

CONVERGENCE ANALYSIS OF FINITE VOLUME SCHEME FOR NONLINEAR TENSOR ANISOTROPIC DIFFUSION IN IMAGE PROCESSING*

OLGA DRBLÍKOVÁ[†] AND KAROL MIKULA[†]

Abstract. In this article we design the semiimplicit finite volume scheme for coherence enhancing diffusion in image processing and prove its convergence to the weak solution of the problem. The finite volume methods are natural tools for image processing applications since they use piecewise constant representation of approximate solutions similarly to the structure of digital images. They have been successfully applied in image processing, e.g., for solving the Perona–Malik equation or curvature-driven level set equations, where the nonlinearities are represented by a scalar function dependent on a solution gradient. Design of suitable finite volume schemes for tensor diffusion is a nontrivial task here we present the first such scheme with a convergence proof for the practical nonlinear model used in coherence-enhancing image smoothing. We provide basic information about this type of nonlinear diffusion including a construction of its diffusion tensor, and we derive a semi-implicit finite volume scheme for this nonlinear model with the help of covolume mesh. This method is well known as the diamond-cell method owing to the choice of covolume as a diamond-shaped polygon. Further, we prove a convergence of a discrete solution given by our scheme to the weak solution of the problem. The proof is based on Kolmogorov’s compactness theorem and a bounding of a gradient in the tangential direction by using a gradient in the normal direction. Finally computational results illustrated in figures are discussed.

Key words. image processing, nonlinear tensor diffusion, numerical solution, semiimplicit scheme, diamond-cell finite volume method, convergence

AMS subject classifications. 74S10, 65M12, 35K60, 68U10

DOI. 10.1137/070685038

1. Introduction. Nonlinear diffusion models are widely used nowadays in many practical tasks of image processing. In this paper we deal with the numerical solution of the model of tensor nonlinear anisotropic diffusion introduced by Weickert (see [23, 24, 22]) in the following form:

$$(1.1) \quad \frac{\partial u}{\partial t} - \nabla \cdot (D \nabla u) = 0 \quad \text{in } Q_T \equiv I \times \Omega,$$

$$(1.2) \quad u(x, 0) = u_0(x) \quad \text{in } \Omega,$$

$$(1.3) \quad (D \nabla u) \cdot \mathbf{n} = 0 \quad \text{on } I \times \partial\Omega,$$

where D is a matrix depending on the eigenvalues and eigenvectors of the so-called (regularized) structure tensor, $u_0 \in L^2(\Omega)$, and \mathbf{n} is the outer normal unit vector to $\partial\Omega$. Such a model is useful in any situation, where strong smoothing is desirable in a preferred direction and a low smoothing is expected in the perpendicular direction, e.g., for images with interrupted coherence of structures. To that goal the matrix

$$(1.4) \quad J_0(\nabla u_i) = \nabla u_i \otimes \nabla u_i = \nabla u_i \nabla u_i^T,$$

*Received by the editors March 12, 2007; accepted for publication (in revised form) September 10, 2007; published electronically December 12, 2007. This work was supported by EC projects Embryonics and BioEmergences of the 6th Framework program and by the projects VEGA 1/3321/06 and APVV-RPEU-0004-06.

<http://www.siam.org/journals/sinum/46-1/68503.html>

[†]Department of Mathematics, Slovak University of Technology, Radlinského 11, 813 68 Bratislava, Slovakia (drblikov@math.sk, mikula@math.sk).

where

$$(1.5) \quad u_{\tilde{t}}(x, t) = (G_{\tilde{t}} * u(\cdot, t))(x) \quad (\tilde{t} > 0)$$

is used. The matrix J_0 is symmetric and positive semidefinite, and its eigenvectors are parallel and orthogonal to $\nabla u_{\tilde{t}}$, respectively. We can average J_0 by applying another convolution with Gaussian G_ρ and define

$$(1.6) \quad J_\rho(\nabla u_{\tilde{t}}) = G_\rho * (\nabla u_{\tilde{t}} \otimes \nabla u_{\tilde{t}}) \quad (\rho > 0).$$

In computer vision the matrix $J_\rho = \begin{pmatrix} a & b \\ b & c \end{pmatrix}$ is known as a structure tensor, interest operator, or second moment matrix (see [9]). It is again symmetric and positive semidefinite, and its eigenvalues are given by

$$(1.7) \quad \mu_{1,2} = \frac{1}{2} \left(a + c \pm \sqrt{(a - c)^2 + 4b^2} \right), \quad \mu_1 \geq \mu_2.$$

Since the eigenvalues integrate the variation of the gray values within a neighborhood of size $O(\rho)$, they describe the average contrast in the eigendirections v and w .

With the help of the eigenvalues of J_ρ we can obtain useful information on the coherence. The expression $(\mu_1 - \mu_2)^2$ is large for anisotropic structures and tends to zero for isotropic structures; constant areas are characterized by $\mu_1 = \mu_2 = 0$, straight edges by $\mu_1 \gg \mu_2 = 0$, and corners by $\mu_1 \geq \mu_2 \gg 0$.

The corresponding orthogonal set of eigenvectors (v, w) to eigenvalues (μ_1, μ_2) is given by

$$(1.8) \quad \begin{aligned} v &= (v_1, v_2), & w &= (w_1, w_2), \\ v_1 &= 2b, & v_2 &= c - a + \sqrt{(a - c)^2 + 4b^2}, \\ w &\perp v, & w_1 &= -v_2, & w_2 &= v_1. \end{aligned}$$

The orientation of the eigenvector w which corresponds to the smaller eigenvalue μ_2 is called coherence orientation. This orientation has the lowest fluctuations.

One can use the above-mentioned structure tensor information in a construction of specific nonlinear diffusion filter [23, 24, 22]. The idea of the tensor nonlinear diffusion filtering is as follows. We get a processed version $u(x, t)$ of an original image $u_0(x)$ with a scale parameter $t \geq 0$ as the solution of mathematical model (1.1)–(1.3), where matrix D depends on solution u , satisfies smoothness, symmetry, and uniform positive definiteness properties, and steers a filtering process such that diffusion is strong along the coherence direction w and increases with the coherence $(\mu_1 - \mu_2)^2$. To that goal D must possess the same eigenvectors v and w as the structure tensor $J_\rho(\nabla u_{\tilde{t}})$, and we choose the eigenvalues of D as

$$(1.9) \quad \begin{aligned} \kappa_1 &= \alpha, \quad \alpha \in (0, 1), \quad \alpha \ll 1, \\ \kappa_2 &= \begin{cases} \alpha & \text{if } \mu_1 = \mu_2, \\ \alpha + (1 - \alpha) \exp\left(\frac{-C}{(\mu_1 - \mu_2)^2}\right), & C > 0 \quad \text{otherwise.} \end{cases} \end{aligned}$$

The matrix D then has the following form:

$$(1.10) \quad D = ABA^{-1},$$

where $A = \begin{pmatrix} v_1 & -v_2 \\ v_2 & v_1 \end{pmatrix}$ and $B = \begin{pmatrix} \kappa_1 & 0 \\ 0 & \kappa_2 \end{pmatrix}$. The exponential function is used in (1.9) because it ensures that the smoothness of the structure tensor carries over to the

diffusion tensor and that κ_2 does not exceed 1. The positive parameter α guarantees that the process never stops. Even if $(\mu_1 - \mu_2)^2$ tends to zero so the structure becomes isotropic, there still remains some small linear diffusion with diffusivity $\alpha > 0$. Such α is a regularization parameter, which keeps the diffusion tensor uniformly positive definite. C has the role of a threshold parameter. If $(\mu_1 - \mu_2)^2 \gg C$, then $\kappa_2 \approx 1$, and, in contrast, if $(\mu_1 - \mu_2)^2 \ll C$, then $\kappa_2 \approx \alpha$. Due to the convolutions in (1.5) and (1.6), the elements of matrix D are C^∞ functions. Such a model is a nontrivial extension of the regularized Perona–Malik equation [17, 1, 15], and, as well as further PDEs employing tensor diffusion, it is used in many practical image processing applications; see, e.g., [23, 24, 22, 6, 13, 19, 18]. In section 5 of this paper we also illustrate its usefulness by smoothing and segmenting the cell membrane images obtained by a confocal microscope. We show that after application of the nonlinear tensor anisotropic diffusion using our numerical scheme the coherent structures are attenuated. If such improved edge information is used in the so-called subjective surface segmentation method [20, 16, 2], the cell boundaries are correctly segmented.

There are only a few purely finite volume methods designed and studied from the numerical analysis point of view for solving tensor diffusion problems; see, e.g., [3, 4, 5, 26] devoted to discretization of the elliptic operators. On the other hand, finite volume schemes for nonlinear parabolic problems as arising in image analysis are natural since they use piecewise constant representation of approximate solutions similarly to the structure of digital images. Finite and complementary volume schemes have been used successively in image processing for solving the Perona–Malik equation and its generalizations (see, e.g., [15, 11, 12, 10, 7, 21]) and for solving the generalized curvature-driven level set equations (see, e.g., [8, 16, 2]) where the nonlinearities are represented by a scalar function dependent on a solution gradient. Here we present the first finite volume scheme with a convergence proof for the highly nonlinear anisotropic tensor diffusion model arising in coherence-enhancing image smoothing.

The next section is devoted to derivation of our numerical scheme, in section 3 we study the existence and uniqueness of discrete solutions, section 4 contains our convergence proof, and finally, in section 5, we discuss numerical experiments.

2. Finite volume scheme for nonlinear tensor anisotropic diffusion. The aim of this section is to derive our computational method. Let the image be represented by $n_1 \times n_2$ pixels (finite volumes) such that it looks like a mesh with n_1 rows and n_2 columns. Let $\Omega = (0, n_1 h) \times (0, n_2 h)$, h is a pixel size, and let the image $u(x)$ be given by a bounded mapping $u : \Omega \rightarrow R$. The filtering process is considered in a time interval $I = [0, T]$. Let $0 = t_0 \leq t_1 \leq \dots \leq t_{N_{max}} = T$ denote the time discretization, with $t_n = t_{n-1} + k$, where k is the length of a discrete time step. In our scheme we will look for u^n , an approximation of solution at time t_n , for every $n = 1, \dots, N_{max}$. As usual in finite volume methods, we integrate (1.1) over finite volume K , then provide a semiimplicit time discretization, and use a divergence theorem to get

$$(2.1) \quad \frac{u_K^n - u_K^{n-1}}{k} m(K) - \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{int}} \int_{\sigma} (D^{n-1} \nabla u^n) \cdot \mathbf{n}_{K,\sigma} ds = 0,$$

where u_K^n , $K \in \mathcal{T}_h$, represents the mean value of u^n on K . \mathcal{T}_h is an admissible finite volume mesh (see [4]), and further quantities and notations are described as follows: $m(K)$ is the measure of the finite volume K with boundary ∂K , and $\sigma_{KL} = K \cap L = K|L$ is an edge of the finite volume K , where $L \in \mathcal{T}_h$ is an adjacent finite volume to K such that $m(K \cap L) \neq 0$. Due to simplifying notation, we use σ instead of

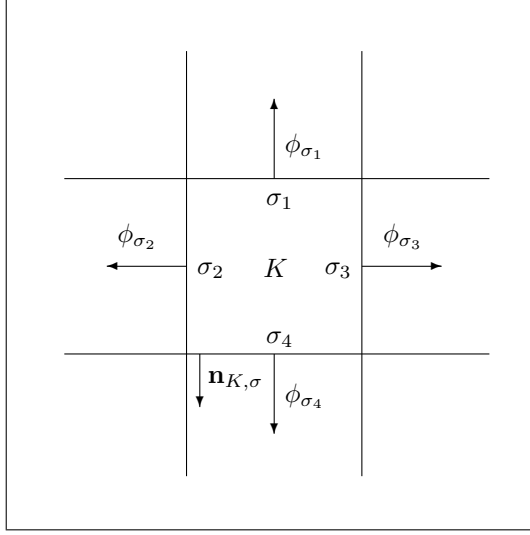


FIG. 2.1. A detail of a finite volume mesh—a finite volume K , its boundaries σ_i , $i = 1, 2, 3, 4$, and the fluxes outward to a finite volume K .

σ_{KL} at several places if no confusion can appear. \mathcal{E}_K is a set of edges such that $\partial K = \bigcup_{\sigma \in \mathcal{E}_K} \sigma$ and $\mathcal{E} = \bigcup_{K \in \mathcal{T}_h} \mathcal{E}_K$. The set of boundary edges is denoted by \mathcal{E}_{ext} , that is, $\mathcal{E}_{ext} = \mathcal{E} \cap \partial\Omega$, and let $\mathcal{E}_{int} = \mathcal{E} \setminus \mathcal{E}_{ext}$. Υ is the set of pairs of adjacent finite volumes, defined by $\Upsilon = \{(K, L) \in \mathcal{T}_h^2, K \neq L, m(K|L) \neq 0\}$, and $\mathbf{n}_{K,\sigma}$ is the normal unit vector to σ outward to K . See Figure 2.1.

Let us define our discrete numerical solution by

$$(2.2) \quad u_{h,k}(x, t) = \sum_{n=0}^{N_{max}} \sum_{K \in \mathcal{T}_h} u_K^n \chi\{x \in K\} \chi\{t_{n-1} < t \leq t_n\},$$

where the function $\chi(A)$ is defined as

$$(2.3) \quad \chi_{\{A\}} = \begin{cases} 1 & \text{if } A \text{ is true,} \\ 0 & \text{elsewhere.} \end{cases}$$

The extension of the function (2.2) outside Ω is given first by its periodic mirror reflection in $\Omega_{\tilde{t}}$, where \tilde{t} is the width of the smoothing kernel:

$$(2.4) \quad \Omega_{\tilde{t}} = \Omega \cup B_{\tilde{t}}(x), \quad x \in \partial\Omega,$$

$B_{\tilde{t}}(x)$ is a ball centered at x with radius \tilde{t} , and then we extend this periodic mirror reflection by 0 outside $\Omega_{\tilde{t}}$ and denote it by $\tilde{u}_{h,k}$.

In our scheme we will start computation by defining initial values

$$(2.5) \quad u_K^0 = \frac{1}{m(K)} \int_K u_0(x) dx, \quad K \in \mathcal{T}_h$$

and letting $u_{h,k}^n(x) = \sum_{K \in \mathcal{T}_h} u_K^n \chi\{x \in K\}$ denote a finite volume approximation at the n th time step. In order to get the scheme we write (2.1) in the form $\frac{u_K^n - u_K^{n-1}}{k} - \frac{1}{m(K)} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{int}} \phi_\sigma^n(u_{h,k}^n) m(\sigma) = 0$, where $m(\sigma)$ is the measure of edge σ and $\phi_\sigma^n(u_{h,k}^n)$

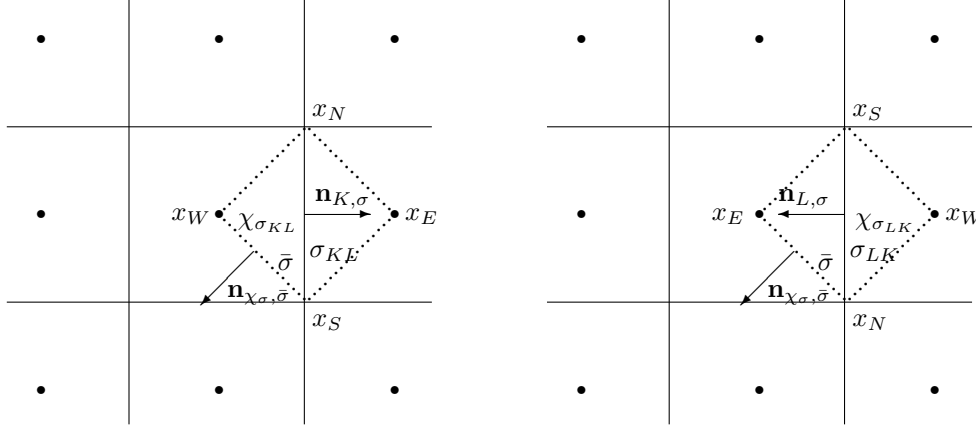


FIG. 2.2. A detail of a mesh. The covolumes associated to edges $\sigma = \sigma_{KL}$ (left) and $\sigma = \sigma_{LK}$ (right).

denotes an approximation of the exact averaged flux $\frac{1}{m(\sigma)} \int_{\sigma} (D^{n-1} \nabla u^n) \cdot \mathbf{n}_{K,\sigma} ds$ for any K and $\sigma \in \mathcal{E}_K$.

We construct $\phi_{\sigma}^n(u_{h,k}^n)$ with the help of a covolume mesh (see, e.g., [3]). The covolume χ_{σ} associated to σ is constructed around each edge by joining end points of this edge and midpoints of finite volumes which are common to this edge; see Figure 2.2. We denote the end points of an edge $\bar{\sigma} \subset \partial \chi_{\sigma}$ by $N_1(\bar{\sigma})$ and $N_2(\bar{\sigma})$ and let $\mathbf{n}_{\chi_{\sigma},\bar{\sigma}}$ be the normal unit vector to $\bar{\sigma}$ outward to χ_{σ} . In order to approximate diffusion flux, using the divergence theorem, we first derive an approximation of the averaged gradient on χ_{σ} , namely, $\frac{1}{m(\chi_{\sigma})} \int_{\chi_{\sigma}} \nabla u^n dx = \frac{1}{m(\chi_{\sigma})} \int_{\partial \chi_{\sigma}} u^n \mathbf{n}_{\chi_{\sigma},\bar{\sigma}} ds$, and then we approximate it by $p_{\sigma}^n(u) = \frac{1}{m(\chi_{\sigma})} \sum_{\bar{\sigma} \in \partial \chi_{\sigma}} \frac{1}{2} (u_{N_1(\bar{\sigma})}^n + u_{N_2(\bar{\sigma})}^n) m(\bar{\sigma}) \mathbf{n}_{\chi_{\sigma},\bar{\sigma}}$. Let the values at x_E and x_W be taken as u_E and u_W , and let the values u_S and u_N at the vertices x_N and x_S be computed as the arithmetic mean of u_K , where K are finite volumes which are common to this vertex.

Since our mesh is uniform and squared, we can use the following relations: $m(\chi_{\sigma}) = \frac{h^2}{2}$ and $m(\bar{\sigma}) = \frac{\sqrt{2}}{2}h$, and after a short calculation we are ready to write

$$(2.6) \quad p_{\sigma}^n(u) = \frac{u_E^n - u_W^n}{h} \mathbf{n}_{K,\sigma} + \frac{u_N^n - u_S^n}{h} \mathbf{t}_{K,\sigma},$$

where $\mathbf{t}_{K,\sigma}$ is a unit vector parallel to σ such that $(x_N - x_S) \cdot \mathbf{t}_{K,\sigma} > 0$. Although such u_N^n , u_W^n , u_E^n , and u_S^n correspond to a particular edge σ , and so we should denote them by $u_{N_{\sigma}}^n$, $u_{W_{\sigma}}^n$, $u_{E_{\sigma}}^n$, and $u_{S_{\sigma}}^n$, respectively, in (2.6), we will use the above simplified notations. By replacing the exact gradient ∇u^n by the numerical gradient $p_{\sigma}^n(u)$ in the approximation of $\phi_{\sigma}^n(u_{h,k}^n)$, we get the numerical flux in the form

$$(2.7) \quad \phi_{\sigma}^n(u_{h,k}^n) = (D_{\sigma} p_{\sigma}^n(u)) \cdot \mathbf{n}_{K,\sigma},$$

where $D_{\sigma} = D_{\sigma}^{n-1} = \begin{pmatrix} \bar{\lambda}_{\sigma} & \bar{\beta}_{\sigma} \\ \bar{\beta}_{\sigma} & \bar{\nu}_{\sigma} \end{pmatrix}$ is an approximation of the mean value of matrix D along σ evaluated at the previous time step. To that goal we take $u_{h,k}^{n-1}$ for constructing the structure and diffusion tensor and evaluate them at x_{KL} , where x_{KL} is a point of $\sigma_{KL} = K|L$ intersecting the segment $x_K x_L$. From an implementation point of view, the structure and then diffusion tensor evaluation can be done in two ways.

Either we can replace gradients of u appearing in structure tensor by their numerical approximation $p_\sigma^n(u)$ and then smooth them by weighted average (convolution), or we can evaluate $\nabla G_{\tilde{t}} * u_{h,k}^{n-1}$ using weights given by $\nabla G_{\tilde{t}}$ applied to discrete piecewise constant values of $u_{h,k}^{n-1}$ as convolution realization. In the latter way we do not introduce additional approximation into the scheme, and in the part devoted to convergence analysis we use the latter approach, although both are realizable computationally.

It is important to note that in (2.7) we always consider the matrix D_σ written in the basis $(\mathbf{n}_{K,\sigma}, \mathbf{t}_{K,\sigma})$; cf. [3]. Although it may look artificial, it will simplify further considerations. In practice it means that (cf. Figure 2.1) if the matrix D is given in standard basis on edge σ by $\begin{pmatrix} \lambda_\sigma & \beta_\sigma \\ \beta_\sigma & \nu_\sigma \end{pmatrix}$, then $D_\sigma = \begin{pmatrix} \lambda_\sigma & \beta_\sigma \\ \beta_\sigma & \nu_\sigma \end{pmatrix}$, i.e., $\bar{\lambda}_\sigma = \lambda_\sigma$, $\bar{\beta}_\sigma = \beta_\sigma$, $\bar{\nu}_\sigma = \nu_\sigma$ for the two edges $\sigma = \sigma_2$ and σ_3 . On the other hand, $D_\sigma = \begin{pmatrix} \nu_\sigma & -\beta_\sigma \\ -\beta_\sigma & \lambda_\sigma \end{pmatrix}$, i.e., $\bar{\lambda}_\sigma = \nu_\sigma$, $\bar{\beta}_\sigma = -\beta_\sigma$, $\bar{\nu}_\sigma = \lambda_\sigma$ for other two edges $\sigma = \sigma_1$ and σ_4 . By using such a matrix representation, the definition (2.7) can be written in this compact form

$$(2.8) \quad \phi_\sigma^n(u_{h,k}^n) = \left[\begin{pmatrix} \bar{\lambda}_\sigma & \bar{\beta}_\sigma \\ \bar{\beta}_\sigma & \bar{\nu}_\sigma \end{pmatrix} \begin{pmatrix} \frac{u_E^n - u_W^n}{h} \\ \frac{u_N^n - u_S^n}{h} \end{pmatrix} \right] \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \bar{\lambda}_\sigma \frac{u_E^n - u_W^n}{h} + \bar{\beta}_\sigma \frac{u_N^n - u_S^n}{h},$$

since in the basis $(\mathbf{n}_{K,\sigma}, \mathbf{t}_{K,\sigma})$ the formula (2.6) can be written for each edge as

$$(2.9) \quad p_\sigma^n(u) = \begin{pmatrix} \frac{u_E^n - u_W^n}{h} \\ \frac{u_N^n - u_S^n}{h} \end{pmatrix}$$

and $\mathbf{n}_{K,\sigma}$ is equal to $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ in the basis $(\mathbf{n}_{K,\sigma}, \mathbf{t}_{K,\sigma})$ for each edge σ . Because of the convolutions in (1.5) and (1.6), the elements of matrix D_σ are C^∞ functions.

Finally, let us summarize our *semimplicit finite volume scheme*:

$$(2.10) \quad \frac{u_K^n - u_K^{n-1}}{k} - \frac{1}{m(K)} \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{int}} \phi_\sigma^n(u_{h,k}^n) m(\sigma) = 0,$$

where

$$(2.11) \quad \phi_\sigma^n(u_{h,k}^n) = \bar{\lambda}_\sigma \frac{u_E^n - u_W^n}{h} + \bar{\beta}_\sigma \frac{u_N^n - u_S^n}{h}.$$

3. Existence and uniqueness of the solution to the discrete scheme.

In order to fulfill the main goal of this section, to prove the existence and uniqueness of u_K^n , $K \in \mathcal{T}_h$, we estimate the expressions $u_N^n - u_S^n$ by means of $u_E^n - u_W^n$ for all edges σ . To that goal we use mainly the results of [3] in our situation. Let us note that, due to simplification of notation, we do not use upper index n in what follows, and at some places we relate u_E and u_W to particular edge σ using u_{E_σ} , u_{W_σ} , etc. In the following we denote by C_i constants which may depend on the properties of the diffusion tensor.

DEFINITION 3.1. *Let P_σ be the set of all edges δ perpendicular to σ (see Figure 3.1 for two illustrative situations when $\sigma = \sigma_{WE}$ and $\sigma = \sigma_{EW}$), which have a common vertex with σ and fulfill the following conditions: $x_{E_\delta} - x_{W_\delta} > 0$ if $x_{N_\sigma} - x_{S_\sigma} > 0$ and $x_{E_\delta} - x_{W_\delta} < 0$ if $x_{N_\sigma} - x_{S_\sigma} < 0$. Let us note that $x_{W_\sigma} = x_{W_\delta}^1 = x_{E_\delta}^3$ for $\sigma = \sigma_{WE}$, $x_{E_\sigma} = x_{W_\delta}^2 = x_{E_\delta}^4$ for $\sigma = \sigma_{WE}$, $x_{W_\sigma} = x_{E_\delta}^2 = x_{W_\delta}^4$ for $\sigma = \sigma_{EW}$, and $x_{E_\sigma} = x_{E_\delta}^1 = x_{W_\delta}^3$ for $\sigma = \sigma_{EW}$.*

Using definitions given in the previous section we can write

$$(3.1) \quad u_N - u_S = \frac{1}{4} [(u_E^1 - u_W^1) + (u_E^3 - u_W^3) + (u_E^2 - u_W^2) + (u_E^4 - u_W^4)],$$

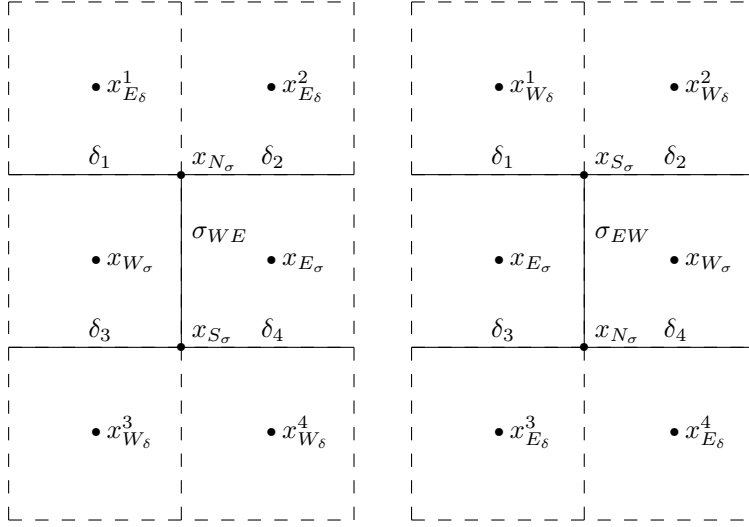


FIG. 3.1. Left: An edge σ_{WE} and edges $\delta_1, \delta_2, \delta_3, \delta_4 \in P_{\sigma_{WE}}$. Right: An edge σ_{EW} and edges $\delta_1, \delta_2, \delta_3, \delta_4 \in P_{\sigma_{EW}}$.

where $u_E^1 = u_{E_{\delta_1}}$ and $u_W^1 = u_{W_{\delta_1}}$ correspond to edge δ_1 and similarly $u_E^2, u_W^2, u_E^3, u_W^3, u_E^4$, and u_W^4 correspond to edges δ_2, δ_3 and δ_4 . Applying the inequality $(a+b)^2 \leq 2a^2 + 2b^2$ we have

$$(3.2) \quad (u_{N_\sigma} - u_{S_\sigma})^2 \leq \sum_{\delta \in P_\sigma \cap \mathcal{E}_{int}} \frac{1}{4} (u_{E_\delta} - u_{W_\delta})^2.$$

Multiplying (3.2) by $(\frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma})^2 \frac{\bar{\lambda}_\sigma}{h^2}$ and summing for all $\sigma \in \mathcal{E}_{int}$ (by σ we mean σ_{WE}) we obtain

$$(3.3) \quad \sum_{\sigma \in \mathcal{E}_{int}} \left(\frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma} \right)^2 \left(\frac{u_{N_\sigma} - u_{S_\sigma}}{h} \right)^2 \bar{\lambda}_\sigma \leq \sum_{\sigma \in \mathcal{E}_{int}} \left(\frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma} \right)^2 \sum_{\delta \in P_\sigma \cap \mathcal{E}_{int}} \frac{1}{4} \left(\frac{u_{E_\delta} - u_{W_\delta}}{h} \right)^2 \bar{\lambda}_\sigma.$$

Then we swap the two sums on the right-hand side of (3.3) to get

$$(3.4) \quad \sum_{\sigma \in \mathcal{E}_{int}} \left(\frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma} \right)^2 \left(\frac{u_{N_\sigma} - u_{S_\sigma}}{h} \right)^2 \bar{\lambda}_\sigma \leq \sum_{\delta \in \mathcal{E}_{int}} \gamma_\delta \left(\frac{u_{E_\delta} - u_{W_\delta}}{h} \right)^2 \bar{\lambda}_\delta,$$

where

$$(3.5) \quad \gamma_\delta = \sum_{\sigma \in P_\delta \cap \mathcal{E}_{int}} \frac{1}{4} \left(\frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma} \right)^2 \frac{\bar{\lambda}_\sigma}{\bar{\lambda}_\delta}.$$

Let us consider the matrix $\begin{pmatrix} \bar{\lambda}_{\sigma^\perp} & \bar{\beta}_{\sigma^\perp} \\ \bar{\beta}_{\sigma^\perp} & \bar{\nu}_{\sigma^\perp} \end{pmatrix}$, which is the matrix D written in the basis $(\mathbf{t}_{K,\sigma}, -\mathbf{n}_{K,\sigma})$ on edge σ . Due to the smoothness of D we get

$$(3.6) \quad \bar{\lambda}_\sigma = \bar{\nu}_{\sigma^\perp} = \bar{\nu}_\delta(1 + O(h)) = \bar{\lambda}_{\delta^\perp}(1 + O(h)), \quad \delta \in P_\sigma,$$

$$(3.7) \quad \bar{\beta}_\sigma = -\bar{\beta}_{\sigma^\perp} = -\bar{\beta}_\delta(1 + O(h)) = \bar{\beta}_{\delta^\perp}(1 + O(h)), \quad \delta \in P_\sigma,$$

$$(3.8) \quad \bar{\nu}_\sigma = \bar{\lambda}_{\sigma^\perp} = \bar{\lambda}_\delta(1 + O(h)) = \bar{\nu}_{\delta^\perp}(1 + O(h)), \quad \delta \in P_\sigma.$$

By applying (3.6)–(3.8) in (3.5) and using $\frac{1}{1-Ch} \leq 1 + (C + 2C^2h)h$ for h sufficiently small ($h \leq \frac{1}{2C}$), we have

$$\gamma_\delta \leq \sum_{\sigma \in P_\delta \cap \mathcal{E}_{int}} \frac{1}{4} \left(\frac{\bar{\beta}_{\delta^\perp}}{\bar{\lambda}_{\delta^\perp}} \right)^2 \frac{\bar{\lambda}_{\delta^\perp}}{\bar{\lambda}_\delta} (1 + O(h)) = \left(\frac{\bar{\beta}_{\delta^\perp}}{\bar{\lambda}_{\delta^\perp}} \right)^2 \frac{\bar{\lambda}_{\delta^\perp}}{\bar{\lambda}_\delta} (1 + O(h)).$$

By using the positive definiteness of the diffusion tensor written in a standard basis as $\begin{pmatrix} \lambda_\delta & \beta_\delta \\ \beta_\delta & \nu_\delta \end{pmatrix}$, we obtain for its determinant

$$(3.9) \quad \lambda_\delta \nu_\delta - \beta_\delta^2 > 0.$$

Now we have two possibilities for γ_δ . Let δ be an arbitrary edge in the mesh parallel to σ_3 (see Figure 2.1). Then $\gamma_\delta \leq \left(\frac{-\beta_\delta}{\nu_\delta} \right)^2 \frac{\nu_\delta}{\lambda_\delta} (1 + O(h)) = \frac{(\beta_\delta)^2}{\lambda_\delta \nu_\delta} (1 + O(h)) < 1$ for h sufficiently small due to (3.9). Similarly, if δ is any edge oriented perpendicularly to σ_3 , we have $\gamma_\delta \leq \left(\frac{\beta_\delta}{\lambda_\delta} \right)^2 \frac{\lambda_\delta}{\nu_\delta} (1 + O(h)) = \frac{(\beta_\delta)^2}{\lambda_\delta \nu_\delta} (1 + O(h)) < 1$ for h sufficiently small. Thus, due to the fact that $\lambda_\sigma \geq C > 0$ and $\nu_\sigma \geq C > 0$, we obtain $0 \leq \gamma_\delta < 1$ for h sufficiently small. Since this condition is fulfilled for each edge δ , we can rewrite (3.4) as

$$(3.10) \quad \sum_{\sigma \in \mathcal{E}_{int}} \left(\frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma} \right)^2 \left(\frac{u_N - u_S}{h} \right)^2 \bar{\lambda}_\sigma \leq \gamma \sum_{\sigma \in \mathcal{E}_{int}} \left(\frac{u_E - u_W}{h} \right)^2 \bar{\lambda}_\sigma,$$

where $0 \leq \gamma < 1$, $\gamma = \max_{\sigma \in \mathcal{E}} \gamma_\sigma$.

Let us now introduce the space of piecewise constant functions associated to our mesh and discrete H^1 norm for this space. This discrete norm will be used to obtain some estimates on the approximate solution given by the finite volume scheme.

DEFINITION 3.2. *Let Ω be an open bounded polygonal subset of R^2 . Let \mathcal{T}_h be an admissible finite volume mesh (see [4]). We define $\mathcal{P}_0(\mathcal{T}_h)$ as the set of functions from Ω to R which are constant over each finite volume K of the mesh \mathcal{T}_h .*

DEFINITION 3.3. *Let Ω be an open bounded polygonal subset of R^2 . For $u \in \mathcal{P}_0(\mathcal{T}_h)$ we define*

$$(3.11) \quad |u_{h,k}^n|_{1,\mathcal{T}_h} = \left(\sum_{(K,L) \in \Upsilon} \frac{(u_L - u_K)^2}{d_{K,L}} m(\sigma) \right)^{\frac{1}{2}},$$

where $d_{K,L}$ is the Euclidean distance between x_K and x_L .

Remark that (3.11) can be rewritten for our uniform mesh into the following form:

$$(3.12) \quad |u_{h,k}^n|_{1,\mathcal{T}_h} = \left(2 \sum_{\sigma \in \mathcal{E}_{int}} \left(\frac{u_E - u_W}{h} \right)^2 m(\chi_\sigma) \right)^{\frac{1}{2}},$$

where $\sigma = \sigma_{WE}$. Let us define a discrete operator \mathcal{L}_h by

$$\mathcal{L}_h(u_{h,k}^n) = u_K^n m(K) - k \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{int}} \phi_\sigma^n(u_{h,k}^n) m(\sigma).$$

Then solution $u_{h,k}^n \in \mathcal{P}_0(\mathcal{T}_h)$ of our scheme at time t_n is given by

$$(3.13) \quad \mathcal{L}_h(u_{h,k}^n) = f_{h,k}(u_{h,k}^{n-1}),$$

where $f_{h,k}(u_{h,k}^{n-1}) = u_K^{n-1}m(K)$, $K \in \mathcal{T}_h$, and u_K^{n-1} is the value of the piecewise constant function $u_{h,k}^{n-1}$ in K . This equality is a linear system of N equations with N unknowns u_K^n , $K \in \mathcal{T}_h$, $N = \text{card}(\mathcal{T}_h)$.

Multiplying $\mathcal{L}_h(u_{h,k})$ by u_K^n , summing over K , and splitting into parts A and B leads to

$$(3.14) \quad \sum_{K \in \mathcal{T}_h} \mathcal{L}_h(u_{h,k}) u_K^n = A + B,$$

with

$$(3.15) \quad A = \sum_{K \in \mathcal{T}_h} (u_K^n)^2 m(K) = \|u_{h,k}^n\|_{L^2(\Omega)}^2$$

and

$$B = k \sum_{K \in \mathcal{T}_h} u_K^n \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{int}} -\phi_\sigma^n(u_{h,k}^n) m(\sigma).$$

The above expression can be written in the following form:

$$(3.16) \quad \begin{aligned} B &= k \sum_{W \in \mathcal{T}_h} u_W^n \sum_{\sigma \in \mathcal{E}_W \cap \mathcal{E}_{int}} -\phi_\sigma^n(u_{h,k}^n) m(\sigma) \\ &= \frac{k}{2} \sum_{\sigma \in \mathcal{E}_{int}} \phi_\sigma^n(u_{h,k}^n) \frac{u_E - u_W}{h} 2m(\chi_\sigma) = Q(u_{h,k}^n) \end{aligned}$$

owing to property $\phi_\sigma^n(u_{h,k}^n) = \phi_{\sigma_{WE}}^n(u_{h,k}^n) = -\phi_{\sigma_{EW}}^n(u_{h,k}^n)$. Since $\phi_\sigma^n(u_{h,k}^n) = 0$ for $\sigma \in \mathcal{E}_{ext}$, we can extend the sum in (3.16) and write

$$Q(u_{h,k}^n) = \frac{k}{2} \sum_{\sigma \in \mathcal{E}} (D_\sigma p_\sigma^*) \cdot p_\sigma 2m(\chi_\sigma) = k(D_h p_h^*, p_h)_{L^2(\Omega)}$$

where $p_\sigma^* = \frac{u_E - u_W}{h} \mathbf{n}_{W,\sigma}$ for $\sigma = \sigma_{WE}$ is the normal component of the gradient and D_h, p_h, p_h^* are piecewise constant functions with values extended from σ to χ_σ . Further, we use the following inequality:

$$(3.17) \quad (D_h p_h^*, p_h)_{L^2(\Omega)} \geq (D_h p_h^*, p_h^*)_{L^2(\Omega)} - |(D_h p_h^*, p_h - p_h^*)_{L^2(\Omega)}|.$$

It is clear that $(D_h p_h^*, p_h^*)_{L^2(\Omega)} = \sum_{\sigma \in \mathcal{E}} \bar{\lambda}_\sigma \left(\frac{u_E - u_W}{h} \right)^2 m(\chi_\sigma)$, due to the fact that $u_E - u_W = 0$ for $\sigma \in \mathcal{E}_{ext}$ thanks to reflexion of $u_{h,k}$ in Ω_i (see the previous section). Applying Young's inequality in the second term on the right-hand side of (3.17) leads to

$$(3.18) \quad \begin{aligned} |(D_h p_h^*, p_h - p_h^*)_{L^2(\Omega)}| &= \left| \sum_{\sigma \in \mathcal{E}} \bar{\beta}_\sigma \frac{u_E - u_W}{h} \frac{u_N - u_S}{h} m(\chi_\sigma) \right| \\ &\leq \sum_{\sigma \in \mathcal{E}_{int}} \frac{1}{2} \left[\left(\frac{u_E - u_W}{h} \right)^2 + \left(\frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma} \right)^2 \left(\frac{u_N - u_S}{h} \right)^2 \right] \bar{\lambda}_\sigma m(\chi_\sigma), \end{aligned}$$

since $\phi_\sigma^n(u_{h,k}^n) = 0$ for $\sigma \in \mathcal{E}_{ext}$. By using inequalities (3.10) we get

$$(3.19) \quad \begin{aligned} |(D_h p_h^*, p_h - p_h^*)_{L^2(\Omega)}| &\leq \frac{1+\gamma}{2} \sum_{\sigma \in \mathcal{E}_{int}} \bar{\lambda}_\sigma \left(\frac{u_E - u_W}{h} \right)^2 m(\chi_\sigma) \\ &= \frac{1+\gamma}{2} (D_h p_h^*, p_h^*)_{L^2(\Omega)}. \end{aligned}$$

Using (3.12), it in turn implies

$$(3.20) \quad Q(u_{h,k}^n) \geq \left(1 - \frac{1+\gamma}{2}\right) k(Dp_h^*, p_h^*)_{L^2(\Omega)} \geq \bar{\lambda}_{\min} \frac{1-\gamma}{2} \frac{k}{2} |u_{h,k}^n|_{1,\mathcal{T}_h}^2,$$

where $\bar{\lambda}_{\min} = \inf_{\sigma \in \mathcal{E}} \bar{\lambda}_\sigma \geq C > 0$. By applying (3.15), (3.16), and (3.20) in (3.14), we get for h sufficiently small and any $u_{h,k}^n \in \mathcal{P}_0(\mathcal{T}_h)$ that

$$\sum_{K \in \mathcal{T}_h} \mathcal{L}_h(u_{h,k}^n) u_K^n \geq \alpha \left(|u_{h,k}^n|_{1,\mathcal{T}_h}^2 + \|u_{h,k}^n\|_{L^2(\Omega)}^2 \right),$$

with $\alpha = \min(\bar{\lambda}_{\min}(1-\gamma)\frac{k}{4}, 1)$.

THEOREM 3.4. *For h sufficiently small, there exists a unique solution $u_{h,k}^n$ given by the scheme (2.10)–(2.11) at any discrete time step t_n .*

Proof. Assume that u_K , $K \in \mathcal{T}_h$, satisfy the linear system (3.13), and let the right-hand side of (3.13) be equal to 0. Then

$$(3.21) \quad \alpha \left(|u_{h,k}^n|_{1,\mathcal{T}_h}^2 + \|u_{h,k}^n\|_{L^2(\Omega)}^2 \right) \leq \sum_{K \in \mathcal{T}_h} \mathcal{L}_h(u_{h,k}^n) u_K^n = \sum_{K \in \mathcal{T}_h} f_{h,k}(u_{h,k}^{n-1}) u_K^n = 0.$$

Due to relation (3.21) and the strict positivity of α we know that $u_K^n = 0$ for all $K \in \mathcal{T}_h$. It means that the kernel of the linear transformation represented by the matrix of the system (3.13) contains only $\bar{0}$ vector, which implies that the matrix is regular. Thus it also implies that there exists a unique solution for any right-hand side. \square

4. Convergence of the scheme to the weak solution.

DEFINITION 4.1. *The weak solution of the problem (1.1)–(1.3) is a function $u \in L^2(0, T; H^1(\Omega))$ which satisfies the identity*

$$(4.1) \quad \int_0^T \int_\Omega u \frac{\partial \varphi}{\partial t}(x, t) dx dt + \int_\Omega u_0(x) \varphi(x, 0) dx - \int_0^T \int_\Omega (D \nabla u) \cdot \nabla \varphi dx dt = 0 \quad \forall \varphi \in \Psi,$$

where $\Psi = \{\varphi \in C^{2,1}(\bar{\Omega} \times [0, T]), (D \nabla \varphi) \cdot \vec{n} = 0 \text{ on } \partial\Omega \times (0, T), \varphi(\cdot, T) = 0\}$.

Remark 1. The existence and uniqueness of the weak solution and extremum principle for the model (1.1)–(1.3) are given in [24]. The proofs are based on theory built in [1].

In the proof of convergence we will use a strategy based on the application of Kolmogorov's compactness criterion in L^2 which gives the relative compactness of the approximate solutions given by the scheme refining the space and time discretization step. By using relative compactness we can choose a convergent subsequence which in the limit gives the weak solution. In order to use Kolmogorov's compactness criterion we shall prove the following four lemmas.

LEMMA 4.2 (uniform boundedness). *There exists a positive constant C such that*

$$(4.2) \quad \|u_{h,k}\|_{L^2(Q_T)} \leq C.$$

LEMMA 4.3 (time translate estimate). *For any $s \in (0, T)$ there exists a positive constant C such that*

$$(4.3) \quad \int_{\Omega \times (0, T-s)} (u_{h,k}(x, t+s) - u_{h,k}(x, t))^2 dx dt \leq Cs.$$

LEMMA 4.4 (space translate estimate I). *There exists a positive constant C such that*

$$(4.4) \quad \int_{\Omega_\xi \times (0, T)} (u_{h,k}(x + \xi, t) - u_{h,k}(x, t))^2 dx dt \leq C |\xi| (|\xi| + 2h)$$

for any vector $\xi \in R^d$, where $\Omega_\xi = \{x \in \Omega, [x, x + \xi] \in \Omega\}$.

LEMMA 4.5 (space translate estimate II). *There exists a positive constant C such that*

$$(4.5) \quad \int_{\Omega \times (0, T)} (u_{h,k}(x + \xi, t) - u_{h,k}(x, t))^2 dx dt \leq C |\xi|$$

for any vector $\xi \in R^d$.

To prove (4.2)–(4.5) we will use the following a priori estimates.

LEMMA 4.6. *The scheme (2.10)–(2.11) leads to the following estimates. For h sufficiently small, there exists a positive constant C which does not depend on h, k such that*

$$(4.6) \quad \max_{0 \leq n \leq N_{\max}} \sum_{K \in \mathcal{T}_h} (u_K^n)^2 m(K) \leq C,$$

$$(4.7) \quad \sum_{n=1}^{N_{\max}} k \sum_{(K,L) \in \Upsilon} \frac{(u_K^n - u_L^n)^2}{d_{K,L}} m(\sigma) \leq C,$$

$$(4.8) \quad \sum_{n=1}^{N_{\max}} \sum_{K \in \mathcal{T}_h} (u_K^n - u_K^{n-1})^2 m(K) \leq C.$$

Proof. We multiply (2.10) by u_K^n , sum it over $K \in \mathcal{T}_h$ and over $n = 1, \dots, m < N_{\max}$, and use the property $(a - b)a = \frac{1}{2}a^2 - \frac{1}{2}b^2 + \frac{1}{2}(a - b)^2$ to obtain

$$(4.9) \quad \begin{aligned} & \frac{1}{2} \sum_{K \in \mathcal{T}_h} (u_K^m)^2 m(K) + \frac{1}{2} \sum_{n=1}^m \sum_{K \in \mathcal{T}_h} (u_K^n - u_K^{n-1})^2 m(K) \\ & - \sum_{n=1}^m k \sum_{K \in \mathcal{T}_h} u_K^n \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{int}} \phi_\sigma^n(u_{h,k}^n) m(\sigma) = \frac{1}{2} \sum_{K \in \mathcal{T}_h} (u_K^0)^2 m(K). \end{aligned}$$

Then by using (3.16) and (3.20) we have

$$(4.10) \quad \begin{aligned} & \frac{1}{2} \sum_{K \in \mathcal{T}_h} (u_K^m)^2 m(K) + \frac{1}{2} \sum_{n=1}^m \sum_{K \in \mathcal{T}_h} (u_K^n - u_K^{n-1})^2 m(K) \\ & + \bar{\alpha} \sum_{n=1}^m k |u_{h,k}^n|_{1, \mathcal{T}_h}^2 \leq \frac{1}{2} \sum_{K \in \mathcal{T}_h} (u_K^0)^2 m(K), \end{aligned}$$

with positive constant $\bar{\alpha} = \bar{\lambda}_{\min} \frac{1-\gamma}{4}$. Since $u_0 \in L^2(\Omega)$, the right-hand side is bounded by a positive constant C . By using the first term of (4.10) we get a priori estimate (4.6), and from the second term of (4.10) we get a priori estimate (4.8). From the strict positiveness of $\bar{\alpha}$ in the third term of (4.10) and from definition (3.11), we obtain a priori estimate (4.7). \square

Proof of Lemma 4.2. It follows from the first $L^2(\Omega)$ —a priori estimate (4.6). \square

Proof of Lemma 4.3. First, for fixed $s \in (0, T)$, we define function

$$f(t) = \int_{\Omega} (u_{h,k}(x, t + s) - u_{h,k}(x, t))^2 dx.$$

By using the fact that $u_{h,k}$ is a piecewise constant function, we get

$$(4.11) \quad f(t) = \sum_{K \in \mathcal{T}_h} (u_K^{n_{t+s}} - u_K^{n_t})^2 m(K),$$

with $n_t = \lceil \frac{t}{k} \rceil$ and $n_{t+s} = \lceil \frac{t+s}{k} \rceil$, where $\lceil \cdot \rceil$ means the upper integer part of a positive real number. We rearrange (4.11) to obtain

$$(4.12) \quad f(t) = \sum_{K \in \mathcal{T}_h} (u_K^{n_{t+s}} - u_K^{n_t}) \sum_{t \leq (n+1)k < t+s} (u_K^n - u_K^{n-1}) m(K).$$

By using the scheme (2.10)–(2.11) in (4.12) (replacing K by W) we get

$$(4.13) \quad f(t) = \sum_{t \leq (n+1)k < t+s} k \sum_{W \in \mathcal{T}_h} \left((u_W^{n_{t+s}} - u_W^{n_t}) \sum_{\sigma \in \mathcal{E}_W \cap \mathcal{E}_{int}} \bar{\lambda}_\sigma (u_E^n - u_W^n) + \bar{\beta}_\sigma (u_N^n - u_S^n) \right),$$

and due to conservativity of numerical fluxes (antisymmetry of term $\bar{\lambda}_\sigma (u_E^n - u_W^n) + \bar{\beta}_\sigma (u_N^n - u_S^n)$) we have

$$(4.14) \quad f(t) = \sum_{t \leq (n+1)k < t+s} \frac{k}{2} \sum_{\sigma \in \mathcal{E}_{int}} (u_W^{n_{t+s}} - u_W^{n_t} - u_E^{n_{t+s}} + u_E^{n_t}) (\bar{\lambda}_\sigma (u_E^n - u_W^n) + \bar{\beta}_\sigma (u_N^n - u_S^n)).$$

Using Young's inequality leads to the relation

$$(4.15) \quad f(t) \leq \sum_{t \leq (n+1)k < t+s} \frac{k}{4} \sum_{\sigma \in \mathcal{E}_{int}} \bar{\lambda}_\sigma (u_W^{n_{t+s}} - u_W^{n_t} - u_E^{n_{t+s}} + u_E^{n_t})^2 + \sum_{t \leq (n+1)k < t+s} \frac{k}{4} \sum_{\sigma \in \mathcal{E}_{int}} \bar{\lambda}_\sigma \left((u_E^n - u_W^n) + \frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma} (u_N^n - u_S^n) \right)^2,$$

where the right-hand side can be further estimated and we get

$$(4.16) \quad f(t) \leq f_1(t) + f_2(t) + f_3(t) + f_4(t),$$

$$(4.17) \quad f_1(t) = \sum_{t \leq (n+1)k < t+s} \frac{k}{2} \sum_{\sigma \in \mathcal{E}_{int}} \bar{\lambda}_\sigma (u_E^{n_t} - u_W^{n_t})^2,$$

$$(4.18) \quad f_2(t) = \sum_{t \leq (n+1)k < t+s} \frac{k}{2} \sum_{\sigma \in \mathcal{E}_{int}} \bar{\lambda}_\sigma (u_E^{n_{t+s}} - u_W^{n_{t+s}})^2,$$

$$(4.19) \quad f_3(t) = \sum_{t \leq (n+1)k < t+s} \frac{k}{2} \sum_{\sigma \in \mathcal{E}_{int}} \bar{\lambda}_\sigma (u_E^n - u_W^n)^2,$$

$$(4.20) \quad f_4(t) = \sum_{t \leq (n+1)k < t+s} \frac{k}{2} \sum_{\sigma \in \mathcal{E}_{int}} \bar{\lambda}_\sigma \left(\frac{\bar{\beta}_\sigma}{\bar{\lambda}_\sigma} \right)^2 (u_N^n - u_S^n)^2.$$

Next we integrate (4.16) in time interval $(0, T-s)$, and by replacing $\sum_{\sigma \in \mathcal{E}_{int}}$ by $\sum_{(K,L) \in \Upsilon}$ (the edge $\sigma \in \mathcal{E}_K$ is an intersection of K and its adjacent finite volume L) we get an estimate of the first integral term

$$(4.21) \quad \int_0^{T-s} f_1(t) dt = \int_0^{T-s} \frac{k}{2} \sum_{(K,L) \in \Upsilon} \bar{\lambda}_\sigma (u_L^{n_t} - u_K^{n_t})^2 \sum_{n \in \mathbf{N}} \chi_{\{t \leq (n+1)k < t+s\}} dt.$$

We substitute the integral over $(0, T - s)$ by the sum of time step intervals and use the property $\chi_{\{t \leq (n+1)k < t+s\}} = \chi_{\{(n+1)k-s < t \leq (n+1)k\}}$ to obtain

$$(4.22) \quad \int_0^{T-s} f_1(t) dt \leq \sum_{n_t=0}^{N_{\max}-1} \frac{k}{2} \sum_{(K,L) \in \Upsilon} \bar{\lambda}_\sigma (u_L^{n_t} - u_K^{n_t})^2 \int_{n_t k}^{(n_t+1)k} \sum_{n \in \mathbf{N}^0} \chi_{\{(n+1)k-s < t \leq (n+1)k\}} dt.$$

Since $\int_{n_t k}^{(n_t+1)k} \sum_{n \in \mathbf{N}^0} \chi_{\{(n+1)k-s < t \leq (n+1)k\}} dt = s$, and $m(\sigma) = d_{K,L}$ for our uniform mesh, the relation (4.22) leads to

$$(4.23) \quad \int_0^{T-s} f_1(t) dt \leq s \sum_{n_t=0}^{N_{\max}} \frac{k}{2} \sum_{(K,L) \in \Upsilon} \frac{m(\sigma)}{d_{K,L}} \bar{\lambda}_\sigma (u_L^{n_t} - u_K^{n_t})^2.$$

The next step is to prove the following relation:

$$(4.24) \quad 0 < C_1 \leq \bar{\lambda}_\sigma \leq C_2 < \infty \text{ for all } \sigma \in \mathcal{E}.$$

Let K be any fixed finite volume. Since at any time step the matrix $D_\sigma = \begin{pmatrix} \bar{\lambda}_\sigma & \bar{\beta}_\sigma \\ \bar{\beta}_\sigma & \bar{\nu}_\sigma \end{pmatrix}$ is uniformly (strictly) positive definite, $\bar{\lambda}_\sigma \geq C_3 > 0$ and $\bar{\nu}_\sigma \geq C_4 > 0$ for all σ . The structure tensor evaluated numerically at point x_{KL} is given by

$$(4.25) \quad J_\rho \left(\nabla u_{h,k}^{n-1} \right)_{\bar{t}} (x_{KL}) = G_\rho * \begin{pmatrix} A & B \\ B & C \end{pmatrix},$$

where

$$(4.26) \quad A = \left(\left(\frac{\partial G_{\bar{t}}}{\partial x} * \tilde{u}_{h,k}^{n-1} \right) (x_{KL}) \right)^2, \quad C = \left(\left(\frac{\partial G_{\bar{t}}}{\partial y} * \tilde{u}_{h,k}^{n-1} \right) (x_{KL}) \right)^2,$$

$$(4.27) \quad B = \left(\frac{\partial G_{\bar{t}}}{\partial x} * \tilde{u}_{h,k}^{n-1} \right) (x_{KL}) \left(\frac{\partial G_{\bar{t}}}{\partial y} * \tilde{u}_{h,k}^{n-1} \right) (x_{KL}).$$

By using Young's inequality, a priori estimate (4.6), and the definition of extension $\tilde{u}_{h,k}^n$ (see (2.4)), we subsequently get for $i = 1, 2$ ($x_1 = x, x_2 = y$)

$$(4.28) \quad \begin{aligned} & \left| \left(\frac{\partial}{\partial x_i} G_{\bar{t}} * \tilde{u}_{h,k}^n \right) (x_{KL}) \right| \leq \int_{R^d} \left| \frac{\partial}{\partial x_i} G_{\bar{t}}(x_{KL} - \xi) \tilde{u}_{h,k}^n(\xi) \right| d\xi \\ & \leq \frac{1}{2} \int_{R^d} \left| \frac{\partial}{\partial x_i} G_{\bar{t}}(x_{KL} - \xi) \right|^2 d\xi + \frac{1}{2} \int_{R^d} |\tilde{u}_{h,k}^n(\xi)|^2 d\xi \leq C_{\bar{t}} \\ & + C_5 \int_{\Omega_{\bar{t}}} |\tilde{u}_{h,k}^n(\xi)|^2 d\xi \leq C_{\bar{t}} + C_5 \sum_{K \in \mathcal{T}_h} (u_K^n)^2 m(K) \leq C_6. \end{aligned}$$

By inspecting relations (1.6)–(1.10) we may observe that if the elements of matrix J_ρ are finite, then also the elements of matrix D_σ are finite, which gives (4.24). By applying (4.24) and (4.7) in (4.23), we get $\int_0^{T-s} f_1(t) dt \leq Cs$. By using similar approach as in [15] and relation (3.10), all further integrals can be estimated in the same way, which ends the proof. \square

Proof of Lemma 4.4. Let us define $\xi_{K,L} = \frac{\xi}{|\xi|} \cdot n_{K,\sigma}$ for all $(K, L) \in \Upsilon$ and let for all $x \in \Omega_\xi$

$$E(x, K, L) = \begin{cases} 1 & \text{if } [x, x + \xi] \text{ intersects } \sigma = \sigma_{KL}, K \text{ and } L; \text{ and } \xi_{K,L} > 0 \\ 0 & \text{otherwise.} \end{cases}$$

For any $t \in (0, T)$ there exists $n \in N$ which satisfies $(n-1)k < t \leq nk$. Then for almost all $x \in \Omega_\xi$ we can see that

$$u_{h,k}(x+\xi, t) - u_{h,k}(x, t) = u_{K(x+\xi)}^n - u_{K(x)}^n = \sum_{(K,L) \in \Upsilon} E(x, K, L) (u_L^n - u_K^n),$$

where $K(x)$ denotes the volume $K \in \mathcal{T}_h$, where $x \in K$. By using these notations we get the proof in similar lines as in [15]—proof of Lemma 3.2. \square

Proof of Lemma 4.5. In this proof, for simplicity, let us consider that $\Omega_{\tilde{t}} = \Omega$; i.e., we extend $u_{h,k}$ outside Ω by 0. The results which are obtained below can be straightforwardly adjusted to the situation with reflexion to $\Omega_{\tilde{t}}$; the derivation is just technically more complicated, and for details we refer to [10]. Let us define the set

$$\mathcal{E}_{ext} = \{\varpi, \text{ such that there exists } K \in \mathcal{T}_h, \varpi \subset \partial K \cap \partial\Omega\}$$

and let $u_\varpi := u_K$, where $K \in \mathcal{T}_h$, $\varpi \subset \partial K \cap \partial\Omega$. Since now for $x \in \Omega - \Omega_\xi$ the point $x + \xi$ can be outside Ω , we see that

$$(4.29) \quad \begin{aligned} & u_{h,k}(x+\xi, t) - u_{h,k}(x, t) \\ &= \sum_{(K,L) \in \Upsilon} E(x, K, L) (u_L^n - u_K^n) - \sum_{\varpi \in \mathcal{E}_{ext}} \chi([x, x+\xi] \cap \varpi) u_\varpi^n. \end{aligned}$$

By using the Cauchy–Schwarz and Young inequalities we obtain

$$(4.30) \quad \begin{aligned} & (u_{h,k}(x+\xi, t) - u_{h,k}(x, t))^2 \\ & \leq 2 \left(\sum_{(K,L) \in \Upsilon} E(x, p, q) \xi_{K,L} d_{K,L} \right) \left(\sum_{(K,L) \in \Upsilon} E(x, K, L) \frac{(u_L^n - u_K^n)^2}{\xi_{K,L} d_{K,L}} \right) \\ & \quad + 2 \sum_{\varpi \in \mathcal{E}_{ext}} \chi([x, x+\xi] \cap \varpi) (u_\varpi^n)^2, \end{aligned}$$

$$(4.31) \quad \begin{aligned} & \int_{\Omega \times (0, T)} (u_{h,k}(x+\xi, t) - u_{h,k}(x, t))^2 dx dt \\ & \leq (2h + |\xi|) |\xi| C + 2 \sum_{n=0}^{N_{\max}} k \int_{\Omega} \sum_{\varpi \in \mathcal{E}_{ext}} \chi([x, x+\xi] \cap \varpi) (u_\varpi^n)^2 dx dt, \end{aligned}$$

which can be written as

$$(4.32) \quad \begin{aligned} & \int_{\Omega \times (0, T)} (u_{h,k}(x+\xi, t) - u_{h,k}(x, t))^2 dx dt \\ & \leq (2h + |\xi|) |\xi| C + 2 |\xi| \sum_{n=0}^{N_{\max}} k \sum_{\varpi \in \mathcal{E}_{ext}} (u_\varpi^n)^2 m(\varpi). \end{aligned}$$

For the last term in (4.32) we use the trace inequality given in [4].

LEMMA 4.7. *Let Ω be an open bounded polygonal connected subset of R^d . Let $\bar{\gamma}(u_{h,k})$ be defined by $\bar{\gamma}(u_{h,k}) = u_\varpi$ a.e. for the $(d-1)$ -Lebesgue measure on $\varpi \in \mathcal{E}_{ext}$. Then there exists positive C depending only on Ω such that*

$$\|\bar{\gamma}(u_{h,k}^n)\|_{L^2(\partial\Omega)} \leq C \left(\|u_{h,k}^n\|_{1, \mathcal{T}_h} + \|u_{h,k}^n\|_{L^2(\Omega)} \right).$$

By using the trace operator $\bar{\gamma}(u_{h,k}) = u_{\varpi}$ we can write

$$(4.33) \quad \begin{aligned} & \int_{\Omega \times (0,T)} (u_{h,k}(x+\xi, t) - u_{h,k}(x, t))^2 dx dt \\ & \leq (2h + |\xi|) |\xi| C + 2 |\xi| \sum_{n=0}^{N_{\max}} k \|\bar{\gamma}(u_{h,k}^n)\|_{L^2(\partial\Omega)}^2 \end{aligned}$$

and applying the trace inequality implies that

$$(4.34) \quad \begin{aligned} & \int_{\Omega \times (0,T)} (u_{h,k}(x+\xi, t) - u_{h,k}(x, t))^2 dx dt \\ & \leq (2h + |\xi|) |\xi| C + 4C |\xi| \sum_{n=0}^{N_{\max}} k \left(|u_{h,k}^n|_{1, \mathcal{T}_h}^2 + \|u_{h,k}^n\|_{L^2(\Omega)}^2 \right). \end{aligned}$$

Then using a priori estimates (4.6) and (4.7) ends the proof. \square

Lemmas 4.2, 4.5, and 4.3 guarantee that sequence $u_{h,k}$ is relatively compact in $L^2(Q_T)$, which implies following convergence result.

LEMMA 4.8. *There exists $u \in L^2(Q_T)$ such that $u_{h,k} \rightarrow u$ in $L^2(Q_T)$ as $h, k \rightarrow 0$ in the sense of subsequences.*

For the sake of simplicity, we denote the subsequence converging to u again by $u_{h,k}$, and we are going to prove that its limit u is the weak solution of (1.1)–(1.3) in the sense of Definition 4.1.

To that goal, let $\varphi \in \Psi$ be given, and multiply the scheme (2.10) by $\varphi(x_K, t_n)$. Then we sum it over all $K \in \mathcal{T}_h$ and for $n = 1, \dots, N_{\max}$ to get

$$(4.35) \quad \begin{aligned} & \sum_{n=1}^{N_{\max}} k \sum_{K \in \mathcal{T}_h} \frac{(u_K^n - u_K^{n-1})}{k} \varphi(x_K, t_{n-1}) m(K) \\ & = \sum_{n=1}^{N_{\max}} k \sum_{K \in \mathcal{T}_h} \varphi(x_K, t_{n-1}) \sum_{\sigma \in \mathcal{E}_K \cap \mathcal{E}_{int}} \phi_{\sigma_{KL}}^n(u_{h,k}^n) m(\sigma). \end{aligned}$$

In order to have a structure similar to the weak solution identity (4.1), we rearrange (4.35) by using a discrete integration by parts and gathering the sums over $K \in \mathcal{T}_h$ and $\sigma \in \mathcal{E}_K \cap \mathcal{E}_{int}$, and we get

$$(4.36) \quad \begin{aligned} & \sum_{n=1}^{N_{\max}} k \sum_{K \in \mathcal{T}_h} u_K^n \frac{\varphi(x_K, t_n) - \varphi(x_K, t_{n-1})}{k} m(K) + \sum_{K \in \mathcal{T}_h} u_K^0 \varphi(x_K, 0) m(K) \\ & - \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{\sigma \in \mathcal{E}_{int}} \phi_{\sigma}^n(u_{h,k}^n) m(\sigma) (\varphi(x_L, t_{n-1}) - \varphi(x_K, t_{n-1})) = 0. \end{aligned}$$

In the same way as in [15] we can prove that

$$(4.37) \quad \sum_{n=1}^{N_{\max}} k \sum_{K \in \mathcal{T}_h} u_K^n \frac{\varphi(x_K, t_n) - \varphi(x_K, t_{n-1})}{k} m(K) \rightarrow \int_0^T \int_{\Omega} u \frac{\partial \varphi}{\partial t}(x, t) dx dt,$$

$$(4.38) \quad \sum_{K \in \mathcal{T}_h} u_K^0 \varphi(x_K, 0) m(K) \rightarrow \int_{\Omega} u_0(x) \varphi(x, 0) dx$$

as $h, k \rightarrow 0$ for all $\varphi \in \Psi$. The main point in proving convergence of the scheme is to get that

$$(4.39) \quad \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{\sigma \in \mathcal{E}_{int}} \phi_{\sigma}^n(u_{h,k}^n) m(\sigma) (\varphi(x_L, t_{n-1}) - \varphi(x_K, t_{n-1})) \\ \rightarrow \int_0^T \int_{\Omega} -\nabla \cdot (D\nabla \varphi) u dx dt$$

as $h, k \rightarrow 0$ for all $\varphi \in \Psi$. By using then the space translate estimate (4.4) we know (see, e.g., [4] or [14]) that the limit function u is in the space $L^2((0, T), H^1)$, so we can use Green's theorem, and by applying the boundary conditions we have

$$\int_0^T \int_{\Omega} -\nabla \cdot (D\nabla \varphi) u dx dt = \int_0^T \int_{\Omega} (D\nabla \varphi) \cdot \nabla u dx dt.$$

Proving (4.39) thus leads to overall convergence of the scheme to the weak solution in the sense of (4.1). To deal with (4.39) we rewrite it and then split into the sum of five terms

$$(4.40) \quad \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} \phi_{\sigma}^n(u_{h,k}^n) m(\sigma) (\varphi(x_E, t_{n-1}) - \varphi(x_W, t_{n-1})) \\ + \int_0^T \int_{\Omega} \nabla \cdot (D\nabla \varphi(x, t)) u(x, t) dx dt = \sum_{i=1}^5 T_i,$$

where

$$T_1 = \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} [(D_{\sigma} p_{\sigma}(u)) \cdot \mathbf{n}_{W,\sigma} (\varphi_E^{n-1} - \varphi_W^{n-1}) \\ - (u_E^n - u_W^n) (D_{\sigma} p_{\sigma}(\varphi^{n-1})) \cdot \mathbf{n}_{W,\sigma}] m(\sigma), \\ T_2 = \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) m(\sigma) [(D_{\sigma} p_{\sigma}(\varphi^{n-1})) \\ \cdot \mathbf{n}_{W,\sigma} - (D_{\sigma} \nabla \varphi(x_{WE}, t_{n-1})) \cdot \mathbf{n}_{W,\sigma}], \\ T_3 = \frac{1}{2} \sum_{n=1}^{N_{\max}} \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) \left(m(\sigma) k (D_{\sigma} \nabla \varphi(x_{WE}, t_{n-1})) \right. \\ \left. \cdot \mathbf{n}_{W,\sigma} - \int_{t_{n-1}}^{t_n} \int_{\sigma} (D_{\sigma} \nabla \varphi(s, t)) \cdot \mathbf{n}_{W,\sigma} ds dt \right), \\ T_4 = \frac{1}{2} \sum_{n=1}^{N_{\max}} \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) \int_{t_{n-1}}^{t_n} \int_{\sigma} ((D_{\sigma} - D) \nabla \varphi(s, t)) \cdot \mathbf{n}_{W,\sigma} ds dt, \\ T_5 = \int_0^T \int_{\Omega} \nabla \cdot (D\nabla \varphi(x, t)) (u(x, t) - u_{h,k}(x, t)) dx dt,$$

where $\varphi_W^{n-1} = \varphi(x_W, t_{n-1})$, $\varphi_E^{n-1} = \varphi(x_E, t_{n-1})$, and $\varphi^{n-1} = \varphi(x, t_{n-1})$. Since

$$\begin{aligned} & \frac{1}{2} \sum_{n=1}^{N_{\max}} \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) \int_{t_{n-1}}^{t_n} \int_{\sigma} (D\nabla \varphi(s, t)) \cdot \mathbf{n}_{W,\sigma} ds dt \\ &= - \sum_{n=1}^{N_{\max}} \sum_{W \in \mathcal{T}_h} u_W^n \int_{t_{n-1}}^{t_n} \sum_{\sigma \in \mathcal{E}_W} \int_{\sigma} (D\nabla \varphi(s, t)) \cdot \mathbf{n}_{W,\sigma} ds dt \\ &= - \sum_{n=1}^{N_{\max}} \sum_{W \in \mathcal{T}_h} u_W^n \int_{t_{n-1}}^{t_n} \int_W \nabla \cdot (D\nabla \varphi(x, t)) dx dt \\ &= - \int_0^T \int_{\Omega} \nabla \cdot (D\nabla \varphi(x, t)) u_{h,k}(x, t) dx dt, \end{aligned}$$

one can see correspondence of terms in T_4 and T_5 . First, let us deal with T_1 and rewrite it using (2.7)–(2.8) into the form

$$\begin{aligned} T_1 = & \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} ([\bar{\lambda}_{\sigma}(u_E^n - u_W^n) + \bar{\beta}_{\sigma}(u_N^n - u_S^n)](\varphi_E^{n-1} - \varphi_W^{n-1}) \\ & - (u_E^n - u_W^n)[\bar{\lambda}_{\sigma}(\varphi_E^{n-1} - \varphi_W^{n-1}) + \bar{\beta}_{\sigma}(\varphi_N^{n-1} - \varphi_S^{n-1})]), \end{aligned}$$

which can be easily simplified to

$$(4.41) \quad T_1 = \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} [\bar{\beta}_{\sigma}(u_N^n - u_S^n)(\varphi_E^{n-1} - \varphi_W^{n-1}) - \bar{\beta}_{\sigma}(u_E^n - u_W^n)(\varphi_N^{n-1} - \varphi_S^{n-1})].$$

By applying (3.1) for u and similarly for φ we get

$$\begin{aligned} (4.42) \quad T_1 = & \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} \sum_{i=1}^4 \left[\frac{\bar{\beta}_{\sigma}}{4} (\varphi_{E_{\sigma}} - \varphi_{W_{\sigma}})(u_{E_{\delta_i}} - u_{W_{\delta_i}}) \right. \\ & \left. - \frac{\bar{\beta}_{\sigma}}{4} (u_{E_{\sigma}} - u_{W_{\sigma}})(\varphi_{E_{\delta_i}} - \varphi_{W_{\delta_i}})(1 + O(h)) \right], \end{aligned}$$

where we omit time indexes due to simplification; for graphical explanation of notations see Figures 3.1 and 4.1. For each term with positive sign in (4.42) one can find a corresponding term in the group of terms with negative signs ($\bar{\beta}_{\sigma}$ corresponds to some $\bar{\beta}_{\delta}$). We denote these couples by $T_{\sigma\delta}$. For example, for $\sigma = \sigma_{WE}$ and δ as plotted in Figure 4.1 (left) we can write the couple as follows:

$$\begin{aligned} T_{\sigma_{WE}\delta} &= \frac{\bar{\beta}_{\sigma_{WE}}}{4} (\varphi_{E_{\sigma_{WE}}}^{n-1} - \varphi_{W_{\sigma_{WE}}}^{n-1})(u_{E_{\delta}}^n - u_{W_{\delta}}^n) \\ &\quad - \frac{\bar{\beta}_{\delta}}{4} (\varphi_{E_{\sigma_{EW}}}^{n-1} - \varphi_{W_{\sigma_{EW}}}^{n-1})(u_{E_{\delta}}^n - u_{W_{\delta}}^n)(1 + O(h)) \\ &= \frac{\bar{\beta}_{\sigma_{WE}}}{4} (\varphi_{E_{\sigma_{WE}}}^{n-1} - \varphi_{W_{\sigma_{WE}}}^{n-1})(u_{E_{\delta}}^n - u_{W_{\delta}}^n) - \frac{\bar{\beta}_{\delta}}{4} (\varphi_{W_{\sigma_{WE}}}^{n-1} - \varphi_{E_{\sigma_{WE}}}^{n-1})(u_{E_{\delta}}^n - u_{W_{\delta}}^n)(1 + O(h)) \end{aligned}$$

because $\bar{\beta}_{\sigma_{WE}} = \bar{\beta}_{\sigma_{EW}}$, $\varphi_{E_{\sigma_{WE}}}^{n-1} = \varphi_{W_{\sigma_{EW}}}^{n-1}$, and $\varphi_{W_{\sigma_{WE}}}^{n-1} = \varphi_{E_{\sigma_{EW}}}^{n-1}$; see Figure 4.1 (right). Using the previous expression for every $\sigma = \sigma_{WE}$ yields

$$\begin{aligned} (4.43) \quad T_{\sigma\delta} &= \left[\frac{\bar{\beta}_{\sigma}}{4} + \left(-\frac{\bar{\beta}_{\sigma}}{4} (1 + O(h)) \right) \right] (\varphi_{E_{\sigma}}^{n-1} - \varphi_{W_{\sigma}}^{n-1})(u_{E_{\delta}}^n - u_{W_{\delta}}^n) \\ &= \frac{\bar{\beta}_{\sigma}}{4} O(h)(\varphi_{E_{\sigma}}^{n-1} - \varphi_{W_{\sigma}}^{n-1})(u_{E_{\delta}}^n - u_{W_{\delta}}^n). \end{aligned}$$

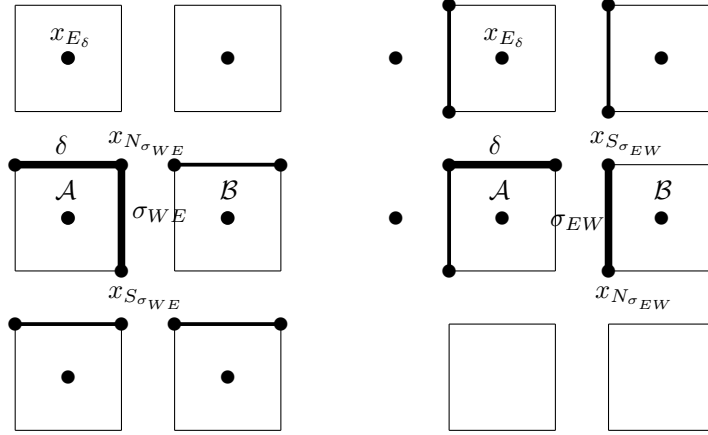


FIG. 4.1. Left: The edge $\sigma = \sigma_{WE}$ and one $\delta \in P_\sigma$. Right: The same edge δ and one corresponding $\sigma = \sigma_{EW} \in P_\delta$. The thickest lines represent one particular couple $T_{\sigma\delta}$ in (4.42), which vanishes up to the $O(h)$ term. $\mathcal{A} = x_{W_{\sigma_{WE}}} = x_{E_{\sigma_{EW}}} = x_{W_\delta}$, $\mathcal{B} = x_{E_{\sigma_{WE}}} = x_{W_{\sigma_{EW}}}$.

Then T_1 can be estimated as follows:

$$|T_1| \leq C_1 h \left| \sum_{n=1}^{N_{\max}} k \sum_{\sigma \in \mathcal{E}_{int}} \sum_{\delta \in P_\sigma \cap \mathcal{E}_{int}} (\varphi_{E_\sigma}^{n-1} - \varphi_{W_\sigma}^{n-1})(u_{E_\delta}^n - u_{W_\delta}^n) \right|,$$

with positive constant C_1 due to the fact that $\bar{\beta}_\sigma$ for each edge of mesh is finite (see (4.24)). By the Cauchy–Schwarz inequality we have

$$(4.44) \quad |T_1| \leq C_2 h \left(\sum_{n=1}^{N_{\max}} k \sum_{\sigma \in \mathcal{E}_{int}} \sum_{\delta \in P_\sigma \cap \mathcal{E}_{int}} (\varphi_{E_\sigma}^{n-1} - \varphi_{W_\sigma}^{n-1})^2 \right)^{\frac{1}{2}} \left(\sum_{n=1}^{N_{\max}} k \sum_{\sigma \in \mathcal{E}_{int}} \sum_{\delta \in P_\sigma \cap \mathcal{E}_{int}} (u_{E_\delta}^n - u_{W_\delta}^n)^2 \right)^{\frac{1}{2}}.$$

It comes from the regularity of φ that there exists a positive constant C_3 such that $(\varphi_{E_\sigma}^{n-1} - \varphi_{W_\sigma}^{n-1})^2 \leq C_3 h^2$. Thanks to geometrical arguments, we know that

$$(4.45) \quad \sum_{\sigma \in \mathcal{E}_{int}} d_{WE} m(\sigma) \leq C_4 |\Omega|,$$

which straightforwardly gives for our uniform square mesh $\sum_{\sigma \in \mathcal{E}_{int}} h^2 \leq C_4 |\Omega|$. The above-mentioned facts lead to

$$\left(\sum_{n=1}^{N_{\max}} k \sum_{\sigma \in \mathcal{E}_{int}} \sum_{\delta \in P_\sigma \cap \mathcal{E}_{int}} (\varphi_{E_\sigma}^{n-1} - \varphi_{W_\sigma}^{n-1})^2 \right) \leq C_5 T |\Omega|,$$

which together with a priori estimate (4.7) gives that $|T_1| \leq C_6 h$, which implies

$$(4.46) \quad |T_1| \rightarrow 0 \text{ as } h, k \rightarrow 0.$$

Term T_2 can be written as $T_2 = \frac{1}{2} \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) m(\sigma) T_{2WE}$, where

$$T_{2WE} = \bar{\lambda}_\sigma \frac{\varphi_E^{n-1} - \varphi_W^{n-1}}{h} + \bar{\beta}_\sigma \frac{\varphi_N^{n-1} - \varphi_S^{n-1}}{h} - \bar{\lambda}_\sigma (\varphi_{WE}^{n-1})_x - \bar{\beta}_\sigma (\varphi_{WE}^{n-1})_y$$

and $\left(\frac{\varphi_{WE}^{n-1}}{(\varphi_{WE}^{n-1})_y}\right)_x = \nabla \varphi(x_{WE}, t_{n-1})$ in the basis $(\mathbf{n}_{W,\sigma}, \mathbf{t}_{W,\sigma})$. Since $\varphi \in C^{2,1}(\bar{\Omega} \times [0, T])$, there exist positive constants C_7 and C_8 such that

$$\left| \frac{\varphi_E^{n-1} - \varphi_W^{n-1}}{h} - (\varphi_{WE}^{n-1})_x \right| \leq C_7 h, \quad \left| \frac{\varphi_N^{n-1} - \varphi_S^{n-1}}{h} - (\varphi_{WE}^{n-1})_y \right| \leq C_8 h.$$

From there and the property that elements of D_σ are finite (see (4.28)) we have $|T_{2WE}| \leq C_9 h$, with a positive constant C_9 , which implies

$$(4.47) \quad |T_2| \leq C_{10} h \left| \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) m(\sigma) \right|.$$

Then we multiply the right-hand side of (4.47) by $\frac{\sqrt{d_{WE}}}{\sqrt{d_{WE}}}$ and apply the Cauchy–Schwarz inequality to obtain

$$|T_2| \leq C_{10} h \left(\sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n)^2 \frac{m(\sigma)}{d_{WE}} \right)^{\frac{1}{2}} \left(\sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} m(\sigma) d_{WE} \right)^{\frac{1}{2}}.$$

A priori estimate (4.7) together with (4.45) gives $|T_2| \leq C_{10}(C_{11}C_4|\Omega|T)^{\frac{1}{2}}h$ and finally

$$(4.48) \quad |T_2| \rightarrow 0 \text{ as } h, k \rightarrow 0.$$

We consider the third term in the form $T_3 = \frac{1}{2} \sum_{n=1}^{N_{\max}} \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) T_{3WE}$, where

$$T_{3WE} = m(\sigma) k (D_\sigma \nabla \varphi(x_{WE}, t_{n-1})) \cdot \mathbf{n}_{W,\sigma} - \int_{t_{n-1}}^{t_n} \int_\sigma (D_\sigma \nabla \varphi(s, t)) \cdot \mathbf{n}_{W,\sigma} ds dt.$$

Due to the smoothness of φ , the mean value theorem, and the finiteness of D_σ , we have $|T_3| \leq C_{12}(h+k) \left| \sum_{n=1}^{N_{\max}} k \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) m(\sigma) \right|$, and similarly as above by using (4.7) and (4.45) we get

$$(4.49) \quad |T_3| \rightarrow 0 \text{ as } h, k \rightarrow 0.$$

Let us define the fourth term as $T_4 = \frac{1}{2} \sum_{n=1}^{N_{\max}} \sum_{(W,E) \in \Upsilon} (u_E^n - u_W^n) T_\sigma$, where T_σ is represented by the relation $T_\sigma = \int_{t_{n-1}}^{t_n} \int_\sigma (\bar{\lambda}_\sigma - \bar{\lambda}) \varphi_x + (\bar{\beta}_\sigma - \bar{\beta}) \varphi_y ds dt$ and $\bar{\lambda}(x, t)$ and $\bar{\beta}(x, t)$ are the elements of matrix $D(x, t)$ in the basis $(\mathbf{n}_{W,\sigma}, \mathbf{t}_{W,\sigma})$. Owing to the continuity of function $\nabla \varphi$, one can show that there exist positive constants C_{13} and C_{14} such that $|\varphi_x| \leq C_{13}$ and $|\varphi_y| \leq C_{14}$. By using these facts it will be sufficient to prove that

$$(4.50) \quad |\lambda_\sigma - \lambda| = |\lambda(x_{WE}, t_{n-1}) - \lambda(x, t)| \leq C_{15} h, \quad x \in \sigma, \quad t \in [t_{n-1}, t_n],$$

$$(4.51) \quad |\beta_\sigma - \beta| = |\beta(x_{WE}, t_{n-1}) - \beta(x, t)| \leq C_{16} h, \quad x \in \sigma, \quad t \in [t_{n-1}, t_n],$$

$$(4.52) \quad |\nu_\sigma - \nu| = |\nu(x_{WE}, t_{n-1}) - \nu(x, t)| \leq C_{17} h, \quad x \in \sigma, \quad t \in [t_{n-1}, t_n]$$

for each edge $\sigma \in \mathcal{E}$. It comes from (1.10) that

$$|\lambda_\sigma - \lambda| = \left| \frac{\kappa_1 v_{1\sigma}^2 + \kappa_{2\sigma} v_{2\sigma}^2}{v_{1\sigma}^2 + v_{2\sigma}^2} - \frac{\kappa_1 v_1^2 + \kappa_2 v_2^2}{v_1^2 + v_2^2} \right|,$$

where the terms with index σ correspond to λ_σ and the remaining terms to λ . Here κ_2 depends on μ_1 and μ_2 (see (1.9)), and μ_1 and μ_2 together with v_1 and v_2 (see (1.7) and (1.8)) depend on the elements of the structure tensor $J_\rho(\nabla u_{\tilde{t}}) = \begin{pmatrix} a & b \\ b & c \end{pmatrix}$ which is constructed by (1.5)–(1.6). Similarly, λ_σ depends on the elements of tensor $J_\rho(\nabla u_{h,k_{\tilde{t}}})(x_{KL}) = \begin{pmatrix} a_\sigma & b_\sigma \\ b_\sigma & c_\sigma \end{pmatrix}$, evaluated numerically at point x_{KL} , using values of $\tilde{u}_{h,k}$ and weights given by convolutions. Using (4.25)–(4.27) leads to

$$\begin{aligned} a_\sigma &= G_\rho * \left(\frac{\partial G_{\tilde{t}}}{\partial x} * u_{h,k} \right)^2, \quad c_\sigma = G_\rho * \left(\frac{\partial G_{\tilde{t}}}{\partial y} * u_{h,k} \right)^2, \\ b_\sigma &= G_\rho * \left[\left(\frac{\partial G_{\tilde{t}}}{\partial x} * u_{h,k} \right) \left(\frac{\partial G_{\tilde{t}}}{\partial y} * u_{h,k} \right) \right]. \end{aligned}$$

Elements of the matrix $J_\rho(\nabla u_{\tilde{t}})$ are given as follows:

$$a = G_\rho * \left(\frac{\partial G_{\tilde{t}}}{\partial x} * u \right)^2, \quad b = G_\rho * \left[\left(\frac{\partial G_{\tilde{t}}}{\partial x} * u \right) \left(\frac{\partial G_{\tilde{t}}}{\partial y} * u \right) \right], \quad c = G_\rho * \left(\frac{\partial G_{\tilde{t}}}{\partial y} * u \right)^2.$$

To prove (4.50), we use inequality $p^2 - q^2 \leq |p + q||p - q|$, and then we get $a_\sigma - a \leq G_\rho * (|\frac{\partial G_{\tilde{t}}}{\partial x} * (u_{h,k} + u)| |\frac{\partial G_{\tilde{t}}}{\partial x} * (u_{h,k} - u)|)$. Applying (4.28), the boundedness of solution u , and Lemma 4.8 yields that $a_\sigma \rightarrow a$ as $h, k \rightarrow 0$. One can obtain convergences $c_\sigma \rightarrow c$ in the same way, just changing x to y . For b_σ and b we add and subtract $G_\rho * [(\frac{\partial G_{\tilde{t}}}{\partial x} * u)(\frac{\partial G_{\tilde{t}}}{\partial y} * u_{h,k})]$. Then we have

$$b_\sigma - b = G_\rho * \left[\left(\frac{\partial G_{\tilde{t}}}{\partial y} * u_{h,k} \right) \left(\frac{\partial G_{\tilde{t}}}{\partial x} * (u_{h,k} - u) \right) + \left(\frac{\partial G_{\tilde{t}}}{\partial x} * u \right) \left(\frac{\partial G_{\tilde{t}}}{\partial y} * (u_{h,k} - u) \right) \right],$$

and using (4.28) and Lemma 4.8 yields the result. Due to dependence of v_1 , v_2 , μ_1 , μ_2 , and κ_2 on a , b , and c , one can conclude that also $v_{1\sigma} \rightarrow v_1$ as $h, k \rightarrow 0$ and $v_{2\sigma}$, $\mu_{1\sigma}$, $\mu_{2\sigma}$, and $\kappa_{2\sigma}$ correspondingly. Then we use above-mentioned facts to get $\lambda_\sigma \rightarrow \lambda$ as $h, k \rightarrow 0$. Since $|\beta_\sigma - \beta| = |\frac{v_{1\sigma} v_{2\sigma} (\kappa_1 - \kappa_{2\sigma})}{v_{1\sigma}^2 + v_{2\sigma}^2} - \frac{v_1 v_2 (\kappa_1 - \kappa_2)}{v_1^2 + v_2^2}|$ and $|\nu_\sigma - \nu| = |\frac{\kappa_{2\sigma} v_{1\sigma}^2 + \kappa_1 v_{2\sigma}^2}{v_{1\sigma}^2 + v_{2\sigma}^2} - \frac{\kappa_2 v_1^2 + \kappa_1 v_2^2}{v_1^2 + v_2^2}|$, one can prove the inequalities (4.51) and (4.52) in the same way. Then by using the same technique as in estimates of T_2 and T_3 we get

$$(4.53) \quad |T_4| \rightarrow 0 \text{ as } h, k \rightarrow 0.$$

The last term is given by $T_5 = \int_0^T \int_\Omega \nabla \cdot (D \nabla \varphi(x, t)) (u(x, t) - u_{h,k}(x, t)) dx dt$. We use the property that $D \in C^\infty(R^{2 \times 2}, R^{2 \times 2})$ (see [24, p. 115]) to state that $\nabla \cdot (D \nabla \varphi(x, t)) \in L^2(Q_T)$. Then by using the strong convergence of $u_{h,k}$ to u one can see that

$$(4.54) \quad |T_5| \rightarrow 0 \text{ as } h, k \rightarrow 0.$$

Now we can state following convergence result.

THEOREM 4.9. *The sequence $u_{h,k}$ converges strongly in $L^2(Q_T)$ to the unique weak solution u of (1.1)–(1.3) as $h, k \rightarrow 0$.*

Proof. The result comes from (4.46), (4.48), (4.49), (4.53), (4.54), and the fact that the limit u of $u_{h,k}$ is in space $L^2(0, T; H^1(\Omega))$. Due to the uniqueness of the weak solution, which can be found in [24], not only the subsequence but also the sequence $u_{h,k}$ itself converges to u . \square

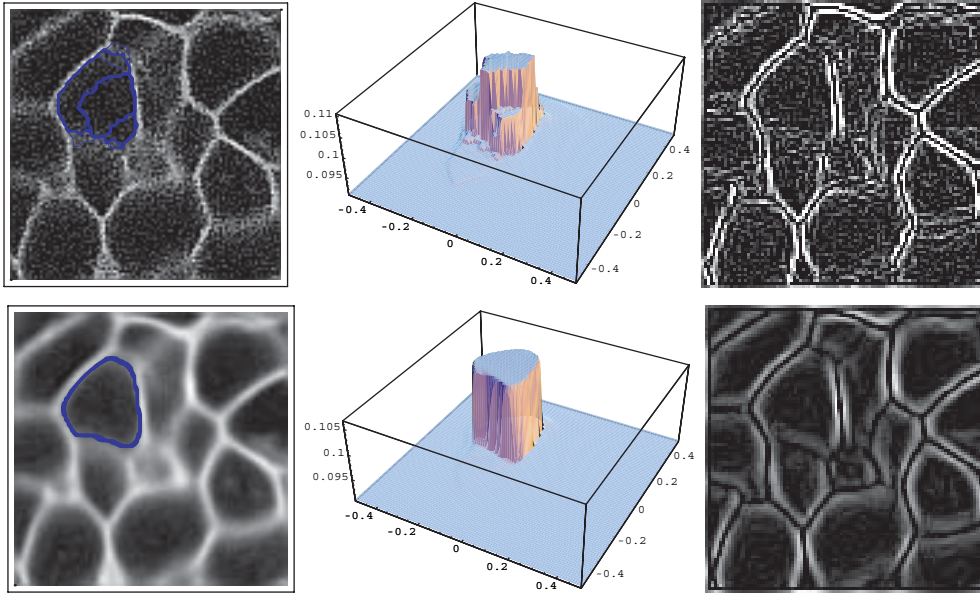


FIG. 5.1. The cell segmentation and edge detection for the original image and the image filtered by 4 diffusion steps. Top row: The original image (100×100 pixels detail) with the isolines of the final state of the segmentation function (left), the graph of the final state of the segmentation function (middle), and the edge detection for the original image (right). Bottom row: The filtered image with the isolines of the final state of the segmentation function (left), the graph of the final state of the segmentation function (middle), and the edge detection for the filtered image (right).

5. Numerical experiments. In this section we present results of several computational examples using real two-dimensional images coming from multiphoton laser scanning microscopy. They represent the membranes and nuclei of cells in the early stages of zebrafish embryogenesis. Especially the images of membranes are well suited for processing by this type of diffusion which is documented by comparing the edge detection and cell segmentation results before and after filtering in Figures 5.1–5.4. In the experiments we use the spatial step $h = 0.01$, time step $k = 0.0001$, $C = 1$, $\alpha = 0.001$, $\tilde{t} = 0.00001$, $\rho = 0.002$. The arising linear systems are solved using Gauss–Seidel iterations. Satisfactory results were obtained after a few filtering steps, so the denoising method is really fast. In the presented experiments we do not observe any stability problems which are a usual drawback of explicit schemes; cf. [25].

The nonlinear tensor anisotropic diffusion smoothes out the noise and improves significantly the connectivity of the coherent structures. Although the filtered image seems to be more blurred compared to the original one—cf. Figure 5.1 (left top and bottom)—the enhancement of the structure connectivity and improvement of the quality of edge detection—cf. Figure (5.1 right top and bottom)—enable us to get much more precise results of segmentation algorithms based on image intensity gradient information such as the subjective surface method [20, 16, 2]. In the subjective surface method, the initial segmentation function in the form of a peak centered in the approximate centroid of the segmented object is created. Then the initial function is evolved by nonlinear PDE, it forms a shock profile during the evolution, and the segmented object is detected as an isoline of the final (numerically steady) state of the segmentation function; for details we refer to [20, 16]. Since many spurious noisy structures can be seen in the original image, which is expressed in a highly noisy edge

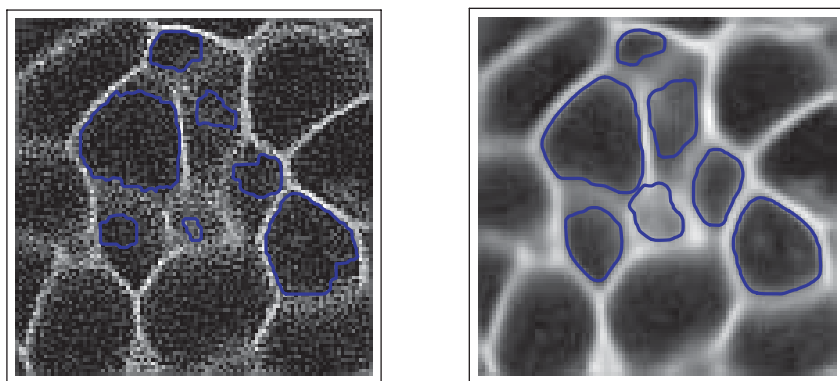


FIG. 5.2. The results of subjective surface cell segmentation when using unfiltered and filtered images, respectively.

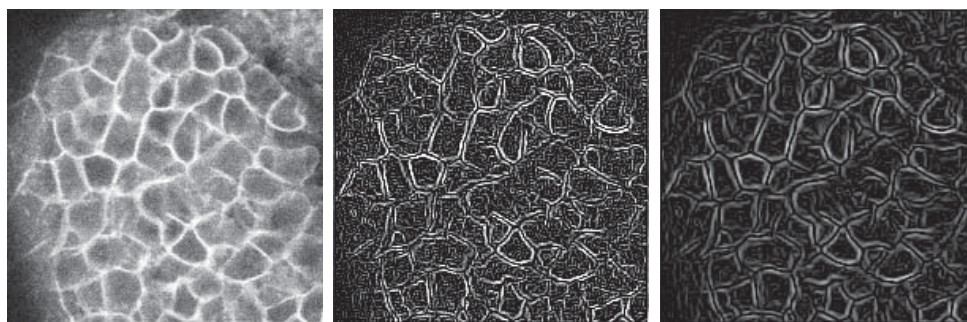


FIG. 5.3. The image of the cell membranes (200×200 pixels, left), its edge detection (middle), and the edge detection for the image filtered by 2 diffusion steps where the strong improvement of structure coherence can be seen (right).

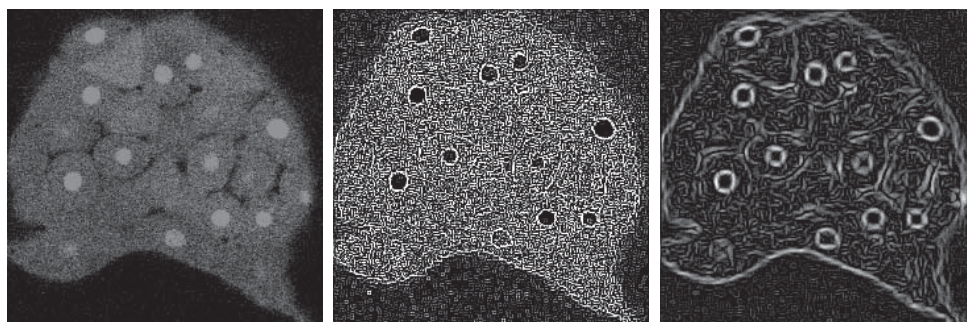


FIG. 5.4. The image of the cell membranes and nuclei (240×240 pixels, left), its edge detection (middle), and the edge detection for the image filtered by 5 diffusion steps (right).

detection result, the segmentation algorithm can hardly find the correct cell boundary using the noisy data. It is difficult to choose the proper isoline when several shocks are formed in the irregular steady state which is created due to a noise in the image—cf. Figure 5.1 (left and middle top). On the other hand, using a few steps of the nonlinear tensor anisotropic diffusion, all level lines are accumulated along the

TABLE 5.1

Error in the $L_2(I, L_2(\Omega))$ norm and the EOC comparing numerical and exact solutions.

n	h	k	Error	EOC
10	0.2	0.04	$1.809572 \cdot 10^{-4}$	
20	0.1	0.01	$0.3835138 \cdot 10^{-4}$	2.2383
40	0.05	0.0025	$0.09159927 \cdot 10^{-4}$	2.06587
80	0.025	0.000625	$0.02238713 \cdot 10^{-4}$	2.03267
160	0.0125	0.00015625	$0.00495121 \cdot 10^{-4}$	2.17682

cell boundary (just one shock is created in the final state of segmentation function)—cf. Figure 5.1 (left and middle bottom), and the cell can be segmented precisely. Now it is easy to choose the isoline for the cell boundary representation, we take the average of minimal and maximal values of the final segmentation function, and in Figure 5.2 we show the segmentation results for several cells visualizing it for both unfiltered and filtered images. We use the same parameters of the subjective surface segmentation method in both cases, and one can see much more precise segmentation results after filtering.

In Figures 5.3–5.4 we show two other real images: originals (left) and edge detection results for originals (middle) and after a few steps of filtering (right). Again, one can clearly see coherence enhancement and edge detection improvement.

In the last experiment we test the experimental order of convergence (EOC) of our method. In the theoretical part we prove its convergence, and the rigorous error estimates will be an objective of a further study; cf. [7]. Here we consider function $u(x, y, t) = t \cos(\pi x) \cos(\pi y)$ which fulfills the boundary conditions in the domain $\Omega = [-1, 1]^2$ and in time interval $I = [0, 1]$. Putting this function into the model equation (1.1), without convolutions, because we do not need to smooth either the function or the structure tensor, we get the nonzero right-hand side, and we modify the scheme accordingly. We take $C = 1$ and $\alpha = 0.001$ so the diffusion matrix D has eigenvalues between α and 1 and the process is strongly anisotropic. Then we take subsequently refined grids with $M = n^2$ finite volumes, $n = 10, 20, 40, 80, 160$, and the time step $k = h^2$, and we measure errors in the $L_2(I, L_2(\Omega))$ norm, which is natural for testing the schemes for solving parabolic problems. In Table 5.1 we report the errors for different grid sizes, and we observe that the EOC of our numerical scheme is equal to 2.

6. Acknowledgments. The authors thank Dr. Nadine Peyrieras from CNRS-DEPSN, Paris, coordinator of the European project Embryomics, for providing data used in our numerical experiments. We thank also Prof. Robert Eymard for giving us a hint on how to treat the T_1 term in the convergence proof.

REFERENCES

- [1] F. CATTÉ, P.-L. LIONS, J.-M. MOREL, AND T. COLL, *Image selective smoothing and edge detection by nonlinear diffusion*, SIAM J. Numer. Anal., 29 (1992), pp. 182–193.
- [2] S. CORSARO, K. MIKULA, A. SARTI, AND F. SGALLARI, *Semi-implicit co-volume method in 3D image segmentation*, SIAM J. Sci. Comput., 28 (2006), pp. 2248–2265.
- [3] Y. COUDIERE, J. P. VILA, AND P. VILLEDIEU, *Convergence rate of a finite volume scheme for a two-dimensional convection-diffusion problem*, M2AN Math. Model. Numer. Anal., 33 (1999), pp. 493–516.
- [4] R. EYMARD, T. GALLOUËT, AND R. HERBIN, *Finite Volume Methods*, in Handbook for Numerical Analysis 7 Ph. Ciarlet and J. L. Lions, eds., Elsevier, New York, 2000.

- [5] R. EYMARD, T. GALLOUËT, AND R. HERBIN, *A cell-centred finite-volume approximation for anisotropic diffusion operators on unstructured meshes in any space dimension*, IMA J. Numer. Anal., 26 (2006), pp. 326–353.
- [6] G. GILBOA, N. SOCHEN, AND Y. Y. ZEEVI, *Forward-and-backward diffusion processes for adaptive image enhancement and denoising*, IEEE Trans. Image Proc., 11 (2002), pp. 689–703.
- [7] A. HANDLOVIČOVÁ AND Z. KRIVÁ, *Error estimates for finite volume scheme for Perona-Malik equation*, Acta Math. Univ. Comenian. (N.S.), 74 (2005), pp. 79–94.
- [8] A. HANDLOVIČOVÁ, K. MIKULA, AND F. SGALLARI, *Semi-implicit complementary volume scheme for solving level set like equations in image processing and curve evolution*, Numer. Math., 93 (2003), pp. 675–695.
- [9] B. JÄHNE, *Spatio-temporal Image Processing*, Lecture Notes in Comput. Sci. 751, Springer, Berlin, 1993.
- [10] Z. KRIVÁ, *Explicit finite volume scheme for the Perona-Malik equation*, Comput. Methods Appl. Math., 5 (2005), pp. 170–200.
- [11] Z. KRIVÁ AND K. MIKULA, *An adaptive finite volume method in processing of color images*, in Proceedings of ALGORITMY 2000, Conference on Scientific Computing, Podbanské, 2000, pp. 174–187.
- [12] Z. KRIVÁ AND K. MIKULA, *An adaptive finite volume scheme for solving nonlinear diffusion equations arising in image processing*, J. Visual Comm. Image Repres., 13 (2002), pp. 22–35.
- [13] C. LANGE AND K. POLTHIER, *Anisotropic smoothing of point sets*, Com. Geom. Design, 22 (2005), pp. 680–692.
- [14] K. MIKULA, *Image processing with partial differential equations*, Modern Methods in Scientific Computing and Applications, A. Bourlioux and M. J. Gander, eds., Springer, Berlin, 2002, pp. 283–321.
- [15] K. MIKULA AND N. RAMAROSY, *Semi-implicit finite volume scheme for solving nonlinear diffusion equations in image processing*, Numer. Math., 89 (2001), pp. 561–590.
- [16] K. MIKULA, A. SARTI, AND F. SGALLARI, *Semi-implicit co-volume level set method in medical image segmentation*, in Handbook of Biomedical Image Analysis: Segmentation and Registration Models, J. Suri et al., eds., Springer, New York, 2005, pp. 583–626.
- [17] P. PERONA AND J. MALIK, *Scale space and edge detection using anisotropic diffusion*, in Proceedings of the IEEE Computer Society Workshop on Computer Vision, 1987.
- [18] T. PREUSSER AND M. RUMPF, *An adaptive finite element method for large scale image processing*, J. Visual Comm. Image Repres., 11 (2002), pp. 183–195.
- [19] M. RUMPF AND T. PREUSSER, *A level set method for anisotropic geometric diffusion in 3D image processing*, SIAM J. Appl. Math., 62 (2002), pp. 1772–1793.
- [20] A. SARTI, R. MALLADI, AND J. A. SETHIAN, *Subjective surfaces: A method for completing missing boundaries*, Proc. Nat. Acad. Sci. U.S.A., 12 (2000), pp. 6258–6263.
- [21] A. SARTI, K. MIKULA, AND F. SGALLARI, *Nonlinear multiscale analysis of 3D echocardiographic sequences*, IEEE Trans. Medical Imaging, 18 (1999), pp. 453–467.
- [22] J. WEICKERT, *Anisotropic Diffusion in Computer Vision*, Teubner, Stuttgart, 1998.
- [23] J. WEICKERT, *Coherence-enhancing diffusion of colour images*, Image and Vision Computing, 17 (1999), pp. 201–212.
- [24] J. WEICKERT, *Coherence-enhancing diffusion filtering*, Internat. J. Comput. Vision, 31 (1999), pp. 111–127.
- [25] J. WEICKERT AND H. SCHARR, *A scheme for coherence-enhancing diffusion filtering with optimized rotation invariance*, J. Visual Comm. Image Repres., 13 (2002), pp. 103–118.
- [26] X. YE, *A new discontinuous finite volume method for elliptic problems*, SIAM J. Numer. Anal., 42 (2004), pp. 1062–1072.