# A monotone nonlinear finite volume method for advection–diffusion equations on unstructured polyhedral meshes in 3D

K. NIKITIN[*] and  Yu. VASSILEVSKI[*]

**Abstract** — We present a new monotone finite volume method for the advection–diffusion equation with a full anisotropic discontinuous diffusion tensor and a discontinuous advection field on 3D conformal polyhedral meshes. The proposed method is based on a nonlinear flux approximation both for diffusive and advective fluxes and guarantees solution non-negativity. The approximation of the diffusive flux uses the nonlinear two-point stencil described in [9]. Approximation of the advective flux is based on the second-order upwind method with a specially designed minimal nonlinear correction [26]. The second-order convergence rate and monotonicity are verified with numerical experiments.

The discrete maximum principle (DMP) and local mass conservation are important properties of a numerical scheme for the approximate solution of the steady state advection–diffusion equation. An accurate discretization method satisfying DMP is hard to develop. We address the monotonicity condition as the simplified version of the DMP, which guarantees only solution non-negativity. A number of physical quantities (concentration, temperature, etc.) are non-negative by their nature and their approximations should be non-negative as well. We present a nonlinear finite volume (FV) method on conformal polyhedral meshes that satisfies the monotonicity condition for a wide range of problem coefficients. We admit a jumping diffusion coefficient represented by full anisotropic tensors, a jumping advection coefficient, which may be produced by the Darcy equation in multimaterial media, and both diffusion-dominated and advection-dominated regimes. The presented method is the extension of numerical schemes [9, 26] developed for the 3D diffusion equation [9] and the 2D advection–diffusion equation with continuous coefficients [26].

The major difficulty encountered in the design of a monotone numerical scheme is suppressing unwanted spurious (non-physical) oscillations in the numerical solution. These oscillations may appear in advection-dominated problems due to internal shocks and boundary layers, and in diffusion-dominated problems in highly anisotropic media due to inappropriate approximations of the diffusive flux.

In the finite element (FE) context, efficient damping of spurious oscillations in advection-dominated regimes has been developed within the streamline upwind

[*]Institute of Numerical Mathematics, Russian Academy of Sciences, Moscow 119333, Russia

Petrov–Galerkin (SUPG) method [5]. However, spurious oscillations around sharp layers may still appear in the SUPG solution. Spurious oscillations at layers diminishing (SOLD) methods [11] are generalizations of SUPG, which satisfy the DMP at least in some model cases. Spurious oscillations of FE solutions in diffusion-dominated regimes are caused by approximation difficulties in the case of general meshes and diffusion tensors. The theoretical analysis of DMP in the FE methods [8, 15, 32] imposes severe restrictions on the coefficients and the computational mesh. An algebraic flux correction [16, 17] is the alternative approach to the design of monotone FE methods. We note, however, that many FE methods are formally not locally conservative on the cells of the original computational mesh.

The finite volume (FV) methods, in contrast, guarantee the local mass conservation by their construction. The development of new FV methods for the advection–diffusion equation has been a popular topic of research, (see [3, 4, 10, 20, 28, 36]) for the steady equation and [35] for the unsteady equation and references therein). The advective fluxes can be approximated via the upwinding approach and controlled with different slope-limiting techniques [4, 7, 23] or the introduction of artificial viscosity [2, 28]. Many advanced second-order accurate *linear* methods for the diffusion equation fail to satisfy the monotonicity condition [1, 24, 30]. *Nonlinear* methods have seemed to be the feasible approach towards monotone and second-order accurate discretization [4, 11]. Nonlinear methods have been developed for the Poisson equation [6] and for the general diffusion equation [9, 14, 18, 21, 22, 24, 27, 29, 35, 37].

Our approximation of the advective flux is the 3D extension of the 2D nonlinear method proposed in [26]. The method follows the idea of the MUSCL method [34] and uses a piecewise linear discontinuous reconstruction of the FV solution on polyhedral cells, whose slope is limited via a three-by-three matrix with nonlinear entries. More precisely, we minimize the deviation of the reconstructed linear function from the given values at selected points subject to some monotonicity constraints, which form a convex set in the space of the function gradient components. The constraints are related to those considered in [12], but differ in the set of selected points and in the norm of deviation.

For the discretization of the diffusive flux we adopt the nonlinear two-point flux approximation on polyhedral meshes proposed in [9]. The original idea was proposed by Le Potier in [21] for the explicit scheme for the unsteady diffusion equation on triangular meshes. Further developments of the method [14, 24, 35, 37] extend it to a wider class of meshes and equations, but inherit the interpolation from *primary* unknowns defined at the mesh cells to *secondary* unknowns at the mesh vertices. The use of interpolation affects the accuracy of the numerical scheme, as well as the properties of nonlinear solvers. An interpolation-free nonlinear FV method on 2D meshes with polygonal cells was developed in [25]. It was extended to polyhedral meshes in [9] using *physical* interpolation for secondary unknowns at boundary faces and faces where the diffusion tensor jumps.

The proposed FV method is exact for linear solutions. Therefore, for problems with smooth solutions, one can expect the second-order asymptotic convergence

rate, which is confirmed in our numerical experiments. The monotone properties of the discrete solution are illustrated by the numerical experiments as well. The two-point stencil for flux approximation results in sparse matrices even on polyhedral meshes. For cubic meshes and a diagonal diffusion tensor these matrices reduce to the conventional 7-point discretization. Although the method is not interpolation-free, most of the interpolation operations are based on physical principles and thus do not affect the numerical properties of the method.

The major computational overhead in nonlinear FV methods is related to two nested iterations in the solution of a nonlinear algebraic problem. The outer iteration is the Picard method, which guarantees the solution non-negativity on each iteration. The set of admissible gradients in linear reconstruction used in the discretization of the advective fluxes is chosen to guarantee the stability of Picard iterations. The inner iteration is the Krylov subspace method for solving linearized problems.

The paper outline is as follows. In Section 1, we state the steady advection–diffusion problem. In Section 2, we describe the construction of discrete fluxes which form the basis of our method. In Section 3, we discuss the properties of the resulting algebraic system and present our algorithm for the generation and solution of that system. In Section 4, we present the numerical properties of the scheme using tetrahedral, hexahedral, and triangular prismatic meshes.

## 1. Steady-state advection–diffusion equation

Let $\Omega$ be a three-dimensional polyhedral domain with the boundary $\Gamma = \Gamma_N \cup \Gamma_D$ where $\Gamma_D \cap \Gamma_N = \varnothing$ and $\Gamma_D$ has a non-zero measure. We consider a model advection–diffusion problem for an unknown concentration $c$ [see 19, 31]:

$$\operatorname{div}\left(\mathbf{v}c - \mathbb{K}\nabla c\right) = g \quad \text{in } \Omega$$
$$c = g_D \quad \text{on } \Gamma_D \tag{1.1}$$
$$-\left(\mathbb{K}\nabla c\right)\cdot\mathbf{n} = g_N \quad \text{on } \Gamma_N$$

where $\mathbb{K}(\mathbf{x})$ is a symmetric positive definite possibly anisotropic piecewise continuous diffusion tensor, $\mathbf{v}(\mathbf{x})$ is a piecewise continuous velocity field, $g \in L^2(\Omega)$ is a source term, $\mathbf{n}$ is the exterior normal vector, and $g_D$, $g_N$ are given boundary data. It is known [19] that under the above assumptions and appropriate restrictions on $g_D$, $g_N$, equation (1.1) has a unique weak solution $c \in W_0^1(\Omega)$. We denote by $\Gamma_{\text{out}}$ the outflow part of $\Gamma$ where $\mathbf{v} \cdot \mathbf{n} \geqslant 0$, and define $\Gamma_{\text{in}} = \Gamma \setminus \Gamma_{\text{out}}$. We assume that $\Gamma_N \subset \Gamma_{\text{out}}$.

The sufficient conditions for the non-negativity of the solution $c(x)$ are $g(x) \geqslant 0$, $g_D \geqslant 0$, and $g_N \leqslant 0$. We assume that these conditions hold. From a physical viewpoint, the requirements $g(x) \geqslant 0$ and $g_N \leqslant 0$ mean that no mass can be taken out of the system.

The Dirichlet boundary condition on $\Gamma_{\text{out}}$ and the discontinuity in boundary data on $\Gamma_{\text{in}}$ may result in parabolic boundary layers. Exponential boundary layers may appear at the part of $\Gamma_{\text{out}}$ where $\mathbf{v} \cdot \mathbf{n} > 0$. An ideal discretization scheme must add

a numerical diffusion, which is small enough to avoid excessive smearing of the boundary layers, but sufficient to damp non-physical oscillations.

## 2. Monotone nonlinear FV scheme on polyhedral meshes

In this section, we derive a FV scheme with a nonlinear two-point flux approximation. Let $\mathbf{q} = -\mathbb{K}\nabla c + c\mathbf{v}$ denote the total flux, which satisfies the mass balance equation:

$$\operatorname{div} \mathbf{q} = g \quad \text{in } \Omega. \tag{2.1}$$

Let $\mathscr{T}$ be a conformal polyhedral mesh composed of $N_{\mathscr{T}}$ shape-regular cells with planar faces. We assume that each cell is a star-shaped 3D domain with respect to its barycenter, and each face is a star-shaped 2D domain with respect to the face's barycenter. We assume that $\mathscr{T}$ is face-connected, i.e. it cannot be split into two meshes having no common faces. We also assume that the tensor function $\mathbb{K}(\mathbf{x})$ and the velocity field $\mathbf{v}(\mathbf{x})$ vary slightly inside each cell and $\operatorname{div} \mathbf{v} \in \mathrm{L}^2(\Omega)$, $\operatorname{div} \mathbf{v} \geqslant 0$ for almost every $\mathbf{x} \in \Omega$; however $\mathbb{K}$ and $\mathbf{v}$ may jump across the mesh faces, as well as may change the direction (the principal directions for $\mathbb{K}$), although the normal component of $\mathbf{v}$ must be continuous on any mesh face. We denote by $\mathbf{n}_T$ the exterior unit normal vector to $\partial T$ and by $\mathbf{n}_f$ the normal vector to face $f$ fixed once and for all. On a boundary face, the vector $\mathbf{n}_f$ is exterior. We assume that $|\mathbf{n}_f| = |f|$ where $|f|$ denotes the area of face $f$.

Let $N_{\mathscr{B}}$ be the number of boundary faces. By $\mathscr{F}_I$, $\mathscr{F}_B$ we denote disjoint sets of interior and boundary faces. The subset $\mathscr{F}_J$ of $\mathscr{F}_I$ includes the faces where $\mathbb{K}(\mathbf{x})$ or $\mathbf{v}(\mathbf{x})$ jump. The set $\mathscr{F}_B$ is further split into subsets $\mathscr{F}_B^D$ and $\mathscr{F}_B^N$ where the Dirichlet and Neumann boundary conditions, respectively, are imposed. Alternatively, the set $\mathscr{F}_B$ is split into subsets $\mathscr{F}_B^{\mathrm{out}}$ and $\mathscr{F}_B^{\mathrm{in}}$ of faces belonging to $\Gamma_{\mathrm{out}}$ and $\Gamma_{\mathrm{in}}$, respectively. Finally, $\mathscr{F}_T$ and $\mathscr{E}_T$ denote the sets of the faces and edges of the polyhedron $T$ respectively, whereas $\mathscr{E}_f$ denotes the set of the edges of the face $f$.

Integrating equation (2.1) over a polyhedron $T$ and using Green's formula we get:

$$\sum_{f \in \partial T} \chi_{T,f} \mathbf{q}_f \cdot \mathbf{n}_f = \int_T f \, \mathrm{d}x, \qquad \mathbf{q}_f = \frac{1}{|f|} \int_f \mathbf{q} \, \mathrm{d}s \tag{2.2}$$

where $\mathbf{q}_f$ is the average flux density for the face $f$, and $\chi_{T,f}$ is either $1$ or $-1$, depending on the mutual orientation of the normal vectors $\mathbf{n}_f$ and $\mathbf{n}_T$.

For each cell $T$, we assign one degree of freedom, $C_T$, for the concentration $c$. Let $C$ be the vector of all discrete concentrations. If two cells $T_+$ and $T_-$ have a common face $f$ and $\mathbf{n}_f$ is exterior to $T_+$, the two-point flux approximation is as follows:

$$\mathbf{q}_f^h \cdot \mathbf{n}_f = M_f^+ C_{T_+} - M_f^- C_{T_-} \tag{2.3}$$

where $M_f^+$ and $M_f^-$ are some coefficients. In a linear FV method, these coefficients are equal and fixed. In the nonlinear FV method, they may be different and depend on the concentrations in the adjacent cells. On a face $f \in \Gamma_D$, the flux has a form

similar to (2.3) with an explicit value for one of the concentrations. For the Dirichlet boundary value problem, $\Gamma_D = \partial\Omega$, substituting (2.3) into (2.2), we obtain a system of $N_{\mathscr{T}}$ equations with $N_{\mathscr{T}}$ unknowns $C_T$.

Therefore, the cornerstone of the cell–centered FV method is the definition of discrete flux (2.3). Our method of the discrete flux definition generalizes the definition of the diffusive flux [9] to the case of an advective–diffusive flux on the basis of the 2D method [26]. In order to present our method, we introduce the notations borrowing them from [9].

For every cell $T$ in $\mathscr{T}$, we define the collocation point $\mathbf{x}_T$ at the barycenter of $T$. For every face $f \in \mathscr{F}_B \cup \mathscr{F}_I$, we denote the face barycenter by $\mathbf{x}_f$ and associate a collocation point with $\mathbf{x}_f$ for $f \in \mathscr{F}_B \cup \mathscr{F}_J$. We also define the collocation points at the centers $\mathbf{x}_e$ of the edges $e \in \mathscr{E}_f$, $f \in \mathscr{F}_B \cup \mathscr{F}_J$.

We shall refer to the collocation points on faces and edges as *auxiliary* collocation points. They are introduced for mathematical convenience and will not enter the final algebraic system, although may affect the system coefficients. In contrast, we shall refer to the other collocation points as *primary* collocation points, whose discrete concentrations form the unknown vector in the algebraic system.

For every cell $T$ we define a set $\Sigma_T$ of nearby collocation points as follows. First, we add to $\Sigma_T$ the collocation point $\mathbf{x}_T$. Then, for every face $f \in \mathscr{F}_T \setminus (\mathscr{F}_J \cup \mathscr{F}_B)$, we add the collocation point $\mathbf{x}_{T'_f}$, where $T'_f$ is the cell, other than $T$, that has the face $f$. Finally, for any other face $f \in \mathscr{F}_T \cap (\mathscr{F}_B \cup \mathscr{F}_J)$, we add the collocation point $\mathbf{x}_f$. Let $N(\Sigma_T)$ denote the number of elements in the set $\Sigma_T$.

Similarly, for every face $f \in \mathscr{F}_B \cup \mathscr{F}_J$ belonging to a cell $T$ we define a set $\Sigma_{f,T}$ of nearby collocation points. We initialize $\Sigma_{f,T} = \{\mathbf{x}_f, \mathbf{x}_T\}$ and add to $\Sigma_{f,T}$ the points from $\Sigma_T$, which are the barycenters of the cells or faces that have common points with $f$. The cardinality of $\Sigma_{f,T}$ is denoted by $N(\Sigma_{f,T})$.

We assume that for every cell–face pair $T \to f$, $T \in \mathscr{T}$, $f \in \mathscr{F}_T$, there exist three points $\mathbf{x}_{f,1}$, $\mathbf{x}_{f,2}$, and $\mathbf{x}_{f,3}$ in the set $\Sigma_T$ such that the following condition holds (see Fig. 1): The co-normal vector $\boldsymbol{\ell}_f = \mathbb{K}(\mathbf{x}_f)\mathbf{n}_f$ starting from $\mathbf{x}_T$ belongs to the trihedral angle formed by the vectors

$$\mathbf{t}_{f,1} = \mathbf{x}_{f,1} - \mathbf{x}_T, \quad \mathbf{t}_{f,2} = \mathbf{x}_{f,2} - \mathbf{x}_T, \quad \mathbf{t}_{f,3} = \mathbf{x}_{f,3} - \mathbf{x}_T \tag{2.4}$$

and

$$\frac{1}{|\boldsymbol{\ell}_f|}\boldsymbol{\ell}_f = \frac{\alpha_f}{|\mathbf{t}_{f,1}|}\mathbf{t}_{f,1} + \frac{\beta_f}{|\mathbf{t}_{f,2}|}\mathbf{t}_{f,2} + \frac{\gamma_f}{|\mathbf{t}_{f,3}|}\mathbf{t}_{f,3} \tag{2.5}$$

where $\alpha_f \geqslant 0$, $\beta_f \geqslant 0$, $\gamma_f \geqslant 0$.

The coefficients $\alpha_f$, $\beta_f$, $\gamma_f$ are computed as follows:

$$\alpha_f = \frac{D_{f,1}}{D_f}, \qquad \beta_f = \frac{D_{f,2}}{D_f}, \qquad \gamma_f = \frac{D_{f,3}}{D_f} \tag{2.6}$$
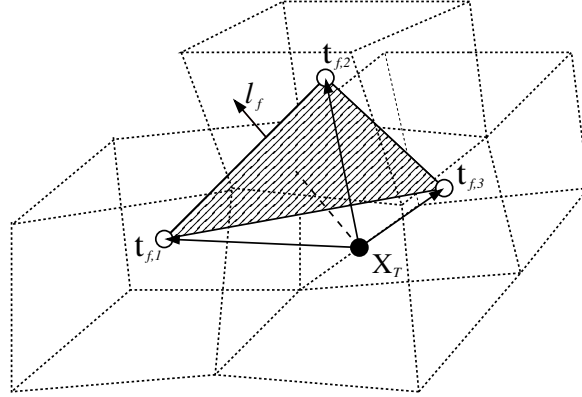
**Figure 1.** Co-normal vector and vector triplet.

where

$$D_f = \frac{|\mathbf{t}_{f,1} \quad \mathbf{t}_{f,2} \quad \mathbf{t}_{f,3}|}{|\mathbf{t}_{f,1}||\mathbf{t}_{f,2}||\mathbf{t}_{f,3}|}, \quad D_{f,1} = \frac{|\boldsymbol{\ell}_f \quad \mathbf{t}_{f,2} \quad \mathbf{t}_{f,3}|}{|\boldsymbol{\ell}_f||\mathbf{t}_{f,2}||\mathbf{t}_{f,3}|}$$

$$D_{f,2} = \frac{|\mathbf{t}_{f,1} \quad \boldsymbol{\ell}_f \quad \mathbf{t}_{f,3}|}{|\mathbf{t}_{f,1}||\boldsymbol{\ell}_f||\mathbf{t}_{f,3}|}, \quad D_{f,3} = \frac{|\mathbf{t}_{f,1} \quad \mathbf{t}_{f,2} \quad \boldsymbol{\ell}_f|}{|\mathbf{t}_{f,1}||\mathbf{t}_{f,2}||\boldsymbol{\ell}_f|}$$

and $|\mathbf{a} \ \mathbf{b} \ \mathbf{c}| = |(\mathbf{a} \times \mathbf{b}) \cdot \mathbf{c}|$.

Similarly, we assume that for every face–cell pair $f \to T$, $f \in \mathscr{F}_B \cup \mathscr{F}_J$, $T : f \in \mathscr{F}_T$ there exist three points $\mathbf{x}_{f,1}$, $\mathbf{x}_{f,2}$, and $\mathbf{x}_{f,3}$ in the set $\Sigma_{f,T}$ such that the vector $\boldsymbol{\ell}_{f,T} = -\mathbb{K}_T(\mathbf{x}_f)\mathbf{n}_f$ starting from $\mathbf{x}_f$ belongs to the trihedral angle formed by the vectors

$$\mathbf{t}_{f,1} = \mathbf{x}_{f,1} - \mathbf{x}_f, \quad \mathbf{t}_{f,2} = \mathbf{x}_{f,2} - \mathbf{x}_f, \quad \mathbf{t}_{f,3} = \mathbf{x}_{f,3} - \mathbf{x}_f \qquad (2.7)$$

and (2.5), (2.6) hold true.

A simple and efficient algorithm for searching triplets for the pairs $T \to f$ and $f \to T$ is presented in [9]. For the sake of brevity we omit the description of the algorithm and refer to [9].

The main idea of the proposed flux definition (2.3) is to define the diffusive and advective fluxes separately.

## 2.1. Nonlinear two-point diffusion flux approximation for an interior face

The definition of the diffusive flux is taken from [9] and is just outlined here.

Let $f$ be an interior face shared by the cells $T_+$ and $T_-$. We assume that $\mathbf{n}_f$ is outward for $T_+$ and $\mathbf{x}_\pm$ (or $\mathbf{x}_{T_\pm}$) is the collocation point of $T_\pm$ and $C_\pm$ (or $C_{T_\pm}$) is the discrete concentration in $T_\pm$.

We begin with the case $f \notin \mathscr{F}_J$ and introduce $\mathbb{K}_f = \mathbb{K}(\mathbf{x}_f)$. Let $T = T_+$. Using

the above notations, the definition of the directional derivative,

$$\frac{\partial c}{\partial \boldsymbol{\ell}_f} |\boldsymbol{\ell}_f| = \nabla c \cdot (\mathbb{K}_f \mathbf{n}_f)$$

and assumption (2.5), we write

$$\mathbf{q}_{f,d} \cdot \mathbf{n}_f = -\frac{|\boldsymbol{\ell}_f|}{|f|} \int_f \frac{\partial c}{\partial \boldsymbol{\ell}_f} \, ds = -\frac{|\boldsymbol{\ell}_f|}{|f|} \int_f \left( \alpha_f \frac{\partial c}{\partial \mathbf{t}_{f,1}} + \beta_f \frac{\partial c}{\partial \mathbf{t}_{f,2}} + \gamma_f \frac{\partial c}{\partial \mathbf{t}_{f,3}} \right) ds. \quad (2.8)$$

To derive a two-point flux approximation, we replace the partial derivatives with the finite differences and derive another approximation of the flux through the face $f$ using $T = T_-$. To distinguish between $T_+$ and $T_-$, we add subscripts $\pm$ and omit the subscript $f$. Since $\mathbf{n}_f$ is the internal normal vector for $T_-$, we have to change the sign of the right-hand side:

$$\mathbf{q}_{\pm,d}^h \cdot \mathbf{n}_f = \mp |\boldsymbol{\ell}_f| \left( \frac{\alpha_\pm}{|\mathbf{t}_{\pm,1}|} (C_{\pm,1} - C_\pm) + \frac{\beta_\pm}{|\mathbf{t}_{\pm,2}|} (C_{\pm,2} - C_\pm) + \frac{\gamma_\pm}{|\mathbf{t}_{\pm,3}|} (C_{\pm,3} - C_\pm) \right)$$
$$(2.9)$$

where $\alpha_\pm$, $\beta_\pm$ and $\gamma_\pm$ are given by (2.6) and $C_{\pm,i}$ denote the concentrations at the points $\mathbf{x}_{\pm,i}$ from $\Sigma_{T_\pm}$.

We define a new discrete flux as a linear combination of $\mathbf{q}_{\pm,d}^h \cdot \mathbf{n}_f$ with non-negative weights $\mu_\pm$:

$$\mathbf{q}_{f,d}^h \cdot \mathbf{n}_f = \mu_+ \mathbf{q}_{+,d}^h \cdot \mathbf{n}_f + \mu_- \mathbf{q}_{-,d}^h \cdot \mathbf{n}_f. \quad (2.10)$$

The weights are chosen so that $\mathbf{q}_{f,d}^h \cdot \mathbf{n}_f$ results in a two-point flux formula and $\mathbf{q}_{f,d}^h \cdot \mathbf{n}_f$ approximates the true diffusive flux. These requirements lead us to the following system

$$\begin{aligned} -\mu_+ d_+ + \mu_- d_- &= 0 \\ \mu_+ + \mu_- &= 1 \end{aligned} \quad (2.11)$$

where

$$d_\pm = |\boldsymbol{\ell}_f| \left( \frac{\alpha_\pm}{|\mathbf{t}_{\pm,1}|} C_{\pm,1} + \frac{\beta_\pm}{|\mathbf{t}_{\pm,2}|} C_{\pm,2} + \frac{\gamma_\pm}{|\mathbf{t}_{\pm,3}|} C_{\pm,3} \right). \quad (2.12)$$

Since coefficients $d_\pm$ depend on both geometry and concentration, the weights $\mu_\pm$ do as well. Thus, the resulting two-point flux approximation is *nonlinear*.

If the collocation point of $C_{+,i}$ ($C_{-,i}$), $i = 1, 2, 3$, coincides with the collocation point of $C_-$ ($C_+$), the terms in (2.12) are changed so that they do not incorporate $C_\pm$.

The solution of (2.11) can be written explicitly. In all cases $d_\pm \geqslant 0$ if $C \geqslant 0$. If $d_\pm = 0$, we set $\mu_+ = \mu_- = 1/2$. Otherwise,

$$\mu_+ = \frac{d_-}{d_- + d_+}, \qquad \mu_- = \frac{d_+}{d_- + d_+}.$$

Substituting this into (2.10), we get the two-point diffusive flux formula

$$\mathbf{q}^h_{f,d} \cdot \mathbf{n}_f = D^+_f C_{T_+} - D^-_f C_{T_-} \tag{2.13}$$

with non-negative coefficients

$$D^\pm_f = \mu_\pm |\boldsymbol{\ell}_f| (\alpha_\pm / |\mathbf{t}_{\pm,1}| + \beta_\pm / |\mathbf{t}_{\pm,2}| + \gamma_\pm / |\mathbf{t}_{\pm,3}|). \tag{2.14}$$

Now we consider the case $f \in \mathscr{F}_J$ when $\mathbb{K}_+(\mathbf{x}_f)$ and $\mathbb{K}_-(\mathbf{x}_f)$ differ, where

$$\mathbb{K}_\pm(\mathbf{x}_f) = \lim_{\mathbf{x} \in T_\pm, \; \mathbf{x} \to \mathbf{x}_f} \mathbb{K}(\mathbf{x}).$$

We derive two-point flux approximations in the cells $T_+$ and $T_-$ independently:

$$(\mathbf{q}^h_{f,d} \cdot \mathbf{n}_f)_+ = N^+ C_+ - N^+_f C_f \tag{2.15}$$

$$-(\mathbf{q}^h_{f,d} \cdot \mathbf{n}_f)_- = N^- C_- - N^-_f C_f. \tag{2.16}$$

Non-negative coefficients $N^+$, $N^+_f$, $N^-$, $N^-_f$ are derived similarly to coefficients (2.14) on the basis of discrete concentrations at collocation points from $\Sigma_{T_\pm}$, $\Sigma_{f,T_\pm}$ and $\boldsymbol{\ell}_\pm = \mp \mathbb{K}_\pm(\mathbf{x}_f) \mathbf{n}_f$, the co-normal vectors to the face $f$ outward with respect to $T_\pm$. The continuity of the normal component of the total flux and the advection field implies the continuity of the normal component of the diffusive flux. This assertion allows us to eliminate $C_f$ from (2.15), (2.16)

$$C_f = (N^+ C_+ + N^- C_-)/(N^+_f + N^-_f) \tag{2.17}$$

and derive the two-point flux approximation (2.3) with coefficients

$$D^\pm_f = N^\pm N^\mp_f / (N^+_f + N^-_f). \tag{2.18}$$

If both $N^\pm_f = 0$, we set $D^\pm_f = N^\pm/2$ and $C_f = (C_+ + C_-)/2$.

## 2.2. Nonlinear advection flux on interior faces

The method of the definition of the discrete advective flux is the generalization of the 2D method [26]. For any interior face $f \in \mathscr{F}_I$ the advective flux

$$\mathbf{q}_{f,a} = \frac{1}{|f|} \int_f c\mathbf{v} \, \mathrm{d}s$$

is approximated via an upwinded linear reconstruction $\mathscr{R}_T$ of the concentration over cell $T$

$$\mathbf{q}^h_{f,a} \cdot \mathbf{n}_f = v^+_f \mathscr{R}_{T_+}(\mathbf{x}_f) + v^-_f \mathscr{R}_{T_-}(\mathbf{x}_f) \tag{2.19}$$

where

$$v_f^+ = \frac{1}{2}(v_f + |v_f|), \qquad v_f^- = \frac{1}{2}(v_f - |v_f|), \qquad v_f = \frac{1}{|f|} \int_f \mathbf{v} \cdot \mathbf{n}_f \, ds.$$

We define the reconstruction $\mathscr{R}_T$ as a linear function

$$\mathscr{R}_T(\mathbf{x}) = C_T + \mathbf{g}_T \cdot (\mathbf{x} - \mathbf{x}_T) \qquad \forall \mathbf{x} \in T \tag{2.20}$$

with a gradient vector $\mathbf{g}_T$. Since $C_T$ is collocated at the barycenter of $T$, this reconstruction preserves the mean value of the concentration for any choice of $\mathbf{g}_T$.

The conventional reconstructions of the gradient target stable approximations of the second order. Let $\mathscr{G}_T$ be a set of admissible gradients $\tilde{\mathbf{g}}_T$, which will be defined below. The gradient vector $\mathbf{g}_T$ is the solution to the following constrained minimization problem:

$$\mathbf{g}_T = \arg \min_{\tilde{\mathbf{g}}_T \in \mathscr{G}_T} \mathscr{J}_T(\tilde{\mathbf{g}}_T) \tag{2.21}$$

where the functional

$$\mathscr{J}_T(\tilde{\mathbf{g}}_T) = \frac{1}{2} \sum_{\mathbf{x}_k \in \tilde{\Sigma}_T} [C_T + \tilde{\mathbf{g}}_T \cdot (\mathbf{x}_k - \mathbf{x}_T) - C_k]^2$$

measures the deviation of the reconstructed function from the targeted values $C_k$ collocated at points $\mathbf{x}_k$ from a set $\tilde{\Sigma}_T$. The set $\tilde{\Sigma}_T$ is built as follows. First, the auxiliary set $\hat{\Sigma}_T$ is defined by eliminating the secondary collocation points $\mathbf{x}_f$, $f \in \mathscr{F}_B^{\mathrm{out}}$, from the set $\Sigma_T$. Second, the set $\hat{\Sigma}_T$ is extended whenever it is either too small or ill-conditioned. More precisely, if $\hat{\Sigma}_T = \{\mathbf{x}_T, \mathbf{x}_{T'}\}$ or $\hat{\Sigma}_T = \{\mathbf{x}_T, \mathbf{x}_{T'}, \mathbf{x}_{T''}\}$, we add to it the elements of $\hat{\Sigma}_{T'}$ and $\hat{\Sigma}_{T''}$ other than $\mathbf{x}_T$. If $\hat{\Sigma}_T = \{\mathbf{x}_T, \mathbf{x}_{T'}, \mathbf{x}_{T''}, \mathbf{x}_{T'''}\}$ and the volume of the tetrahedron formed by these four points is less than $10^{-3}|T|$, we add to it the elements of $\hat{\Sigma}_{T'}$, $\hat{\Sigma}_{T''}$, and $\hat{\Sigma}_{T'''}$ other than $\mathbf{x}_T$. The resulting set forms the set $\tilde{\Sigma}_T$.

The set of admissible gradients $\mathscr{G}_T$ is defined via three constraints suppressing non-physical oscillations. These constraints (as well as the set $\tilde{\Sigma}_T$) have been designed to be practical and at the same time as weak as possible. First, a linear reconstruction defined by the admissible gradient $\tilde{\mathbf{g}}_T$ must be bounded at the collocation points $\mathbf{x}_k \in \hat{\Sigma}_T$:

$$\min \left\{ C_1, C_2, \ldots, C_{N(\hat{\Sigma}_T)} \right\} \leqslant C_T + \tilde{\mathbf{g}}_T \cdot (\mathbf{x}_k - \mathbf{x}_T) \leqslant \max \left\{ C_1, C_2, \ldots, C_{N(\hat{\Sigma}_T)} \right\}. \tag{2.22}$$

Due to (2.22), we get that $\tilde{\mathbf{g}}_T \equiv 0$ in local minima and maxima.

Second, for the sake of the correct sign of the advective flux, we require that the reconstructed function must be non-negative at points $\mathbf{x}_f$ on faces $f \in \mathscr{F}_T$ where $v_f > 0$:

$$C_T + \tilde{\mathbf{g}}_T \cdot (\mathbf{x}_f - \mathbf{x}_T) \geqslant 0. \tag{2.23}$$

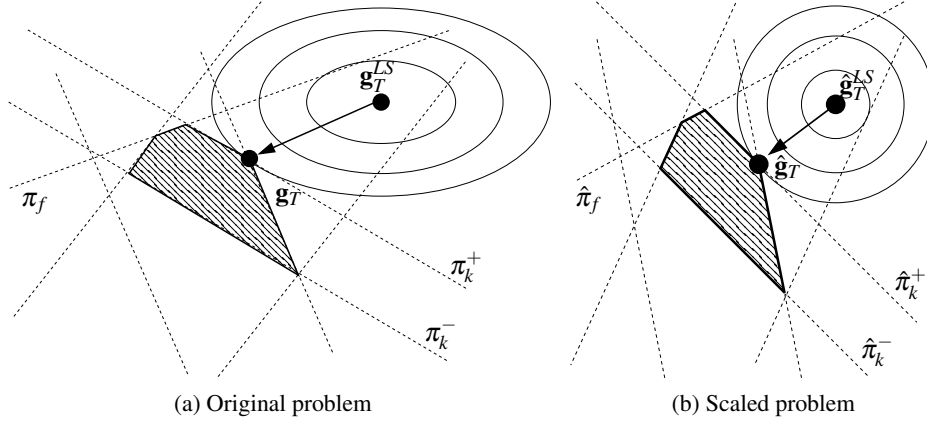(a) Original problem                                (b) Scaled problem

**Figure 2.** Original and scaled constrained problems.

We note that when the face center $\mathbf{x}_f$ lies outside the convex hull of points $\mathbf{x}_k \in \hat{\Sigma}_T$, the reconstructed function may be negative at $\mathbf{x}_f$ even if (2.22) is satisfied.

Third, the reconstructed function must be bounded from below at the secondary collocation points on $\Gamma_{\text{out}}$ (they do not belong to $\hat{\Sigma}_T$):

$$\min\left\{C_1, C_2, \ldots, C_{N(\hat{\Sigma}_T)}\right\} \leqslant C_T + \tilde{\mathbf{g}}_T \cdot (\mathbf{x}_f - \mathbf{x}_T), \qquad f \in \mathscr{F}_T \cap \mathscr{F}_B^{\text{out}}. \qquad (2.24)$$

Ignoring constraint (2.24) may produce instability in the iterative solution of the resulting nonlinear algebraic system.

It can be proved [26] that minimization problem (2.21) with constraints (2.22), (2.23), (2.24) has a unique solution.

Constraint (2.22) describes the slice between two planes in 3D-space:

$$\pi_k^- : C_T + \tilde{\mathbf{g}}_T \cdot (\mathbf{x}_k - \mathbf{x}_T) = \min\left\{C_1, C_2, \ldots, C_{N(\hat{\Sigma}_T)}\right\}$$

$$\pi_k^+ : C_T + \tilde{\mathbf{g}}_T \cdot (\mathbf{x}_k - \mathbf{x}_T) = \max\left\{C_1, C_2, \ldots, C_{N(\hat{\Sigma}_T)}\right\}.$$

Constraints (2.23) and (2.24) define half-spaces bounded by planes $\pi_f$:

$$\pi_f = \begin{cases} C_T + \tilde{\mathbf{g}}_T \cdot (\mathbf{x}_f - \mathbf{x}_T) = 0, & f \in \mathscr{F}_T, \quad v_f > 0 \\ C_T + \tilde{\mathbf{g}}_T \cdot (\mathbf{x}_f - \mathbf{x}_T) = \min\left\{C_1, \ldots, C_{N(\hat{\Sigma}_T)}\right\}, & f \in \mathscr{F}_T \cap \mathscr{F}_B^{\text{out}}. \end{cases}$$

Let $\mathscr{P}$ denote the set of all such planes $\{\pi_k^{\pm} : \mathbf{x}_k \in \hat{\Sigma}_T\}$ and $\mathscr{P}_f$ denote the set of planes $\pi_f$.

The deviation functional $\mathscr{J}_T$ has ellipsoidal isosurfaces in the general case. The scaling operator $\mathscr{S}_T$ transforms the ellipsoids into spheres (see 2D case in Fig. 2), so that minimization of the functional reduces to a simple projection. The same operator maps the planes $\pi \in \mathscr{P} \cup \mathscr{P}_f$ into planes $\hat{\pi}$, the solution of the non-constrained minimization problem (2.21) point $\mathbf{g}_T^{LS}$ into $\hat{\mathbf{g}}_T^{LS}$ and thus reduces problem (2.21) to its scaled counterpart.

Algorithm 2.1 uses the scaled problem for searching the solution $\mathbf{g}_T$ of the original constrained minimization problem (2.21).

**Algorithm 2.1.** Solution of constrained minimization problem (2.21).

Use the least squares method to find the solution $\mathbf{g}_T^{LS}$ of the non-constrained counterpart of (2.21);

**if** $\mathbf{g}_T^{LS}$ satisfies (2.22), (2.23) and (2.24) **then**
$\quad$ $\mathbf{g}_T = \mathbf{g}_T^{LS}$.
$\quad$ Exit;
**end if**
Set $\mathbf{g}_T = \{0,0,0\}$;
Apply the scaling operator $\mathscr{S}_T$ to transform ellipsoidal isosurfaces into spheres and define $\hat{\pi} = \mathscr{S}_T(\pi)$, $\hat{\mathbf{g}}_T^{LS} = \mathscr{S}_T(\mathbf{g}_T^{LS})$ and $\hat{\mathbf{g}}_T = \mathscr{S}_T(\mathbf{g}_T)$;
**for** $\pi \in \mathscr{P} \cup \mathscr{P}_f$ **do**
$\quad$ **if** $\begin{array}{l} \pi = \pi_i^{\pm} \in \mathscr{P} \text{ and } \mathbf{g}_T^{LS} \text{ satisfies (2.22) for } k = i, \text{ or} \\ \pi = \pi_{f'} \in \mathscr{P}_f \text{ and } \mathbf{g}_T^{LS} \text{ satisfies (2.23) or (2.24) for } f = f' \end{array}$ **then**
$\quad\quad$ Continue;
$\quad$ **else**
$\quad\quad$ Project $\hat{\mathbf{g}}_T^{LS}$ onto the plane $\hat{\pi}$ to the point $\hat{\mathbf{g}}_T'$;
$\quad\quad$ **if** $\mathscr{S}_T^{-1}(\hat{\mathbf{g}}_T')$ satisfies (2.22)–(2.24) and $|\hat{\mathbf{g}}_T^{LS} - \hat{\mathbf{g}}_T'| < |\hat{\mathbf{g}}_T^{LS} - \hat{\mathbf{g}}_T|$ **then**
$\quad\quad\quad$ $\hat{\mathbf{g}}_T = \hat{\mathbf{g}}_T'$;
$\quad\quad$ **end if**
$\quad$ **end if**
**end for**
**for** any pair $\pi, \pi' \in \mathscr{P} \cup \mathscr{P}_f$ **do**
$\quad$ Find intersection line $g = \pi \cap \pi'$;
$\quad$ Find the segment $s$ of $g$ where all constrains are satisfied.
$\quad$ **if** the segment $s$ is not empty **then**
$\quad\quad$ Project $\hat{\mathbf{g}}_T^{LS}$ onto the segment $\mathscr{S}_T(s)$ to get the point $\hat{\mathbf{g}}_T''$;
$\quad\quad$ **if** $|\hat{\mathbf{g}}_T^{LS} - \hat{\mathbf{g}}_T''| < |\hat{\mathbf{g}}_T^{LS} - \hat{\mathbf{g}}_T|$ **then**
$\quad\quad\quad$ $\hat{\mathbf{g}}_T = \hat{\mathbf{g}}_T''$;
$\quad\quad$ **end if**
$\quad$ **end if**
**end for**
Apply inverse mapping $\mathbf{g}_T = \mathscr{S}_T^{-1}(\hat{\mathbf{g}}_T)$.

Using (2.19) and (2.20), we represent the advective flux as the sum of a linear part (the first-order approximation) and a nonlinear part (the second-order correction):

$$\mathbf{q}_{f,a}^h \cdot \mathbf{n}_f = A_f^+ C_+ - A_f^- C_- \tag{2.25}$$

where

$$A_f^{\pm} = \pm v_f^{\pm}(1 + \mathbf{g}_{\pm} \cdot (\mathbf{x}_f - \mathbf{x}_{\pm})C_{\pm}^{-1}) \tag{2.26}$$

subscript $\pm$ stands for $T_{\pm}$ and $\mathbf{g}_{\pm} = \mathbf{g}_{T_{\pm}}$.

The coefficients $A_f^{\pm}$ are non-negative for positive concentrations. If $C_T = 0$ in a cell $T$, then $\mathbf{g}_T$ must be the zero vector and $A_f^{\pm} = \pm v_f^{\pm}$.

## 2.3. Fluxes on boundary faces

Consider a Neumann boundary face $f \in \mathscr{F}_B^N$ and a cell $T$ containing $f$. The diffusive flux through this face is

$$\mathbf{q}_{f,d}^h \cdot \mathbf{n}_f = \bar{g}_{N,f}|f| \tag{2.27}$$

where $\bar{g}_{N,f}$ is the mean value of $g_N$ on the face $f$. We can think about $f$ as a cell with zero volume neighbouring $T$. Replacing $C_+$ and $C_-$ by $C_T$ and $C_f$, respectively, we get from formula (2.25) the approximation of the advective flux:

$$\mathbf{q}_{f,a}^h \cdot \mathbf{n}_f = A_f^+ C_T. \tag{2.28}$$

Thus, the equation for the total flux is

$$(\mathbf{q}_{f,d}^h + \mathbf{q}_{f,a}^h) \cdot \mathbf{n}_f = \bar{g}_{N,f}|\mathbf{n}_f| + A_f^+ C_T, \quad f \in \mathscr{F}_B^N \tag{2.29}$$

where coefficient $A_f^+$ is non-negative for non-negative concentrations.

Consider a Dirichlet boundary face $f \in \mathscr{F}_B^D$ and the cell $T$ containing this face. For the face $f$ we define

$$C_f = \bar{g}_{D,f} = \frac{1}{|f|} \int_f g_D \, ds \tag{2.30}$$

and for every edge $e \in \mathscr{E}_f$ of the face $f$:

$$C_e = \bar{g}_{D,e} = \frac{1}{|e|} \int_e g_D \, dx. \tag{2.31}$$

The approximation of the diffusive flux is given by the formula

$$\mathbf{q}_{f,d}^h \cdot \mathbf{n}_f = D_f^+ C_T - D_f^- C_f \tag{2.32}$$

where coefficients $D_f^\pm$ are given by (2.14). The approximation of the advective flux depends on the velocity direction on the face $f$. If $f \in \mathscr{F}_B^{\text{out}}$, the approximation adopts formulas (2.28) and (2.26). If $f \in \mathscr{F}_B^{\text{in}}$, we use

$$\mathbf{q}_{f,a}^h \cdot \mathbf{n}_f = -A_f^- \tag{2.33}$$

where

$$A_f^- = -\bar{g}_{D,f} v_f \equiv -\bar{g}_{D,f} v_f^- \geqslant 0. \tag{2.34}$$

## 2.4. Recovery of discrete solution at auxiliary collocation points

Coefficients $D_f^\pm$ in (2.14), (2.18) may depend on the discrete solution $C_f$ and $C_e$ at auxiliary collocation points $\mathbf{x}_f$, $f \in \mathscr{F}_B \cup \mathscr{F}_J$, and $\mathbf{x}_e$, $e \in \mathscr{E}_f$. On the other hand, the discrete FV system is formulated only for concentrations $C_T$ at the primary collocation points. The values $C_f$, $C_e$, $f \in \mathscr{F}_B^D$, $e \in \mathscr{E}_f$, are computed by (2.30),

(2.31) from the Dirichlet data. The values $C_f$, $f \in \mathscr{F}_J$, are recovered from (2.17). However, the values $C_f$, $C_e$, $f \in \mathscr{F}_B^N$, $e \in \mathscr{E}_f$, $e \notin \Gamma_D$ and the values $C_e$, $e \in \mathscr{E}_f$, $e \notin \Gamma_D$, $f \in \mathscr{F}_J$ have to be recovered from available data.

We recover the concentrations at Neumann faces from $C_T$ using (2.32) and (2.27). The coefficients $D_f^{\pm}$ can depend on the values at primary collocation points and $C_f$, $f \in \mathscr{F}_J \cup \mathscr{F}_B$, and $C_e$, $e \in \mathscr{E}_f$. Therefore, concentrations $C_f$ at mesh faces $f$, $f \in \mathscr{F}_J \cup \mathscr{F}_B$, are interpolated from the cell data on the basis of *physical* relationships, such as the diffusion flux continuity or the given diffusion flux. The coefficients of interpolation can depend on concentrations $C_e$ to be found at $\mathbf{x}_e$, $e \in \mathscr{E}_f$, $f \in \mathscr{F}_J \cup \mathscr{F}_B^N$, $e \notin \Gamma_D$. For every such edge, we suggest to compute $C_e$ by arithmetic averaging of $C_f$ for all faces $f \in \mathscr{F}_B^N \cup \mathscr{F}_J$ sharing $e$. These are the only rare data whose recovery is based on mathematical rather than physical arguments.

## 3. Discrete system and properties of discrete solution

Substituting two-point flux formula (2.3) with non-negative coefficients $M_f^{\pm} = D_f^{\pm} + A_f^{\pm}$ given by (2.14), (2.18) and (2.26) into the mass balance equation (2.2) and eliminating the boundary concentrations, we get a nonlinear system of $N_{\mathscr{T}}$ equations with $N_{\mathscr{T}}$ unknowns:

$$\mathbf{M}(\mathbf{C})\mathbf{C} = \mathbf{G}(\mathbf{C}) \tag{3.1}$$

where $\mathbf{C}$ is the vector of discrete concentrations at the primary collocation points. The matrix $\mathbf{M}(\mathbf{C})$ is assembled from $2 \times 2$ matrices

$$\mathbf{M}_f(\mathbf{C}) = \begin{pmatrix} M_f^+(\mathbf{C}) & -M_f^-(\mathbf{C}) \\ -M_f^+(\mathbf{C}) & M_f^-(\mathbf{C}) \end{pmatrix} \tag{3.2}$$

for the interior faces and $1 \times 1$ matrices $\mathbf{M}_f(\mathbf{C}) = M_f^+(\mathbf{C})$ for the Dirichlet faces. The right-hand side vector $\mathbf{G}(\mathbf{C})$ is generated by the source and the boundary data:

$$\mathbf{G}_T(\mathbf{C}) = \int_T g \, dx + \sum_{f \in \mathscr{F}_B^D \cap \partial T} M_f^-(\mathbf{C}) \bar{g}_{D,f} - \sum_{f \in \mathscr{F}_B^N \cap \partial T} |f| \bar{g}_{N,f} \quad \forall T \in \mathscr{T}. \tag{3.3}$$

For $g(x) \geqslant 0$, $g_D \geqslant 0$, and $g_N \leqslant 0$ the components of vector $\mathbf{G}$ are non-negative. We use the Picard iterations to solve nonlinear system (3.1). Our method of the generation and the solution of the algebraic system is summarized in Algorithm 3.1.

**Algorithm 3.1.** Generation and solution of nonlinear system (3.1).

For each cell–face pair $T \to f$, $f \in \mathscr{F}_T$, and each face–cell pair $f \to T$, $f \in \mathscr{F}_J \cup \mathscr{F}_B$ find vectors $\mathbf{t}_{f,1}$, $\mathbf{t}_{f,2}$, $\mathbf{t}_{f,3}$, satisfying conditions (2.4), (2.5) and (2.7), (2.5), respectively;

Select initial vectors $\mathbf{C}^0 \in \mathfrak{R}^{N_{\mathscr{T}}}$ and $\mathbf{C}_f^0 \in \mathfrak{R}^{N_{\mathscr{F}_J} + N_{\mathscr{F}_B^N}}$ with non-negative entries and a small value $\varepsilon_{\text{non}} > 0$;

Calculate concentrations $\mathbf{C}_e^0$ at the auxiliary collocation points on edges using (2.31) or arithmetic averaging of neighbouring data $\mathbf{C}_f^0$;

**for** $k = 0, \ldots,$ **do**

    Assemble the global matrix $\mathbf{M}_k = \mathbf{M}(\mathbf{C}^k)$ from the face-based matrices $\mathbf{M}_f(\mathbf{C}^k)$. To form $\mathbf{M}_f(\mathbf{C}^k)$, use (2.14), (2.18), (2.26) and concentrations $\mathbf{C}_f^k$, $\mathbf{C}_e^k$ at auxiliary collocation points;

    Calculate the right-hand side vector $\mathbf{G}^k = \mathbf{G}(\mathbf{C}^k)$ using (3.3) and concentrations $\mathbf{C}_f^k$, $\mathbf{C}_e^k$ at auxiliary collocation points;

    Stop if $\|\mathbf{M}_k\mathbf{C}^k - \mathbf{G}^k\| \leqslant \varepsilon_{\mathrm{non}}\|\mathbf{M}_0\mathbf{C}^0 - \mathbf{G}^0\|$.
    Solve $\mathbf{M}_k\mathbf{C}^{k+1} = \mathbf{G}^k$;

    Calculate concentrations $\mathbf{C}_f^{k+1}$ at the auxiliary collocation points on faces $f \in \mathscr{F}_J \cup \mathscr{F}_B$ using (2.17), (2.27), (2.30), (2.32), and data $\mathbf{C}^{k+1}$, $\mathbf{C}_f^k$, $\mathbf{C}_e^k$;

    Calculate concentrations $\mathbf{C}_e^{k+1}$ at the auxiliary collocation points on edges using (2.31) or arithmetic averaging of neighbouring data $\mathbf{C}_f^{k+1}$.

**end for**

The linear system in Step 8 with the non-symmetric matrix $\mathbf{M}(\mathbf{C}^k)$ can be solved, for example, by the preconditioned Bi-Conjugate Gradient Stabilized (BiCGStab) method [33]. The BiCGStab iterations are terminated when the relative norm of the residual becomes smaller than $\varepsilon_{\mathrm{lin}}$.

We note that since our method is exact for linear functions, we can expect the asymptotic second-order convergence rate on all sequences of meshes.

The next two theorems show that the solution to (3.1) is non-negative, provided that it exists, and that the nonlinear iterates $C^k$ are non-negative vectors, provided that $\varepsilon_{\mathrm{lin}} = 0$. The proofs of the theorems can be found in [26].

**Theorem 3.1.** *Let* $\Gamma_N = \varnothing$, $\mathscr{F}_B^D \equiv \mathscr{F}_B$, $g \geqslant 0$, $\mathrm{div}\,\mathbf{v} \geqslant 0$ *in* $\Omega$, $g_D \geqslant 0$ *on* $\Gamma_D \equiv \partial\Omega$ *and the solution* $\mathbf{C}$ *to* (3.1) *exist. Then* $\mathbf{C} \geqslant 0$.

**Theorem 3.2.** *Let* $g \geqslant 0$, $g_D \geqslant 0$, $g_N \leqslant 0$ *and* $\Gamma_D \neq \varnothing$ *in* (1.1). *If* $\mathbf{C}^0 \geqslant 0$ *and linear systems in the Picard method be solved exactly, then* $\mathbf{C}^k \geqslant 0$ *for* $k \geqslant 1$.

**Remark 3.1.** Theorems 3.1 and 3.2 hold true also for linear advective fluxes:

$$\mathbf{q}_{f,a}^h \cdot \mathbf{n}_f = A_f^+ C_+ - A_f^- C_-, \qquad A_f^\pm = \pm v_f^\pm.$$

The above assertions allow us to refer to the presented method as monotone, although it may violate the discrete maximum principle, see Subsection 4.2.

## 4. Numerical experiments

In all but one experiments, we set $\Gamma_N = \varnothing$. For advection-dominated problems, this helps to find more analytical solutions, such that the right-hand side vector is non-negative, $\mathbf{G}(\mathbf{C}) \geqslant 0$, for any non-negative vector $\mathbf{C}$.

We use the following discrete $L^2$-norms to evaluate relative discretization errors for the concentration $c$ and the flux $\mathbf{q}$:

$$
\varepsilon_2^c = \left[ \frac{\sum_{T \in \mathscr{T}} \left( c(\mathbf{x}_T) - C_T \right)^2 |T|}{\sum_{T \in \mathscr{T}} \left( c(\mathbf{x}_T) \right)^2 |T|} \right]^{1/2}, \qquad \varepsilon_2^q = \left[ \frac{\sum_{f \in \mathscr{F}_I \cup \mathscr{F}_B} \left( (\mathbf{q}_f - \mathbf{q}_f^h) \cdot \mathbf{n}_f \right)^2 |V_f|}{\sum_{f \in \mathscr{F}_I \cup \mathscr{F}_B} \left( \mathbf{q}_f \cdot \mathbf{n}_f \right)^2 |V_f|} \right]^{1/2}
$$

where $|V_f|$ is the arithmetic average of the volumes of mesh cells sharing the face $f$. The nonlinear iterations are terminated when the reduction of the initial residual norm becomes smaller then $\varepsilon_{\text{non}} = 10^{-7}$. The convergence tolerance for the linear solver is set to $\varepsilon_{\text{lin}} = 10^{-12}$. The linear regression algorithm has been used for calculating the convergence rates.

We consider three classes of polyhedral meshes for the unit cube $[0,1]^3$ introduced in [9]. All meshes are considered to be quasiuniform.

Hexahedral meshes are constructed from uniform cubic meshes by the distortion of the internal nodes. In each plane $x = 0.5$, $y = 0.5$, and $z = 0.5$ the nodes are randomly shifted along the planes. The position of other nodes is determined by the requirement of the face planarity. The distance and direction in which the nodes are shifted from the original position, are chosen randomly. The shifts of all nodes do not exceed $0.3h$, where $h$ is the cubic mesh size.

Prismatic meshes are constructed as a tensor product of a quasiuniform unstructured triangular $xy$-mesh and 1D $z$-mesh, both meshes having the size $h$. Additionally, $z$-planes are slightly tilted in such a way, that they do not intersect each other and the distance between them is at least $0.75h$. The height of each cell in these meshes lies between $0.75h$ and $1.25h$.

Tetrahedral meshes are quasiuniform unstructured tetrahedral meshes with the mesh size $h$. There is no hierarchical relation between a coarser and a finer meshes.

Representative examples of all three mesh classes are shown in Fig. 3.

### 4.1. Convergence study

At the first stage, the convergence study is performed for a smooth solution on tetrahedral, prismatic and hexahedral mesh sequences. We recall that we consider sequences of distorted meshes and thus perform the most challenging test for a numerical scheme. Let the exact solution, the velocity field and the diffusion tensor be as follows:

$$
c(x,y,z) = x \cos\left( \frac{\pi y}{2} \right) + \frac{\pi y}{2}, \qquad \mathbf{v} = (1,1,1)^T, \qquad \mathbb{K} = \begin{pmatrix} K & 0 & 0 \\ 0 & K & 0 \\ 0 & 0 & K \end{pmatrix}.
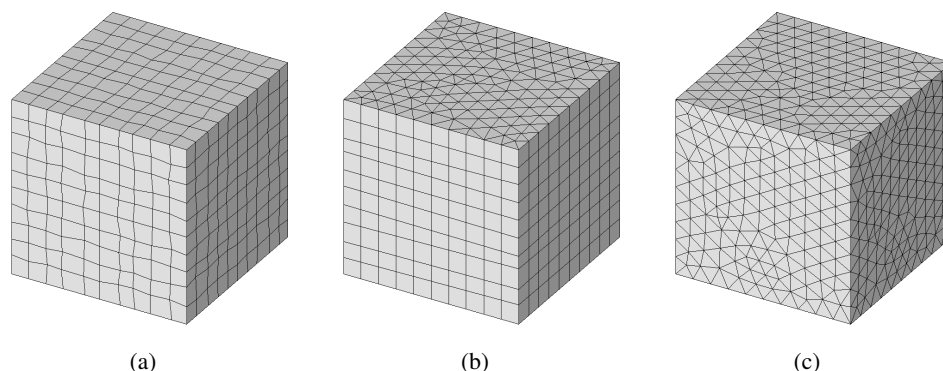$$

(a)                                    (b)                                    (c)

**Figure 3.** Examples of hexahedral (a), triangular prismatic (b), and tetrahedral (c) meshes.

**Table 1.**

Convergence analysis for diffusion-dominated problems ($K = 1$).

| $h$ | Hexahedral | | Prismatic | | Tetrahedral | |
|---|---|---|---|---|---|---|
| | $\varepsilon_2^C$ | $\varepsilon_2^q$ | $\varepsilon_2^C$ | $\varepsilon_2^q$ | $\varepsilon_2^C$ | $\varepsilon_2^q$ |
| 1/10 | 4.54e-4 | 3.64e-3 | 1.68e-4 | 2.01e-3 | 3.69e-4 | 8.59e-3 |
| 1/20 | 1.13e-4 | 1.24e-3 | 4.17e-5 | 8.09e-4 | 1.04e-4 | 3.88e-3 |
| 1/40 | 2.86e-5 | 4.15e-4 | 1.02e-5 | 2.51e-4 | 2.48e-5 | 1.89e-3 |
| rate | 1.99 | 1.57 | 2.02 | 1.50 | 1.95 | 1.09 |

**Table 2.**

Convergence analysis for advection-dominated problems ($K = 0.01$).

| $h$ | Hexahedral | | Prismatic | | Tetrahedral | |
|---|---|---|---|---|---|---|
| | $\varepsilon_2^C$ | $\varepsilon_2^q$ | $\varepsilon_2^C$ | $\varepsilon_2^q$ | $\varepsilon_2^C$ | $\varepsilon_2^q$ |
| 1/10 | 8.14e-4 | 8.53e-4 | 7.76e-4 | 4.80e-4 | 2.04e-3 | 1.88e-3 |
| 1/20 | 1.90e-4 | 2.08e-4 | 1.24e-4 | 9.90e-5 | 3.65e-4 | 3.38e-4 |
| 1/40 | 4.50e-5 | 4.97e-5 | 2.00e-5 | 2.38e-5 | 7.02e-5 | 7.51e-5 |
| rate | 2.09 | 2.05 | 2.64 | 2.17 | 2.43 | 2.32 |

The forcing term $f$ and the Dirichlet boundary data $g_D$ are set according to the exact solution. Table 1 shows the relative $L^2$-norms of the errors for a diffusion-dominated problem ($K = 1$) and Table 2 shows the relative $L^2$-norms of the errors for an advection-dominated problem ($K = 0.01$).

The convergence rate for the concentration is close to the second order, while the convergence rate for the flux is higher than the first order.

At the second stage, we consider the convergence towards the solution of the problem with a jumping diffusion tensor and a jumping velocity field. Let $\Omega = (0,1)^3$ be split into two non-overlapping subdomains $\Omega^{(1)} = \Omega \cap \{x < 0.5\}$, $\Omega^{(2)} = \Omega \cap \{x > 0.5\}$, with the interface defined by the plane $x = 0.5$, the tensor $\mathbb{K}$ and the
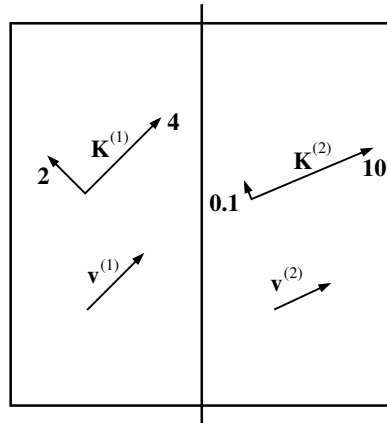
**Figure 4.** Tensor $\mathbb{K}$ and velocity $\mathbf{v}$ jumping across the interface $x = 0.5$.
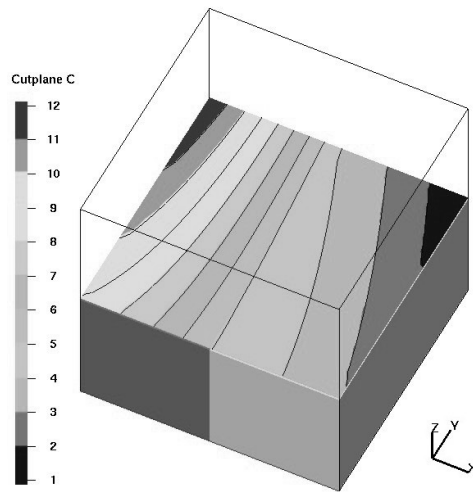


**Figure 5.** The solution isolines in the *xy*-plane for the problem with a jumping diffusion tensor.

**Table 3.**
Convergence analysis for the problem with the jumping diffusion tensor and velocity field.

| $h$ | Hexahedral | | Prismatic | | Tetrahedral | |
|------|------------|------------|------------|------------|------------|------------|
| | $\varepsilon_2^C$ | $\varepsilon_2^q$ | $\varepsilon_2^C$ | $\varepsilon_2^q$ | $\varepsilon_2^C$ | $\varepsilon_2^q$ |
| 1/10 | 1.46e-3 | 2.70e-3 | 6.78e-4 | 2.38e-3 | 8.42e-4 | 5.65e-3 |
| 1/20 | 3.76e-4 | 9.23e-4 | 1.90e-4 | 7.93e-4 | 2.36e-4 | 2.72e-3 |
| 1/40 | 9.58e-5 | 3.16e-4 | 5.08e-5 | 3.05e-4 | 5.96e-5 | 1.35e-3 |
| rate | 1.96 | 1.55 | 1.87 | 1.48 | 1.91 | 1.03 |

velocity **v** jump across the interface (see Fig. 4). Let $\mathbb{K}(\mathbf{x}) = \mathbb{K}^{(i)}$ for $\mathbf{x} \in \Omega^{(i)}$, where

$$\mathbb{K}^{(1)} = \begin{pmatrix} 3 & 1 & 0 \\ 1 & 3 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \qquad \mathbb{K}^{(2)} = \begin{pmatrix} 10 & 3 & 0 \\ 3 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

and the velocity field is $\mathbf{v}(\mathbf{x}) = \mathbf{v}^{(i)}$ for $\mathbf{x} \in \Omega^{(i)}$, where

$$\mathbf{v}^{(1)} = (1, 1, 0), \qquad \mathbf{v}^{(2)} = (1, 0.3, 0).$$

The spectral decomposition $\mathbb{K}^{(i)} = (W^{(i)})^T \Lambda^{(i)} W^{(i)}$ demonstrates the strong jump of the eigenvalues and the orientation of the eigenvectors of $\mathbb{K}(\mathbf{x})$:

$$\Lambda^{(1)} = \text{diag}\{4, 2, 1\}, \qquad \Lambda^{(2)} \approx \text{diag}\{10.908, 0.092, 1\}$$

$$W^{(1)} \approx \begin{pmatrix} 0.707 & 0.707 & 0 \\ -0.707 & 0.707 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \qquad W^{(2)} \approx \begin{pmatrix} 0.957 & 0.290 & 0 \\ -0.290 & 0.957 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

We define the following exact solution of (1.1) with $\Gamma_D = \partial\Omega$:

$$c(\mathbf{x}) = \begin{cases} 9 - 4x^2 - y^2 - 8xy - 6x + 4y, & \mathbf{x} \in \Omega^{(1)} \\ 6 - 2x^2 - y^2 - 2xy - x + y, & \mathbf{x} \in \Omega^{(2)} \end{cases}$$

so that the right-hand side is

$$g(\mathbf{x}) = \begin{cases} 44.0 - 16x - 10y, & \mathbf{x} \in \Omega^{(1)} \\ 53.3 - 4.6x - 2.6y, & \mathbf{x} \in \Omega^{(2)}. \end{cases}$$

The numerical tests were performed on the hexahedral, prismatic and tetrahedral meshes defined above. The meshes were generated so that the interface $x = 0.5$ is approximated by the mesh faces exactly. The solution isolines in the $xy$-plane are shown in Fig. 5. The convergence results presented in Table 3 demonstrate that the discontinuity of the diffusion tensor does not affect the convergence rate for all the considered meshes.

## 4.2. Monotonicity tests

At the first stage, we consider the advection-dominated problem with discontinuous Dirichlet boundary data. The discontinuity produces an internal shock in the solution, in addition to exponential boundary layers. This is a popular test case for the discretization schemes designed for advection-dominated problems (see [11, 13]). Following [11], we set

$$\mathbf{v} = \left(\cos\frac{\pi}{3}, -\sin\frac{\pi}{3}, 0\right), \qquad \mathbb{K} = v\mathbb{I}, \qquad v = 10^{-8}.$$

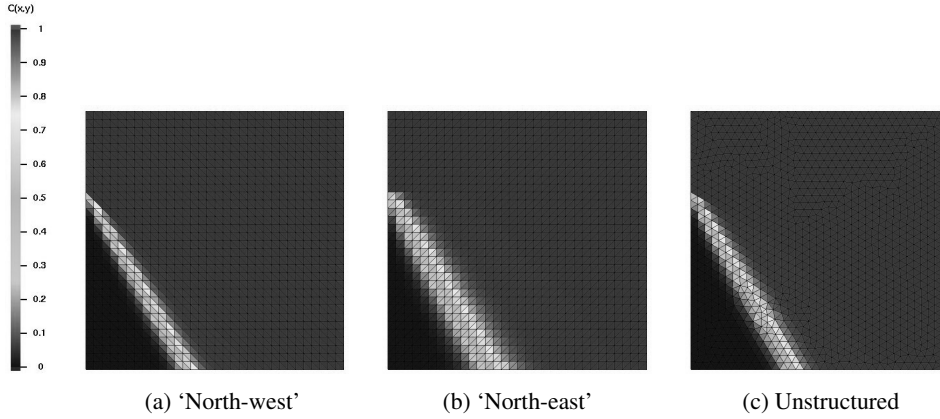(a) 'North-west'          (b) 'North-east'          (c) Unstructured

**Figure 6.** The numerical solutions for monotonicity test 1 ($h = 1/32$).

The Dirichlet boundary conditions are imposed as follows:

$$c(x,y,z) = \begin{cases} 0 & \text{if } x = 1 \text{ or } y \leqslant 0.7 \\ 1 & \text{otherwise.} \end{cases}$$

In order to preserve the solution independence of $z$, the homogeneous Neumann boundary conditions are set for $z = 0$ and $z = 1$.

The exact solution has a boundary layer next to two planes $y = 0$ and $x = 1$. It also has an internal layer along the stream plane passing through the line $(0, 0.7, z)$.

The numerical solution is non-negative in all cases that satisfy Theorems 3.1 and 3.2.

The computations were performed on prismatic meshes with structured (see Figs. 6a and 6b) and unstructured (see Fig. 6c) triangular grids at the base. The effective mesh resolution is set to $h = 1/64$ for measurements and $h = 1/32$ for figures.

In order to measure the quality of the numerical solution, the authors of [11] have proposed several estimates which quantify the solution oscillations and the smearing effects caused by a discretization scheme. We extend 2D measurements to 3D by taking values from a single mesh layer $0.5 - h \leqslant z \leqslant 0.5$. Let

$$\Omega_1 = \{(x,y,z)_h \in \Omega : x \leqslant 0.5, \quad y \geqslant 0.1, \quad 0.5 - h \leqslant z \leqslant 0.5\}$$

$$\Omega_2 = \{(x,y,z)_h \in \Omega : x \geqslant 0.7, \quad 0.5 - h \leqslant z \leqslant 0.5\}$$

and $\Omega_3$ denote a cell strip in the vicinity of the line $y = 0.25$,

$$\Omega_3 = \left\{T \in \mathscr{T} : \mathbf{x}_T = (x_T, y_T, z_T), \quad |y_T - 0.25| \leqslant |T|^{1/3}, \quad 0.5 - h \leqslant z \leqslant 0.5\right\}.$$

First we define estimate (4.1), which characterizes the total value of undershoots and overshoots in $\Omega_1$:

$$\mathrm{osc}_{\mathrm{int}} \equiv \left( \sum_{(x,y,z) \in \Omega_1} \left( \min\{0, c_h(x,y,z)\}\right)^2 + \left(\max\{0, c_h(x,y,z) - 1\}\right)^2 \right)^{1/2}. \quad (4.1)$$

**Table 4.**
'North-west' grid (see Fig. 6a), comparison of oscillations and smearing.

| Name | $osc_{int}$ | $osc_{exp}$ | $smear_{int}$ | $smear_{exp}$ |
|------|-------------|-------------|---------------|---------------|
| SUPG | 5.9e-1 | 2.1e-0 | 3.7e-2 | 5.6e-1 |
| MH85 | 6.1e-13 | 0 | 5.8e-2 | 1.1e-5 |
| FVMON | 7.8e-8 | 1.6e-6 | 4.7e-2 | 1.7e-5 |

**Table 5.**
'North-east' grid (see Fig. 6b), comparison of oscillations and smearing.

| Name | $osc_{int}$ | $osc_{exp}$ | $smear_{int}$ | $smear_{exp}$ |
|------|-------------|-------------|---------------|---------------|
| SUPG | 6.9e-1 | 3.8e-0 | 6.2e-2 | 1.7e-0 |
| MH85 | 0 | 0 | 1.0e-1 | 1.2e-5 |
| FVMON | 2.1e-11 | 1.7e-7 | 1.1e-1 | 1.8e-5 |

**Table 6.**
Unstructured grids (see Fig. 6c), comparison of oscillations and smearing.

| Name | $osc_{int}$ | $osc_{exp}$ | $smear_{int}$ | $smear_{exp}$ |
|------|-------------|-------------|---------------|---------------|
| SUPG | 5.9e-1 | 1.5e-0 | 5.5e-2 | 4.1e-1 |
| MH85 | 4.9e-15 | 1.8e-14 | 9.7e-2 | 5.3e-2 |
| FVMON | 3.5e-6 | 5.0e-7 | 5.9e-2 | 2.2e-5 |

Second, we define estimate (4.2), which quantifies the oscillations near the boundary layer in $\Omega_2$:

$$osc_{exp} \equiv \left( \sum_{(x,y,z) \in \Omega_2} \left( \max\{0, c_h(x,y,z) - 1\} \right)^2 \right)^{1/2}. \qquad (4.2)$$

Third, we define two estimates (4.3) and (4.4), which measure the thickness of the boundary layer and the internal shock, respectively:

$$smear_{exp} \equiv \left( \sum_{(x,y,z) \in \Omega_2} \left( \min\{0, c_h(x,y,z) - 1\} \right)^2 \right)^{1/2} \qquad (4.3)$$

$$smear_{int} \equiv x_2 - x_1 \qquad (4.4)$$

where

$$x_1 = \min_{\mathbf{x}_T \in \Omega_3,\, C(\mathbf{x}_T) \geqslant 0.1} x_T, \qquad x_2 = \max_{\mathbf{x}_T \in \Omega_3,\, C(\mathbf{x}_T) \leqslant 0.9} x_T.$$

For the continuous solution these estimates depend only on the diffusion process, so they are much smaller than the considered mesh size. For the numerical solution,
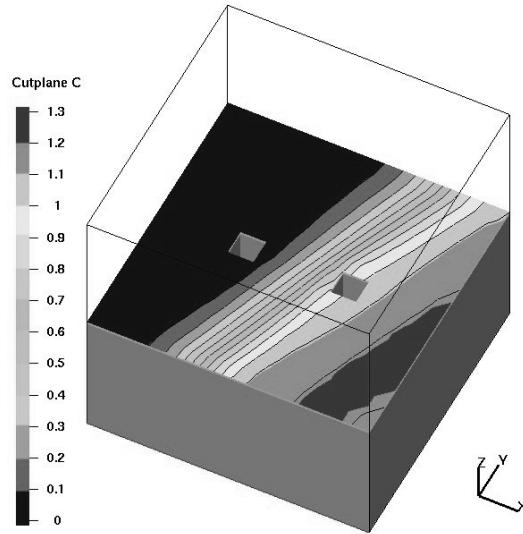
**Figure 7.** The solution isolines in the *xy*-plane for the problem with the no-flow outer boundary conditions.

**Table 7.**
The maximum concentration values for the problem with the no-flow outer boundary conditions.

| $h$ | 1/11 | 1/22 | 1/44 | 1/88 |
|---|---|---|---|---|
| $\max C$ | 1.882 | 1.219 | 1.041 | 1.020 |

small values of estimates (4.1)–(4.4) characterize an almost non-oscillatory and almost non-diffusive discrete solution. The smaller width of the internal shock and the boundary layer in a numerical solution, the less smearing is introduced by the discretization scheme.

The results obtained by the nonlinear FV, SUPG, and MH85 methods are shown in Tables 4, 5, and 6. Our method is competitive with the best 2D results presented in review [11]. The increase in the internal shock width on the 'north-east' prismatic mesh is caused by the larger cell size in the direction normal to the 'shock'.

At the second stage, we demonstrate that the FV discrete solution can violate the DMP even on cubic meshes.

We extend the test case investigated in [1, 9] to the advection–diffusion equation by introducing a velocity field. We consider a unit cube with two vertical holes $P_1$, $P_2$, $\Omega = (0,1)^3 \setminus (P_1 \cup P_2)$, $P_i = S_i \times (0,1)$, $i = 1, 2$, $S_1 = [3/11, 4/11] \times [5/11, 6/11]$, $S_2 = [7/11, 8/11] \times [5/11, 6/11]$. The domain boundary is split into the outer part $\Gamma_N$ where the homogeneous Neumann (no-flow) boundary condition is set, and two inner parts $\Gamma_{D,1}$, $\Gamma_{D,2}$ where the Dirichlet boundary conditions are set: $g_D(\mathbf{x}) = 0$, $\mathbf{x} \in \Gamma_{D,1}$, $g_D(\mathbf{x}) = 1$, $\mathbf{x} \in \Gamma_{D,2}$.

The anisotropic diffusion tensor is

$$\mathbb{K} = R_z(-\theta_z)R_y(-\theta_y)R_x(-\theta_x)\text{diag}(k_1,k_2,k_3)R_x(\theta_x)R_y(\theta_y)R_z(\theta_z) \qquad (4.5)$$

where $k_1 = k_3 = 1$, $k_2 = 10^{-3}$, $\theta_x = \theta_y = 0$, $\theta_z = 67.5°$, and $R_a(\alpha)$ is the rotation matrix in the plane orthogonal to *oa* with the angle $\alpha$.

The velocity field is $\mathbf{v} = (10^{-2}, 10^{-2}, 0)$. According to the maximum principle for elliptic PDEs the exact solution should be between 0 and 1 and have no extrema on the no-flow boundary $\Gamma_N$.

The FV discrete solution on the cubic mesh with $h = 1/32$ is shown in Fig. 7. It is non-negative in agreement with Theorem 3.2, but demonstrates overshoots near the no-flow boundary. These overshoots decrease rapidly as we refine the mesh, see Table 7.

## Conclusion

We have presented a new monotone finite volume method on polyhedral meshes for the advection–diffusion equation with a piecewise continuous full anisotropic diffusion tensor and a piecewise continuous advection field. The method is the natural extension of the 3D scheme for the diffusion equation [9] and the 2D scheme for the advection–diffusion equation [26]. The numerical solution is non-negative, provided that the source term and the Dirichlet boundary data are non-negative and the flux on the Neuman boundary is non-positive. The method does not require to interpolate the solution to the mesh nodes and can be applied to unstructured polyhedral meshes. The numerical experiments demonstrate the second-order convergence rate for the concentration and the first-order convergence rate for the flux on *randomly distorted meshes* in both advection-dominated and diffusion-dominated regimes.

## References

1. I. Aavatsmark, G. Eigestad, B. Mallison, and J. Nordbotten, A compact multipoint flux approximation method with improved robustness. *Numer. Meth. Partial Diff. Equations* (2008) **24**, No. 5, 1329–1360.

2. D. Benson, A new two-dimensional flux-limited shock viscosity for impact calculations. *Comp. Meth. Appl. Mech. Engrg.* (1991) **93**, 39–95.

3. E. Bertolazzi and G. Manzini, A cell-centered second-order accurate finite volume method for convection–diffusion problems on unstructured meshes. *Math. Models Meth. Appl. Sci.* (2004) **14**, No. 8, 1235–1260.

4. E. Bertolazzi and G. Manzini, A second-order maximum principle preserving finite volume method for steady convection–diffusion problems. *SIAM J. Numer. Anal.* (2005) **43**, No. 5, 2172–2199.

5. A. N. Brooks and T. J. R. Hughes, Streamline upwind/Petrov–Galerkin formulations for convection dominated flows with particular emphasis on the incompressible Navier–Stokes equations. *Comp. Meth. Appl. Mech. Engrg.* (1982) **32**, No. 1–3, 199–259.

6. E. Burman and A. Ern, Discrete maximum principle for Galerkin approximations of the Laplace operator on arbitrary meshes. *Comptes Rendus Mathematique* (2004) **338**, No. 8, 641–646.

7. G. Chavent and J. Jaffré, *Mathematical Models and Finite Elements for Reservoir Simulation.* Elsevier Science Publishers, B.V., Netherlands, 1986.

8. P. G. Ciarlet and P.-A. Raviart, Maximum principle and uniform convergence for the finite element method, *Comp. Meth. Appl. Mech. Engrg.* (1973) **2**, 17–31.

9. A. A. Danilov and Yu. V. Vassilevski, A monotone nonlinear finite volume method for diffusion equations on conformal polyhedral meshes. *Russ. Numer. Anal. Math. Modelling* (2009) **24**, No. 3, 207–227.

10. F. Gao, Y. Yuan, and D. Yang, An upwind finite-volume element scheme and its maximum-principle-preserving property for nonlinear convection–diffusion problem. *Int. J. Numer. Meth. Fluids* (2008) **56**, No. 12, 2301–2320.

11. V. John and P. Knobloch, On spurious oscillations at layers diminishing (SOLD) methods for convection–diffusion equations: Part I. A review. *Comp. Meth. Appl. Mech. Engrg.* (2007) **196**, 2197–2215.

12. M. E. Hubbard, Multidimensional slope limiters for MUSCL-type finite volume schemes on unstructured grids. *J. Comp. Phys.* (1999) **155**, No. 1, 54–74.

13. T. J. R. Hughes, M. Mallet, and A. Mizukami, A new finite element formulation for computational fluid dynamics. II. Beyond SUPG. *Comp. Meth. Appl. Mech. Engrg.* (1986) **54**, 341–355.

14. I. Kapyrin, A family of monotone methods for the numerical solution of three-dimensional diffusion problems on unstructured tetrahedral meshes. *Dokl. Math.* (2007) **76**, No. 2, 734–738.

15. S. Korotov, M. Křížek, and P. Neittaanmäki, Weakened acute type condition for tetrahedral triangulations and the discrete maximum principle. *Math. Comp.* (2001) **70**, No. 233, 107–119 (electronic).

16. D. Kuzmin, On the design of general-purpose flux limiters for finite element schemes, I. Scalar convection. *J. Comp. Phys.* (2006) **219**, No. 2, 513–531.

17. D. Kuzmin and M. Moller, Algebraic flux correction, I. Scalar conservation laws. In: *Flux-Corrected Transport: Principles, Algorithms, and Applications* (Eds. D. Kuzmin, R. Lohner, and S. Turek). Springer-Verlag, Berlin, 2005, pp. 155–206.

18. D. Kuzmin, M. J. Shashkov, and D. Svyatskiy, A constrained finite element method satisfying the discrete maximum principle for anisotropic diffusion problems. *J. Comp. Phys.* (2009) **228**, No. 9, 3448–3463.

19. O. Ladyzhenskaya and N. Uraltseva, *Linear and Quasilinear Elliptic Equations.* Nauka, Moscow, 1973 (in Russian).

20. S. Lamine and M. G. Edwards, Higher-resolution convection schemes for flow in porous media on highly distorted unstructured grids. *Int. J. Numer. Meth. Engrg.* (2008) **76**, No. 8, 1139–1158.

21. C. LePotier, Schema volumes finis monotone pour des operateurs de diffusion fortement anisotropes sur des maillages de triangle non structures. *C. C. Acad. Sci., Paris* (2005) **341**, 787–792.

22. C. LePotier, Finite volume scheme satisfying maximum and minimum principles for anisotropic diffusion operators. In: (Eds. R. Eymard and J.-M. Hérard). *Finite Volumes for Complex Applications V* (2008), 103–118.

23. R. J. LeVeque, *Finite Volume Methods for Hyperbolic Problems.* Cambridge Univ. Press, Cambridge, 2002.

24. K. Lipnikov, D. Svyatskiy, M. Shashkov, and Y. Vassilevski, Monotone finite volume schemes for diffusion equations on unstructured triangular and shape-regular polygonal meshes. *J. Comp. Phys.* (2007) **227**, 492–512.

25. K. Lipnikov, D. Svyatskiy, and Y. Vassilevski, Interpolation-free monotone finite volume method for diffusion equations on polygonal meshes. *J. Comp. Phys.* (2009) **228**, No. 3, 703–716.

26. K. Lipnikov, D. Svyatskiy, and Y. Vassilevski, A monotone finite volume method for advection–diffusion equations on unstructured polygonal meshes. *J. Comp. Phys.* (2010) **229**, 4017–4032.

27. R. Liska and M. Shashkov, Enforcing the discrete maximum principle for linear finite element solutions of second-order elliptic problems. *Commun. Comp. Phys.* (2008) **3**, No. 4, 852–877.

28. G. Manzini and A. Russo, A finite volume method for advection–diffusion problems in convection-dominated regimes. *Comp. Meth. Appl. Mech. Engrg.* (2008) **197**, No. 13–16, 1242–1261.

29. K. B. Nakshatrala and A. J. Valocchi, Non-negative mixed finite element formulations for a tensorial diffusion equation, *J. Comp. Phys.* (2008) **228**, No. 18, 6726–6752.

30. J. M. Nordbotten, I. Aavatsmark, and G. T. Eigestad, Monotonicity of control volume methods. *Numer. Math.* (2007) **106**, No. 2, 255–288.

31. A. Quarteroni and A. Valli, *Numerical Approximation of Partial Differential Equations.* Springer-Verlag, Heidelberg, SCM Series No. 23, 1994.

32. V. Ruas Santos, On the strong maximum principle for some piecewise linear finite element approximate problems of nonpositive type. *J. Fac. Sci. Univ. Tokyo Sect. IA Math.* (1982) **29**, No. 2, 473–491.

33. H. van der Vorst, Bi-CGSTAB: a fast and smoothly converging variant of Bi-CG for the solution of non-symmetric linear systems. *SIAM J. Sci. Stat. Comp.* (1992) **13**, No. 2, 631–644.

34. B. van Leer, Towards the ultimate conservative difference scheme. V. A second-order sequel to Godunov's method. *J. Comp. Phys.* (1979) **32**, No. 1, 101–136.

35. Yu. Vassilevski and I. Kapyrin, Two splitting schemes for nonstationary convection–diffusion problems on tetrahedral meshes. *Comp. Math. Math. Phys.* (2008) **48**, No. 8, 1349–1366.

36. D. Wollstein, T. Linss, and R. Hans-Gorg, Uniformly accurate finite volume discretization for a convection–diffusion problem *Electronic Transactions on Numer. Anal.* (2002) **13**, 1–11.

37. A. Yuan and Z. Sheng, Monotone finite volume schemes for diffusion equations on polygonal meshes. *J. Comp. Phys.* (2008) **227**, No. 12, 6288–6312.