular and rectangular grids does not completely eliminate the non-physical oscillations. Decreasing the value of $\alpha$ will smear oscillations, however the scheme becomes more diffusive. The simplified Runge–Kutta method is used for the time integration. The time step is chosen so that the CFL $= \beta_x \Delta t / \Delta x$ condition is equal to 0.9 for rectangular grid and 0.4 for triangular grid. It should be noted that all temporal schemes introduced in the previous section produce similar results.

## 4. MODIFIED SLOPE LIMITER

A remedy for this drawback is possible by preventing the reconstruction to produce any new extrema at the midpoints of edges within each cell. This approach has an important physical interpretation since it limits the interelement numerical fluxes rather than the function values at the grid vertices. However, the previous slope limiter does not satisfy this property. In what follows, we introduce two new slope limiters for rectangular and triangular unstructured grids that satisfy this property.

### 4.1. Slope limiting for rectangular elements

Up to our knowledge not many slope limiters are available in literature for quadrangular elements. The slope limiter proposed by Cockburn and Shu [16] for rectangular grids is essentially designed for linear cell approximations ($P^1$). Our slope limiter is related in some way to that limiter, however, we use piecewise quadratic polynomials ($Q^1$) to approximate the solution.

In order to formulate our slope limiter, we choose an arbitrary rectangular element $K_0$ surrounded by its neighbours $K_i$, $i = 1, \ldots, 4$, as illustrated in Figure 6.

We denote by $A_{i,j}$ the midpoints of the edge $[A_i, A_j]$ and by $u_{i,j}$ the state average, in $K_0$, over the edge $[A_i, A_j]$. Indeed, $u_{i,j}$ is the midpoint of $[u_i, u_j]$.

Thus, non-physical oscillations at the midpoints $A_{i,j}$ can be avoided by enforcing the edge average $u_{i,j}$ to be within the averages of cells containing $[A_i, A_j]$ as a common edge, that is,

$$(1-\alpha)\overline{u}_{K_0} + \alpha \min(\overline{u}_{K_1}, \overline{u}_{K_0}) \leqslant u_{1,2} \leqslant (1-\alpha)\overline{u}_{K_0} + \alpha \max(\overline{u}_{K_1}, \overline{u}_{K_0})$$
$$(1-\alpha)\overline{u}_{K_0} + \alpha \min(\overline{u}_{K_0}, \overline{u}_{K_3}) \leqslant u_{3,4} \leqslant (1-\alpha)\overline{u}_{K_0} + \alpha \max(\overline{u}_{K_0}, \overline{u}_{K_3})$$
(8)

$$(1-\alpha)\overline{u}_{K_0} + \alpha \min(\overline{u}_{K_2}, \overline{u}_{K_0}) \leqslant u_{2,3} \leqslant (1-\alpha)\overline{u}_{K_0} + \alpha \max(\overline{u}_{K_2}, \overline{u}_{K_0})$$
$$(1-\alpha)\overline{u}_{K_0} + \alpha \min(\overline{u}_{K_0}, \overline{u}_{K_4}) \leqslant u_{4,1} \leqslant (1-\alpha)\overline{u}_{K_0} + \alpha \max(\overline{u}_{K_0}, \overline{u}_{K_4})$$
(9)

Therefore, a natural choice for the slope limiter would be to add these constraints to system (7). Unfortunately, the resolution of the resulting minimization problem is computationally very expensive. A key solution to overcome this difficulty is to replace the inequality constraints given in (8) and (9) by equality constraints. This can be done by using the following splitting technique:

1. We first reconstruct the state averages at the midpoints of the edges by using a direction splitting. Let us start in the $x$ direction, for example. Here, only information about the neighbour cells $K_1$ and $K_3$ is needed. Due to mass balance, we should have:

$$\tfrac{1}{2}(u_{1,2} + u_{3,4}) = \overline{u}_{K_0}$$

By combining the local constraints on the edge averages $u_{1,2}$ and $u_{3,4}$ (8), the resulting problem is thus reduced to a one-dimensional linear reconstruction. Consequently, we apply the one-dimensional slope limiter $\mathcal{M}$ previously discussed.

$$u_{1,2} = \overline{u}_{K_0} - \mathcal{M}(\overline{u}_{K_0} - \widetilde{u}_{1,2}, \alpha(\overline{u}_{K_0} - \overline{u}_{K_1}), \alpha(\overline{u}_{K_3} - \overline{u}_{K_0}))$$

$$u_{3,4} = \overline{u}_{K_0} + \mathcal{M}(\widetilde{u}_{2,3} - \overline{u}_{K_0}, \alpha(\overline{u}_{K_0} - \overline{u}_{K_1}), \alpha(\overline{u}_{K_3} - \overline{u}_{K_0}))$$
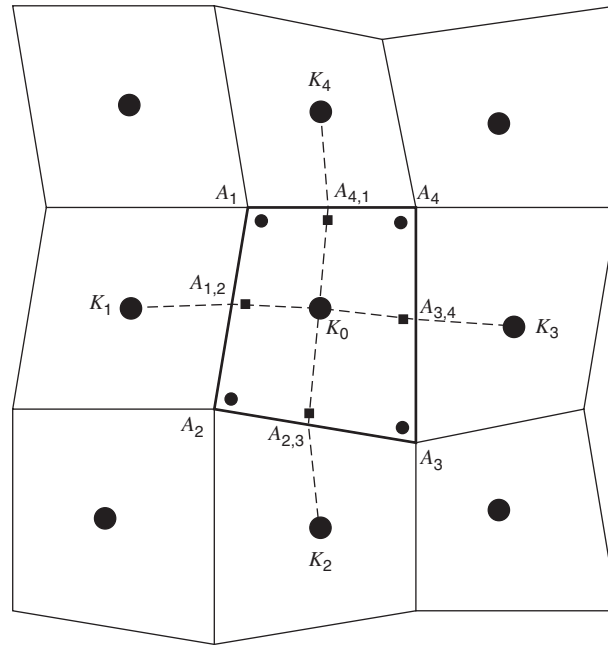
Figure 6. Illustration of limiting for quadrilateral elements.

Similarly, we reconstruct $u_{2,3}$ and $u_{4,1}$ in the $y$ direction by applying:

$$u_{2,3} = \overline{u}_{K_0} - \mathcal{M}(\overline{u}_{K_0} - \widetilde{u}_{2,3}, \alpha(\overline{u}_{K_0} - \overline{u}_{K_2}), \alpha(\overline{u}_{K_4} - \overline{u}_{K_0}))$$

$$u_{4,1} = \overline{u}_{K_0} + \mathcal{M}(\widetilde{u}_{4,1} - \overline{u}_{K_0}, \alpha(\overline{u}_{K_0} - \overline{u}_{K_2}), \alpha(\overline{u}_{K_4} - \overline{u}_{K_0}))$$

2. The aim now is to reconstruct the cell values at the vertices. Indeed, the midpoint edge averages $u_{i,j}$, which are already computed in the previous step, are not sufficient to uniquely reconstruct the nodal values $u_i$ of the function $u_h$ over $K$. Consequently, we use the following constraints which ensure the mass conservation over each edge within the cell:

$$
\begin{aligned}
u_1 + u_2 &= 2u_{1,2} \\
u_2 + u_3 &= 2u_{2,3} \\
u_3 + u_4 &= 2u_{3,4} \\
u_4 + u_1 &= 2u_{4,1}
\end{aligned}
\tag{10}
$$

These constraints are simply due to the linear approximation of the function $u_h$ along the edges. The values at the vertices $u_i$ are thus reconstructed by combining the equality constrains (10) with those given in system (7). Therefore, the resulting slope limiter guarantees that no new extrema can be created at the vertices as well as at the midpoints of the edges within each cell. On the other hand, the obtained optimization problem

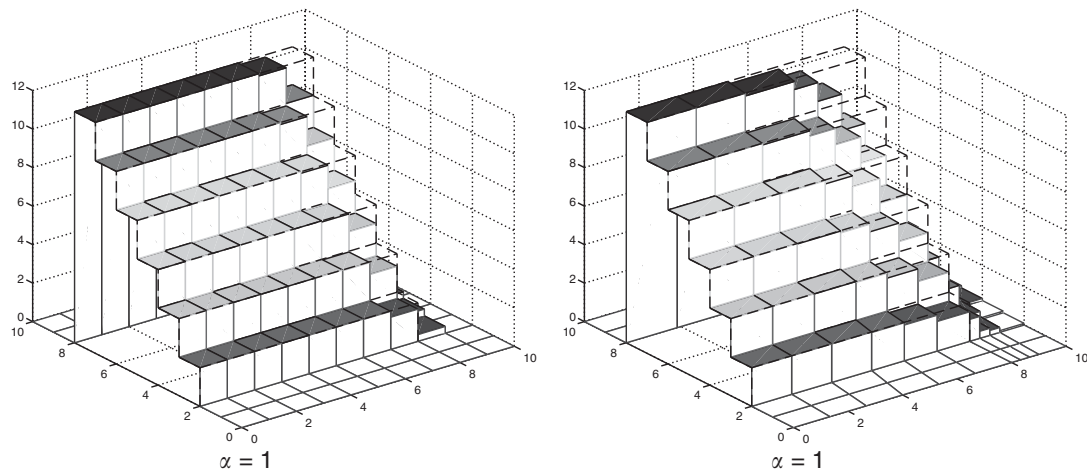*Int. J. Numer. Meth. Engng* 2004; **61**:2566–2593

Figure 7. DG results obtained when using the modified slope limiter for
a structured and an unstructured grid.

is very simple to solve. The system of equality constraints (10) is of rank 3, thus the problem can be reduced to a minimization problem with only one variable. Further, the optimal solution can be attained without iterating.

*4.1.1. Numerical test.* The same test problem given in the previous section is now approximated by using the DG method with our modified slope limiter. Computations are done for a structured grid as well as for an unstructured grid of trapezoids. The time step is taken to be 0.62 for the unstructured mesh, so that the Courant number does not exceed 0.9 in the active elements of the grid. Results depicted in Figure 7 show that our modified slope limiter completely eliminates spurious oscillations with minimal numerical smearing, that is, with $\alpha = 1$.

### 4.2. Slope limiting for triangular elements

The construction of the slope limiting operator for triangular elements follows a similar approach used for rectangular elements. It is proved in Reference [25] (see also Reference [18]) that an appropriate local maximum principle is satisfied by ensuring that no new extrema are created at the midpoints of the grid edges. Consequently, the proposed slope limiting operator aims to eliminate oscillations at the midpoints of edges within each cell. To describe the slope limiting procedure, let us consider a triangular element $K_0$ surrounded by its neighbourhoods $K_i$, $i = 1, \ldots, 3$. (Figure 8). The notations for the vertices and for the midpoints of the edges are the same as those used for quadrangles (Figure 6). The slope limiting process consists of two main operations:

1. The aim in the first stage is to reconstruct the average values $\widetilde{u}_{i,j}$ at the midpoints of the edges. A necessary condition that must be satisfied is the local mass conservation. To obey a local maximum principle, some constraints are imposed to ensure that each reconstructed $u_{i,j}$ is between the cell averages of the two adjacent elements. To have a
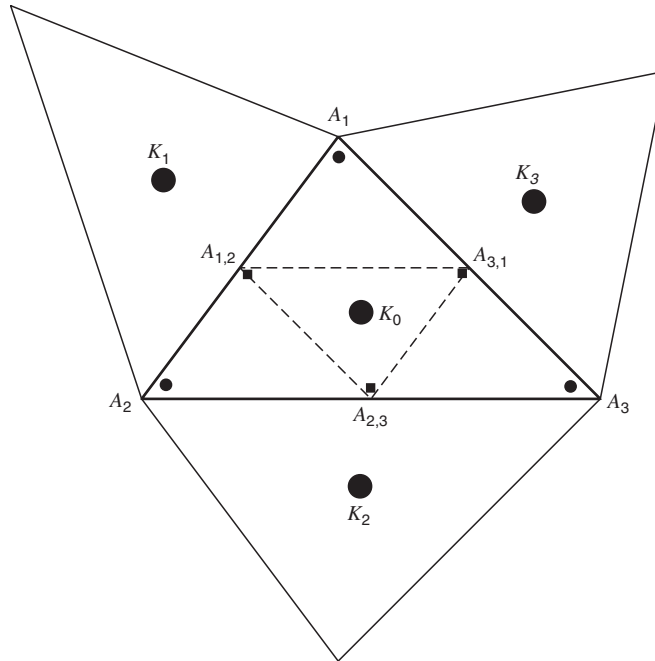
Figure 8. Illustration of limiting for triangular elements.

less restrictive limiting, the reconstructed $u_{i,j}$ are kept as close as possible to the initial state values $\widetilde{u}_{i,j}$. The resulting optimization problem to solve is therefore:

For given initial state values $\widetilde{U}_{K_0} = (\widetilde{u}_{1,2}, \widetilde{u}_{2,3}, \widetilde{u}_{3,1})$, find $U_{K_0}$ the solution of the problem:

$$\min_{W} \|W - \widetilde{U}_{K_0}\|_2 \text{ subject to the linear constraints}$$

$$\overline{w} = \tfrac{1}{3}(w_{1,2} + w_{2,3} + w_{3,1}) = \overline{u}_{K_0}$$

$$(1 - \alpha)\overline{u}_{K_0} + \alpha \min(\overline{u}_{K_1}, \overline{u}_{K_0}) \leqslant w_{1,2} \leqslant (1 - \alpha)\overline{u}_{K_0} + \alpha \max(\overline{u}_{K_1}, \overline{u}_{K_0}) \qquad (11)$$

$$(1 - \alpha)\overline{u}_{K_0} + \alpha \min(\overline{u}_{K_2}, \overline{u}_{K_0}) \leqslant w_{2,3} \leqslant (1 - \alpha)\overline{u}_{K_0} + \alpha \max(\overline{u}_{K_2}, \overline{u}_{K_0})$$

$$(1 - \alpha)\overline{u}_{K_0} + \alpha \min(\overline{u}_{K_3}, \overline{u}_{K_0}) \leqslant w_{3,1} \leqslant (1 - \alpha)\overline{u}_{K_0} + \alpha \max(\overline{u}_{K_3}, \overline{u}_{K_0})$$

See Appendix A for the solution of this problem.

2. Unlike the case of rectangular elements, the state values at the cell vertices can be directly computed by using the reconstructed midpoint edge values $u_{i,j}$. Therefore, a quite simple system of linear equation has to be solved:

$$u_1 + u_2 = 2u_{1,2}$$

$$u_2 + u_3 = 2u_{2,3} \qquad (12)$$

$$u_3 + u_1 = 2u_{3,1}$$

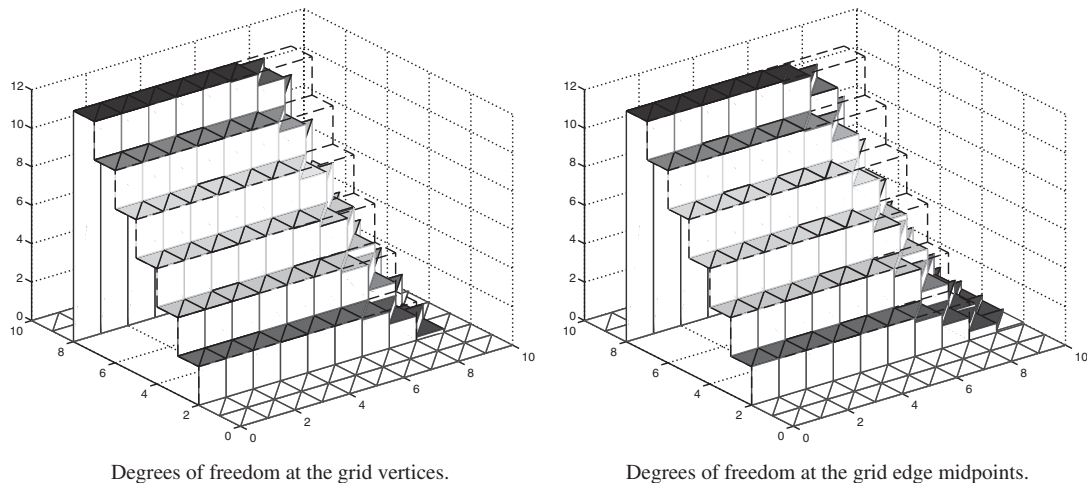Degrees of freedom at the grid vertices.                    Degrees of freedom at the grid edge midpoints.

Figure 9. DG results for triangular grid with different degrees of freedom.

This procedure insures the maximum principle at the triangle vertices since the resulting re-constructions at the vertices satisfy the constraints of Chavent–Jaffré slope limiter given in system (7).

*4.2.1. Different degrees of freedom.* Depending on the choice of the basis functions $\varphi_{K,i}$, formula (4) will produce, together the chosen time discretization, different vectors $\widetilde{U}_K$. By defining the degrees of freedom at the midpoints of the grid edges the proposed slope limiter is more convenient since, in this case, no construction at the vertices is required. Indeed, local approximation of the solution obtained by taking the degrees of freedom either at the vertices or at the midpoints of edges has the same order of accuracy. However, we have found that the numerical solution obtained by the later approximation is more diffusive. This drawback is due to the upwinding. In fact by using the midpoints of edges as degrees of freedom the only information transmitted from one element to its neighbours, in the upwind direction, is the value at the midpoint of their common edge. On the other hand, by taking the degrees of freedom at the vertices, the two nodal values at the extremities of the common edges are transmitted. This approach is indeed more precise since it provides information about the state gradient along the edge.

*4.2.2. Numerical test.* In Figure 9, we present the numerical results obtained by the DG method which is stabilized by using our modified slope limiter (limiter previously described). The obtained solution is free from any spurious oscillations even with minimal artificial diffusion ($\alpha = 1$). However, comparison between solutions obtained by the DG method with degrees of freedom at the vertices and at the midpoints of the edges, as illustrated in Figure 9, shows that the later approach is more diffusive.

## 5. EXISTING SLOPE LIMITERS

In this section, we briefly review two slope limiters introduced by Cockburn and Shu [16] for rectangular and triangular grids. We restrict the presentation for $P^1$ piecewise approximation functions.

### 5.1. Rectangular grids

The approximation solution $u_h(x, y, t)$ over each rectangular element $[x_{i-1/2}, x_{i+1/2}] \times [x_{i-1/2}, x_{i+1/2}]$ in a cartesian grid is approximated by means of $P^1$ polynomials. A convenient choice of the degrees of freedom is the cell average $\overline{u}$ and the two slopes of the state function $u_x$ and $u_y$ in the $x$ and $y$ directions, respectively. Thus, over each element we have:

$$u_h(x, y, t) = \overline{u}(t) + u_x(t)\phi_i(x) + u_y(t)\psi_i(t) \tag{13}$$

where

$$\phi_i(x) = \frac{x - x_i}{\Delta x/2}, \quad \psi_i(x) = \frac{y - y_i}{\Delta y/2}$$

*5.1.1. Limiting.* The reconstruction of $u_x$ and $u_y$ is carried out sequentially in the $x$ and $y$ directions by applying a one-dimensional slope limiter. Cockburn and Shu [7, 8] proposed to use the TVB *generalized* slope limiter for the reconstruction. Since our aim is to have the numerical solution free from any spurious oscillation, we will use the one-dimensional slope limiter proposed in this paper (6). Therefore, within each cell $u_x$ and $u_y$ are, respectively, replaced by:

$$\mathcal{M}(u_x, \alpha(\overline{u}_{i+1,j} - \overline{u}_{i,j}), \alpha(\overline{u}_{i,j} - \overline{u}_{i-1,j}))$$

$$\mathcal{M}(u_y, \alpha(\overline{u}_{i,j+1} - \overline{u}_{i,j}), \alpha(\overline{u}_{i,j} - \overline{u}_{i,j-1}))$$

### 5.2. Triangular grids

The approximation solution $u_h(x, y, t)$ is approximated by means of piecewise linear polynomials. The degrees of freedom are the state values at the midpoints of the grid edges.

### 5.3. Limiting

To describe Cockburn and Shu's slope limiter, we use the same notations as in Reference [16]. For an arbitrary triangle $K_0$ and its surrounding neighbours $K_i$, $i = 1, \ldots, 3$, the notations $b_i$, $i = 0, \ldots, 3$ and $m_i$, $i = 1, \ldots, 3$ refer, respectively, to the barycentres of the triangles and the midpoints of the edges within $K_0$ (see Figure 10).

Choosing any edge midpoint $m_1$, we get:

$$m_1 - b_0 = \alpha_1(b_1 - b_0) + \alpha_2(b_2 - b_0) \quad \text{for some } \alpha_1, \alpha_2 \in \mathbb{R}^2$$

Then, for any *linear function* $u_h$ we can write:

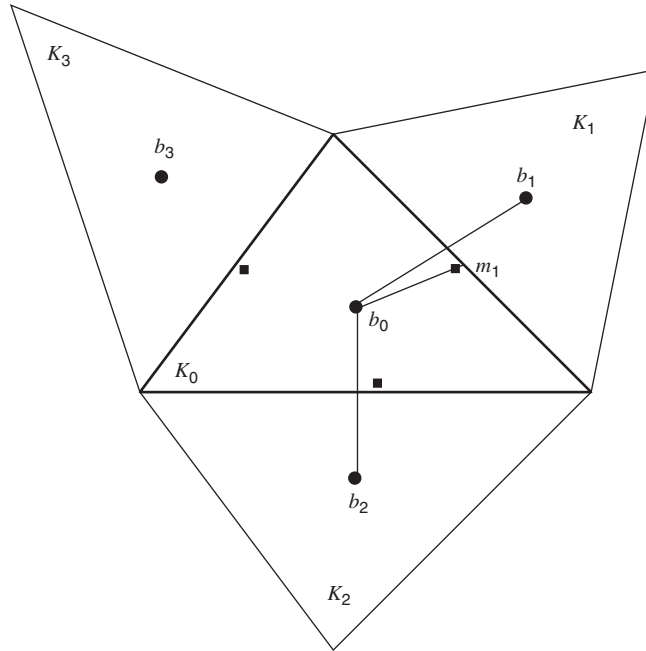$$u_h(m_1) - u_h(b_0) = \alpha_1(u_h(b_1) - u_h(b_0)) + \alpha_2(u_h(b_2) - u_h(b_0)) \tag{14}$$

Figure 10. Cockburn and Shu limiting illustration for triangular elements.

Since the cell average $\overline{u}_{K_i}$ is nothing but the value of the function at the barycentre $u_h(b_i)$, (14) is rewritten as follows:

$$
\begin{aligned}
\widetilde{u}_h(m_1, K_0) &\equiv u_h(m_1) - \overline{u}_{K_0} \\
&= \alpha_1(\overline{u}_{K_1} - \overline{u}_{K_0}) + \alpha_2(\overline{u}_{K_2} - \overline{u}_{K_0}) \\
&\equiv \Delta\overline{u}(m_1, K_0)
\end{aligned}
\tag{15}
$$

To describe the slope limiting operator $\Lambda\Pi_h$, we consider any *piecewise linear function* $u_h$. Indeed Equation (15) is no more valid. By using some basis functions $\phi_i$, $u_h$ can be expressed over $K_0$ as follows:

$$
u_h(x, y) = \sum_{i=1}^3 u_h(m_i)\phi_i(x, y) = \overline{u}_{K_0} + \sum_{i=1}^3 \widetilde{u}_h(m_i, K_0)\phi_i(x, y)
$$

First, we compute the quantities

$$
\Delta_i = \mathcal{M}(\widetilde{u}_h(m_i, K_0), v\Delta\overline{u}(m_i, K_0)) \quad \text{for some } v > 1
$$

by using the *minmod* function described above. Note that Cockburn and Shu used instead their modified TVB minmod function. Consequently, reconstruction is carried out according to the

following two cases:

1. If $\sum_{i=1}^{3} \Delta_i = 0$, we set

$$\Lambda \Pi_h u_h = \bar{u}_{K_0} + \sum_{i=1}^{3} \Delta_i \phi_i(x, y)$$

2. If $\sum_{i=1}^{3} \Delta_i \neq 0$, we compute

$$\text{pos} = \sum_{i=1}^{3} \max(0, \Delta_i), \quad \text{neg} = \sum_{i=1}^{3} \max(0, -\Delta_i)$$

and define

$$\theta^+ = \min\left(1, \frac{\text{neg}}{\text{pos}}\right), \quad \theta^- = \min\left(1, \frac{\text{pos}}{\text{neg}}\right)$$

Finally, we set

$$\Lambda \Pi_h u_h = \bar{u}_{K_0} + \sum_{i=1}^{3} \widehat{\Delta}_i \phi_i(x, y)$$

where

$$\widehat{\Delta}_i = \theta^+ \max(0, \Delta_i) - \theta^- \max(0, -\Delta_i)$$

This limiting operator conserves the mass within each element and guarantees that the reconstructed gradient of $\Lambda \Pi_h u_h$ is not larger than that of $u_h$.

## 6. NUMERICAL EXPERIENCES

In order to test the behaviour of the introduced slope limiters, we present two classical numerical experiments for linear convection equations with either rectangular or triangular grids. All the presented numerical tests are solved by using the simplified Runge–Kutta method for the time integration.

### 6.1. Diagonally moving prism

The first test is a solid shifting of a square along the diagonal of a computational domain $\Omega = (0, 50) \times (0, 50)$. The initial profile is given by:

$$u_h(x, y, 0) = \begin{cases} 1 & (x, y) \in [1, 10] \times [1, 10] \\ 0 & \text{elsewhere} \end{cases}$$

A grid of $50 \times 50$ rectangular elements is considered to test the different reconstruction techniques. Periodic boundary conditions are applied. A parallel flow is taken diagonal to the grid such that $\beta = (1/2, 1/2)$. The time step is chosen to fix the condition $\mathscr{C} = \beta_x \Delta t / \Delta x = \beta y \Delta t / \Delta_y$. In Figures 11, 12 and 13, we present the DG results obtained by using Chavent–Jaffré, Cockburn–Shu and our modified slope limiters, respectively. The dashed lines represents
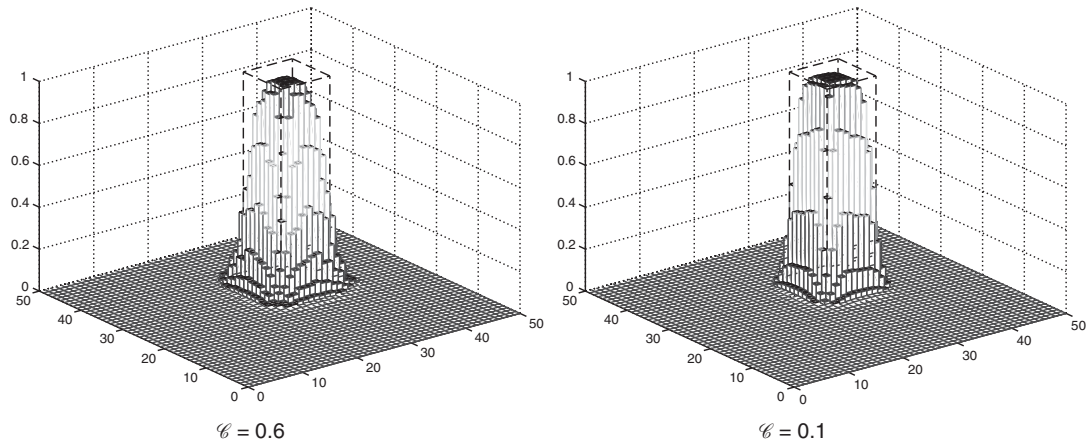
Figure 11. Results obtained by using the slope limiter introduced by Chavent and Jaffré with two different time steps.
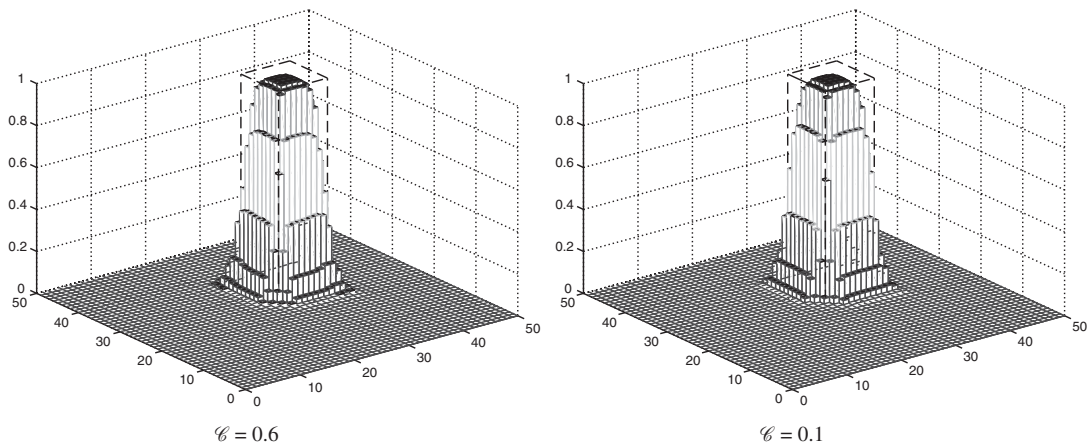


Figure 12. Results obtained by using the slope limiter introduced by Cockburn and Shu with two different time steps.

the shifted profile after a simulation time $T = 50$. It is clear that the slope limiter introduced by Chavent and Jaffré (Figure 11) produces some oscillations. However decreasing the time step leads to some smoothing in the solution. One the other hand, results obtained by Cockburn–Shu and our modified slope limiter seam to be very comparable. Table I presents the $L_1$ and $L_2$ errors against the resolution for the DG numerical solutions by using the three slope limiters. The modified slope limiter gives slightly better accuracy than the others.
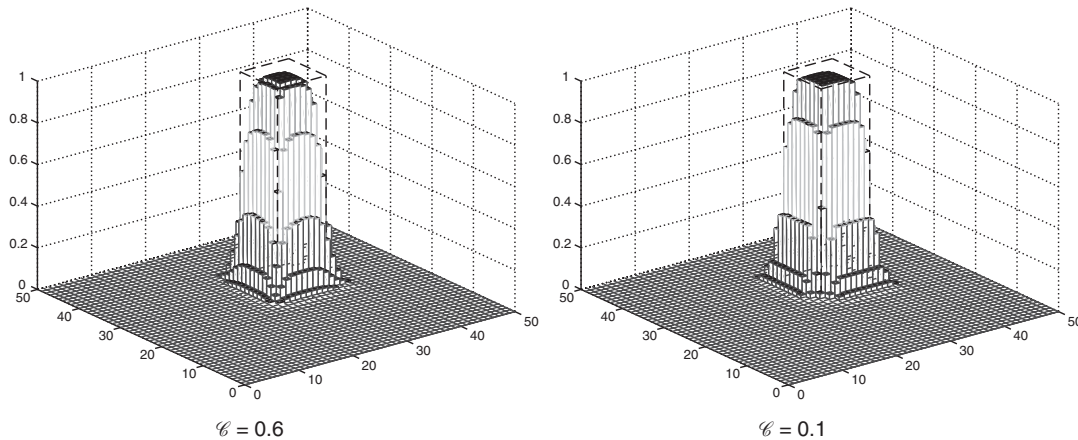
Figure 13. Results obtained by using the modified slope limiter with two different time steps.

Table I. $L_1$ and $L_2$ relative errors for different slope limiters.

| Slope limiter | $\mathscr{C} = 0.6$ $10^2$ error | | $\mathscr{C} = 0.1$ $10^2$ error | |
|---|---|---|---|---|
| | $L_1$—error | $L_2$—error | $L_1$—error | $L_2$—error |
| Chavent–Jaffré | 3.41 | 37.95 | 2.80 | 27.46 |
| Cockburn–Shu | 2.82 | 28.03 | 2.67 | 26.67 |
| Modified-limiter | 2.72 | 27.27 | 2.50 | 23.34 |

### 6.2. Rotating cylinder

A classical test for multi-dimensional scalar convection equation is the rotating cylinder (see, e.g. Reference [15]). A circular peak is moved in a rotating flow field. The results after four rotations are compared with the exact solution which is simply the initial conditions. Three grids made of rectangles, parallelograms and arbitrary triangles are used to examine the behaviour of the described slope limiters. The rotational velocity field $\beta(x) = r(x)(-\sin\theta, \cos\theta)$ is one rotation in $T = 2\pi$, where $r(x) = \|x - x_0\|$ is the rotation radius around the centre $x_0$. The time discretization step is taken to be 0.01; so that the Courant number does not exceed unity everywhere in the domain.

In the first test, the computational domain $\Omega = (0, 50) \times (0, 50)$ is discretized into a cartesian grid of $100 \times 100$ cells. Figures 14, 15 and 16 display the profile and the isolines of the DG solutions obtained by using Chavent–Jaffré, Cockburn–Shu and our modified slope limiters. The first limiter gives the least accurate solution. On the other hand, solutions obtained by the two other limiters seem to be very similar. Nevertheless, the $L_1$ and $L_2$ errors given in Figure 17 show a slight better accuracy of our modified limiter.

In the second test, the domain is discretized into a grid of parallelograms. In Figure 18, we present the results obtained after four rotations of the initial profile obtained by using
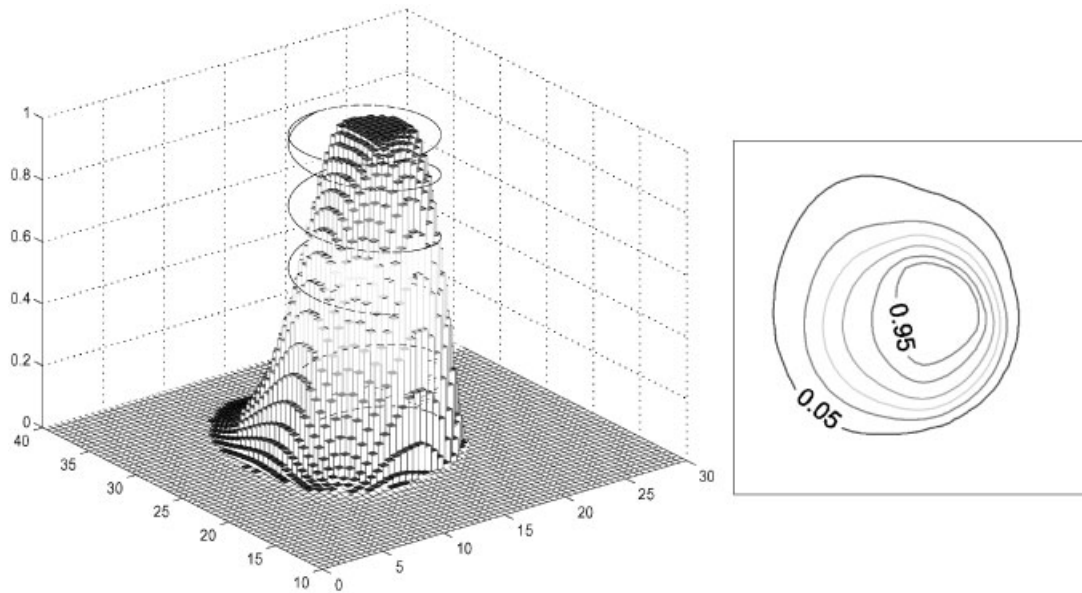
Figure 14. Profile and isolines of the DG solution obtained by using Chavent–Jaffré
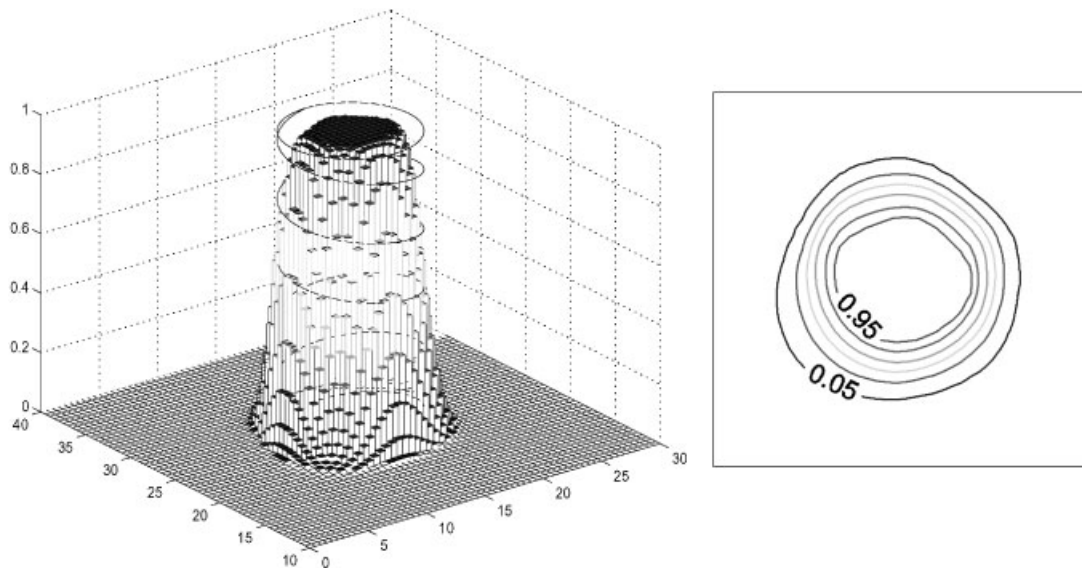slope limiter after four rotations.



Figure 15. Profile and isolines of the DG solution obtained by using Cockburn–Shu
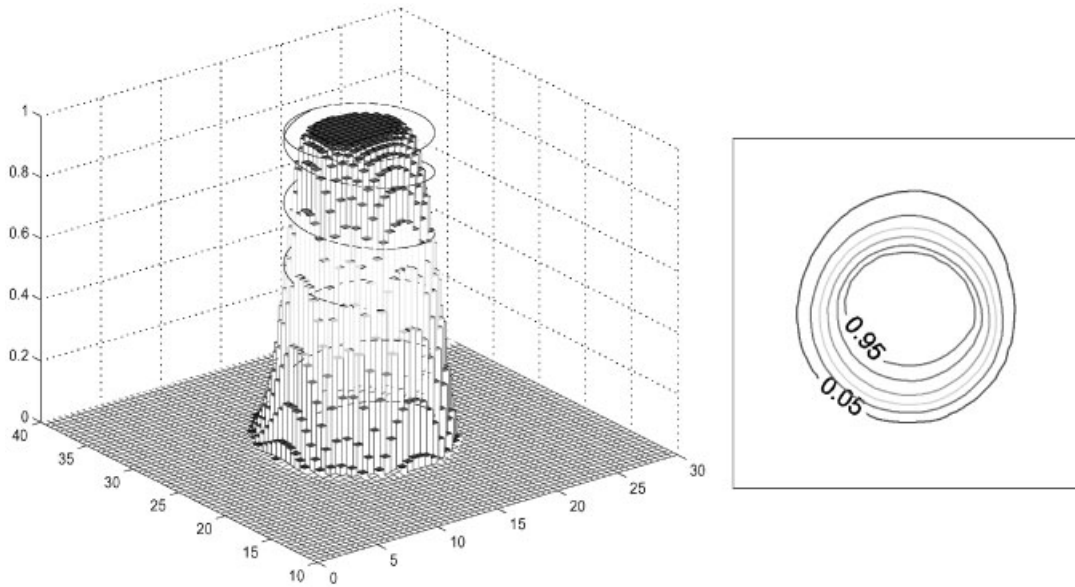slope limiter after four rotations.

Figure 16. Profile and isolines of the DG solution obtained by using the modified
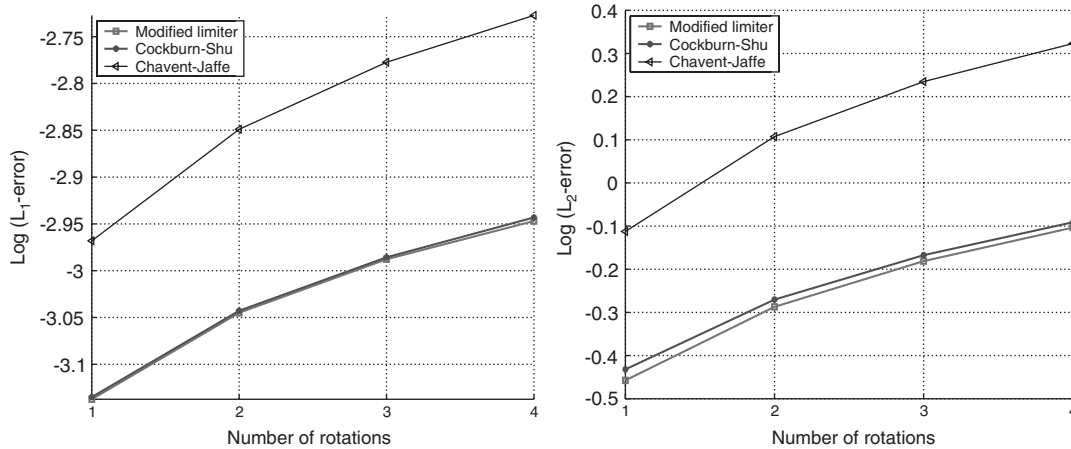slope limiter after four rotations.



Figure 17. $L_1$ and $L_2$ errors for the rotating cylinder.

Chavent–Jaffré and the modified slope limiters. The first limiter clearly suffers from dispersive
errors.

In the final test, computations are carried out on an arbitrary grid of triangular elements
made of 8000 cells. The results obtained by Chavent–Jaffré, Cockburn–Shu and the modified
slope limiters are depicted in Figures 19, 20 and 21, respectively. Errors presented in Figure 23
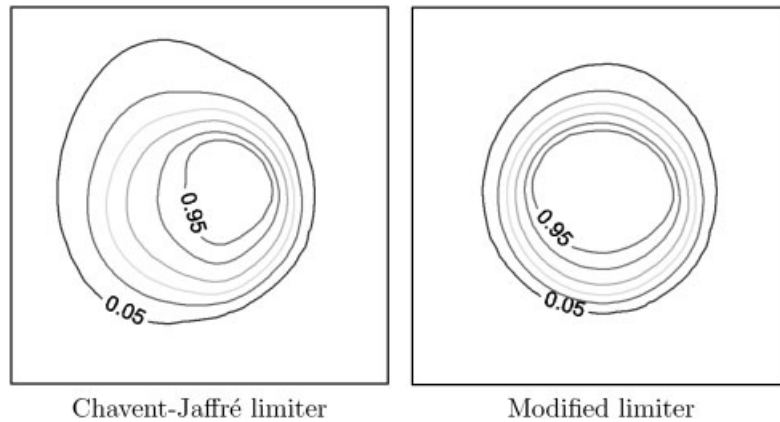
Figure 18. Results after four rotations obtained by using Chavent–Jaffré and the modified slope limiters on a grid made of parallelograms.
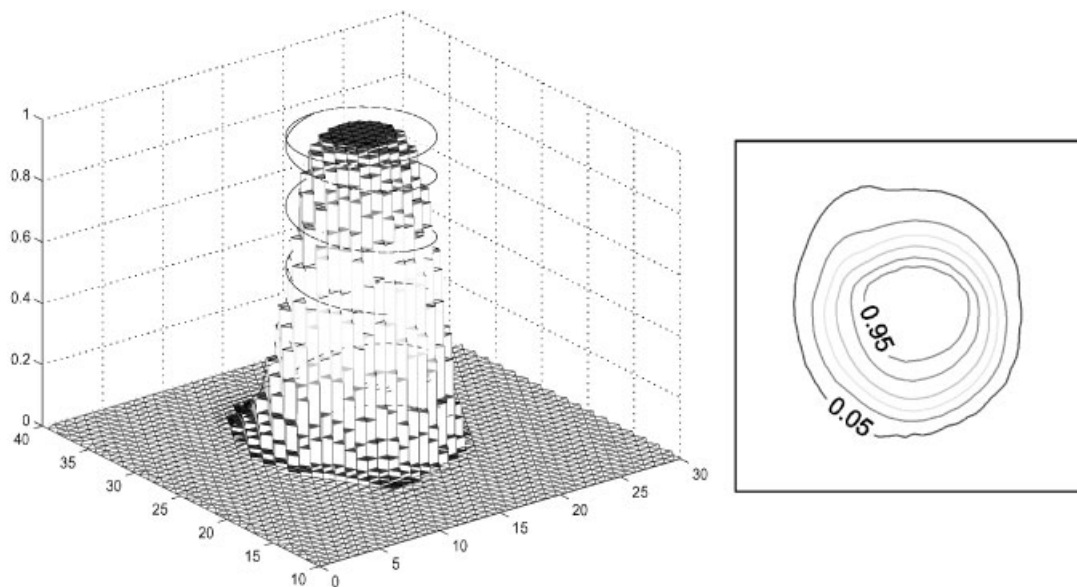


Figure 19. Profile and isolines of the solution obtained by using Chavent–Jaffré slope limiter over a triangular grid.

show that the limiter for triangular grids introduced by Cockburn and Shu is the least accurate. It should be noted that the degrees of freedom are chosen at the grid vertices. As we have previously mentioned, the DG method generates excessive smearing when degrees of freedom of the state function are sought at the midpoints of the grid edges. This is seen in Figure 22 for the DG approximation with the modified slope limiter. The two other slope limiters produce similar results.
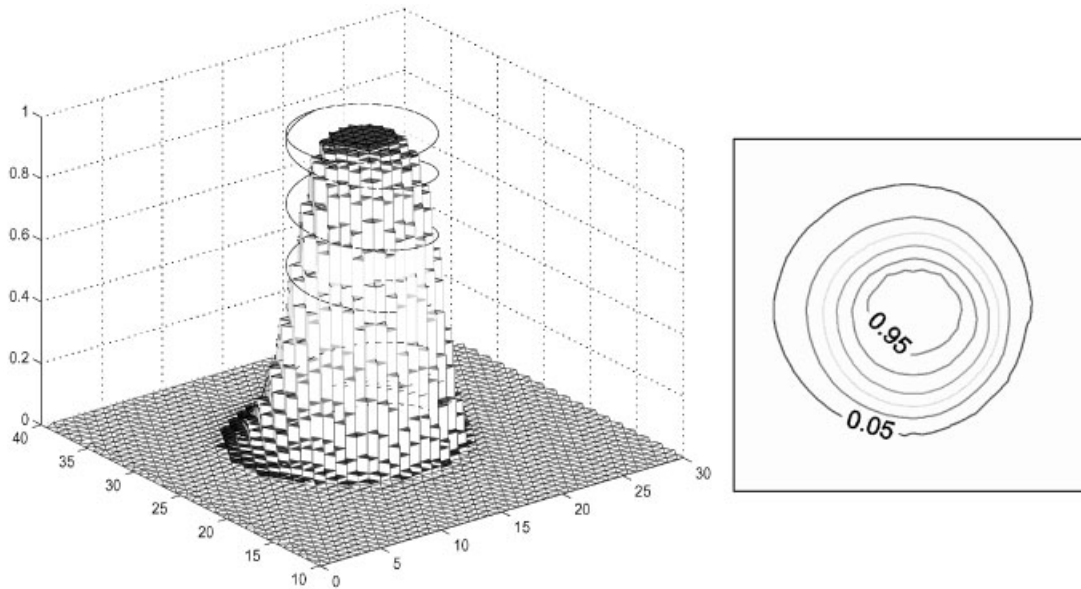
Figure 20. Profile and isolines of the solution obtained by using Cockburn–Shu slope limiter over a triangular grid.
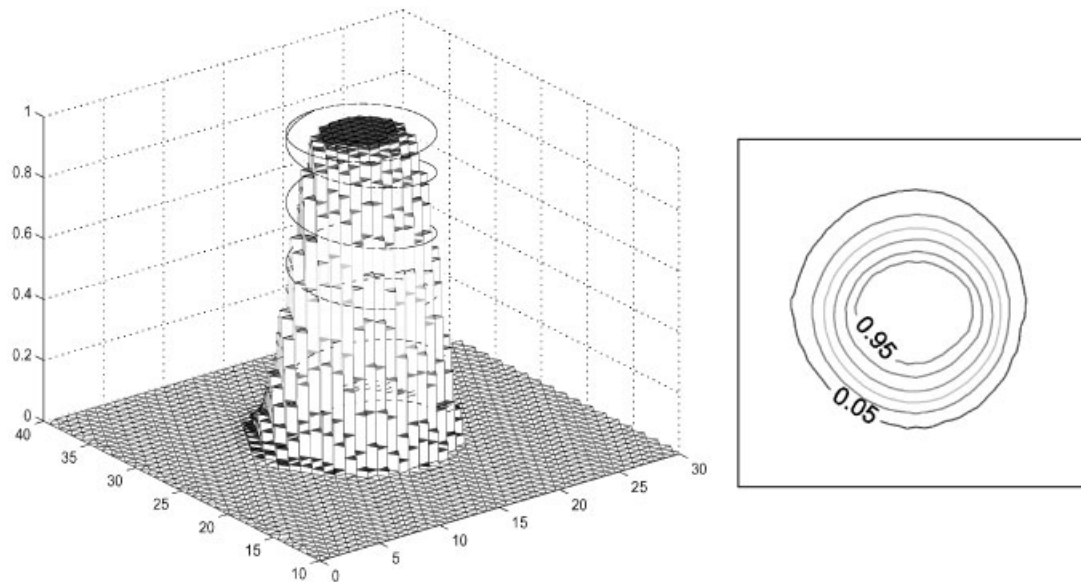


Figure 21. Profile and isolines of the solution obtained by the modified limiter over a triangular grid.
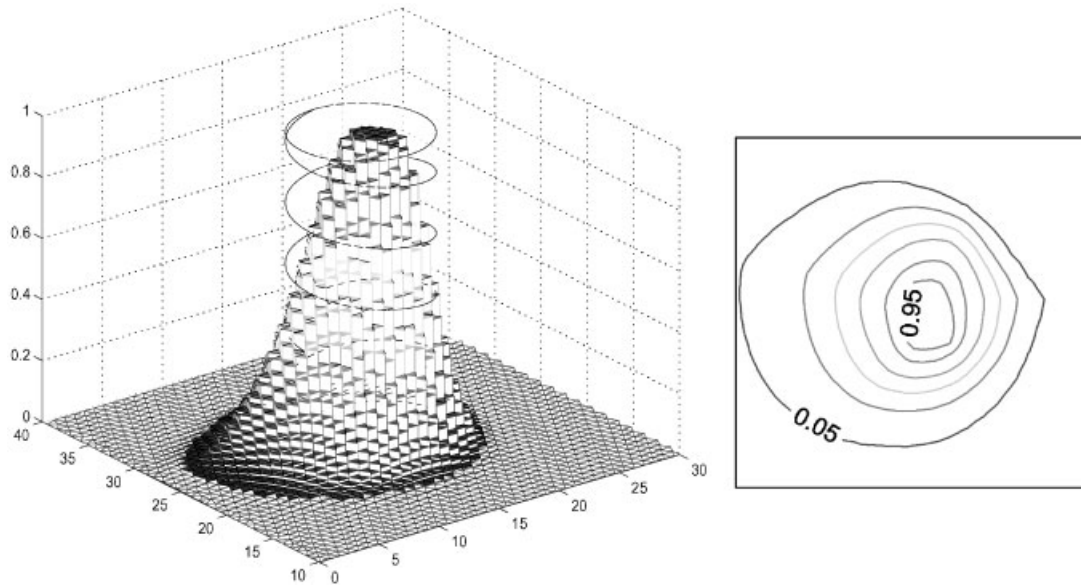
Figure 22. The approximation solution obtained by the DG method with degrees of freedom at the midpoints of the grid edges.



Figure 23. $L_1$ and $L_2$ errors for the rotating cylinder on the triangular mesh.

## 7. CONCLUSION

Data reconstruction is crucial for the stabilization of high-order discontinuous Galerkin methods. In one-dimensional space, many successful slope limiters have been developed. However, in higher dimensions, specially on unstructured grids, the construction of reliable slope limiters that preserve the accuracy of the scheme is still a challenge.

The multi-dimensional slope limiter introduced by Chavent and Jaffré which is an extension of the Van Leer's MUSCL slope limiter sometimes fails to eliminate all spurious oscillations for both rectangular and triangular discretizations. This drawback is due to the fact that the reconstruction of data by means of local constraints at the vertices is insufficient to prevent non-physical values at the midpoints of the edges. The proposed remedy is to reconstruct data by using constraints applied at the midpoints of the grid edges. This approach has an important physical interpretation since it limits the numerical fluxes across the interelements rather than the function values at the grid vertices.

For rectangular elements, where piecewise quadratic functions are used for the approximation space, the solution is reconstructed first at the midpoints of the cell edges by means of a dimension splitting technique. The nodal values are then reconstructed by solving a minimization problem. A similar approach is used for triangular grids. However, we have found that by taking the degrees of freedom of the approximation solution at the midpoints of the edges, the scheme becomes more diffusive.

Numerical comparisons with other slope limiters show a good improvement of the proposed reconstruction techniques.

## APPENDIX A

The minimization problem described in Section 3.2 may introduce some difficulties for the resolution. This problem is rewritten as follows:

For a given vector $\widetilde{U}_K \in \mathscr{P}_K$, find $U_K \in \mathscr{P}_K \cap \mathscr{Q}_K$ the solution of the least squares problem:

$$\min_{W \in \mathscr{P}_K \cap \mathscr{Q}_K} \mathscr{J}(W)$$

where $\mathscr{J}(W) = \frac{1}{2} \| W - \widetilde{U}_K \|_2$ is the objective function, $\mathscr{P}_K$ and $\mathscr{Q}_K$ are the hyperplane and the hypercube describing, respectively, the linear equality and inequality constraints as follows:

$$\mathscr{P}_K = \left\{ W \in \mathbb{R}^{\mathscr{N}_K}; \sum_{i=1}^{\mathscr{N}_K} w_i = \mathscr{N}_K \overline{w}_K \right\}$$

$$\mathscr{Q}_K = \prod_{i=1}^{\mathscr{N}_K} [\gamma_i, \mu_i]$$

with $\gamma_i = (1 - \alpha)\overline{u}_K + \alpha \overline{u}_{\min,i}$, $\mu_i = (1 - \alpha)\overline{u}_K + \alpha \overline{u}_{\max,i}$.

It is easy to check that the convex closed set $\mathscr{P}_K \cap \mathscr{Q}_K$ is non-empty since it contains the point $W = (w_i = \overline{u}_k, i = 1, \ldots, \mathscr{N}_K)$. Thus, the convex property of the objective function guarantees the existence and the uniqueness of the solution.

In this appendix, we present an efficient algorithm which is based on the so-called *active set algorithm* [26]. This algorithm is not iterative in nature, but rather it decreases the value of the objective function so that the optimal solution is reached in finitely many steps.

For any $W \in \mathbb{R}^{\mathscr{N}_K}$, $W$ is a *feasible point* if it satisfies all the inequality constraints, i.e. $W \in \mathscr{Q}_K$. Let us denote by $I = I(W)$ the set of indices of active constraints at $W$, i.e. constraints satisfied with equality.

In the active set algorithm a sequence of equality-constrained problems are solved corresponding to a prediction of the active set. At each step one constraint is added to or dropped

from the active set. The stepping process is described as follows:

1. Initialization:
   Choose $W^{(0)} = (\overline{w}_K, i = 1, \ldots, \mathcal{N}_K)$ a feasible point.
2. Stepping process:
   Let $W^{(k)}$ be the iterate at the $k$th step and $I^{(k)}$ the corresponding active set. The process seeks $(W^{(k+1)}, I^{(k+1)})$ according to the following steps:

   3. Solve the equality-constrained problem:

   $$\min_Z \mathscr{J}(Z) \quad \text{subject to}$$

   $$Z \in \mathscr{P}_K$$

   $$z_i = w_i^{(k)} \quad \text{for } i \in I^{(k)}$$

   This problem can be easily solved by using Lagrange multipliers $\lambda_i$.
   4. If $Z$ is a feasible point, then
      Check for optimality such that:

      $$\forall i \in I^{(k)}, \lambda_i \text{ is optimal} \; \Leftrightarrow \; \begin{cases} (w_i^{(k)} = \gamma_i \quad \text{and} \quad \lambda_i \leqslant 0) \\ \qquad\qquad\qquad \text{or} \\ (w_i^{(k)} = \mu_i \quad \text{and} \quad \lambda_i \geqslant 0) \end{cases}$$

      4.1. If $I^{(k)} = \emptyset$ or $\lambda_i$, $i \in I^{(k)}$, are optimal then
           the optimal solution $Z$ is reached.
      4.2. Otherwise, there exists $i \in I^{(k)}$ such that $\lambda_i$ is non-optimal, then one active constraint is dropped such that
           $I^{(k+1)} = I^{(k)} \backslash \{i\}$, where $|\lambda_i| = \max\{|\lambda_j|; \lambda_j\}$ is non-optimal.
           Set $W^{(k+1)} = Z$.
           Go to step 2.

   5. If $Z$ is not feasible then
      Choose $\delta = \min\{\delta_i; i \notin I^{(k)}\} \in [0, 1[$, such that

      $$\delta_i = \begin{cases} \dfrac{w_i^{(k)} - \gamma_i}{w_i^{(k)} - z_i} & \text{if } z_i < \gamma_i \\[3mm] \dfrac{w_i^{(k)} - \mu_i}{w_i^{(k)} - z_i} & \text{if } z_i < \mu_i \end{cases}$$

      Set $W^{(k+1)} = W^{(k)} + \delta(Z - W^{(k)})$.
      Update $I^{(k+1)}$ by checking the active constraints.
      Go to step 2.

This algorithm is not expensive from a computational point of view. Numerical observations showed that the optimal solution is reached with at most $2\mathcal{N}_K$ steps.

# REFERENCES

1. Godunov S. Finite difference methods for numerical computation of discontinuous solutions of the equations of fluid dynamics. *Mathematics Sbornik* 1959; **47**:271–306.
2. Van Leer B. Towards the ultimate conservative scheme: II. *Journal of Computational Physics* 1974; **14**: 361–376.
3. Van Leer B. Towards the ultimate conservative scheme: IV. A new approach to numerical convection. *Journal of Computational Physics* 1977; **23**:276–299.
4. Van Leer B. Towards the ultimate conservative scheme: V. A second order Godunov's method. *Journal of Computational Physics* 1979; **32**:101–136.
5. Toro E. *Riemann Solvers and Numerical Methods for Fluid Dynamics*. Springer: Berlin, 1997.
6. Hirsch C. *Numerical Computation of Internal and External Flows*. Wiley-Interscience Publication: New York, 1990.
7. Shu CW. TVB uniformly high order schemes for conservative laws. *Mathematics of Computation* 1987; **49**:105–121.
8. Cockburn B, Shu CW. TVB Runge Kutta local projection discontinuous Galerkin finite element method for conservative laws II: general frame-work. *Mathematics of Computation* 1989; **52**:411–435.
9. Cockburn B, Hou S, Shu CW. TVB Runge Kutta local projection discontinuous Galerkin finite element method for conservative laws III: one dimensional systems. *Journal of Computational Physics* 1989; **84**:90–113.
10. Goodman J, LeVeque R. On the accuracy of stable schemes for 2D conservation laws. *Mathematics of Computation* 1985; **45**:15–21.
11. Gowda V, Jaffré J. A discontinuous finite element method for scalar nonlinear conservation laws. *Rapport de recherche INRIA, No. 1848*, 1993.
12. Chavent G, Jaffré J. *Mathematical Models and Finite Elements for Reservoir Simulation*. Studies in Mathematics and its Applications. North-Holland: Amsterdam, 1986.
13. Gowda V. Discontinuous finite elements for nonlinear scalar conservation laws. *Thèse de Doctorat*, Université Paris IX, 1988.
14. Chavent G, Salzano. A finite-element method for the 1-D water flooding problem with gravity. *Journal of Computational Physics* 1982; **45**:307–344.
15. Johnson C. *Numerical Solution of Partial Differential Equations by the Finite Element Method*. Cambridge University Press: Cambridge, 1995.
16. Cockburn B, Shu CW. The Runge–Kutta discontinuous Galerkin method for conservative laws V: multidimentional systems. *Journal of Computational Physics* 1998; **141**:199–224.
17. Chavent G, Cockburn B. The local projection $P^0$ $P^1$-discontinuous Galerkin finite element method for scalar conservation laws. *Math. Modelling and Numerical Analysis* 1989; **23**:565–592.
18. Hubbard ME. Multidimensional slope limiters for MUSCL-type finite volume schemes on unstructured grids. *Journal of Computational Physics* 1999; **155**:54–74.
19. Kaddouri L. Une méthode d'élément finis discontinus pour les équations d'Euler des fluides compressibles. *Thèse de Doctorat*, Université Paris VI, 1988.
20. Siegel P, Mosé R, Ackerer Ph, Jaffré J. Solution of the advection dispersion equation using a combination of discontinuous and mixed finite elements. *Journal for Numerical Methods in Fluids* 1997; **24**:595–613.
21. Buès M, Oltean C. Numerical simulations for saltwater intrusion by mixed hybrid finite element method and discontinuous finite element method. *Transport in Porous Media* 2000; **40**:171–200.
22. Harten A. On a class of high resolution total-variation-stable finite-difference schemes. *SIAM Journal on Numerical Analysis* 1984; **21**:1–23.
23. Osher S. Convergence of generalized MUSCL schemes. *SIAM Journal on Numerical Analysis* 1993; **28**: 907–922.
24. Harten A. High resolution schemes for hyperbolic conservation laws. *Journal of Computational Physics* 1983; **49**:357–393.
25. Batten P, Lambert C, Causen D. Positively conservative high-resolution convection schemes for unstructured elements. *International Journal for Numerical Methods in Engineering* 1996; **39**:1821–1838.
26. Bjorck A. *Numerical Methods for Least Squares Problems*. SIAM: Philadelphia, 1996.