



**HAL**  
open science

# The Generalized Finite Volume SUSHI Scheme for the Discretization of Richards Equation

Konstantin Brenner, Danielle Hilhorst, Huy-Cuong Vu-Do

► **To cite this version:**

Konstantin Brenner, Danielle Hilhorst, Huy-Cuong Vu-Do. The Generalized Finite Volume SUSHI Scheme for the Discretization of Richards Equation. Vietnam Journal of Mathematics, 2015, 10.1007/s10013-015-0170-y . hal-01317568

**HAL Id: hal-01317568**

**<https://hal.science/hal-01317568>**

Submitted on 18 May 2016

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

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

# The generalized finite volume SUSHI scheme for the discretization of Richards Equation

Konstantin Brenner · Danielle Hilhorst ·  
Huy-Cuong Vu-Do

Received: date / Accepted: date

**Abstract** In this article, we apply the generalized finite volume method SUSHI to the discretization of Richards equation, an elliptic-parabolic equation modeling groundwater flow, where the diffusion term can be anisotropic and heterogeneous. This class of locally conservative methods can be applied to a wide range of unstructured possibly non-matching polyhedral meshes in arbitrary space dimension. As is needed for Richards equation, the time discretization is fully implicit. We obtain a convergence result based upon a priori estimates and the application of the Fréchet-Kolmogorov compactness theorem. We implement the scheme and present numerical tests.

**Keywords** Richards equation · finite volume scheme · SUSHI scheme

**Mathematics Subject Classification (2000)** 35K15 · 35K65 · 65M08 · 65M12 · 76S05

## 1 Introduction

Let  $\Omega$  be an open bounded polygonal subset of  $\mathbb{R}^d$  ( $d = 2$  or  $3$ ) and let  $T$  be a positive constant; we consider the Richards equation in the space-time domain  $Q_T = \Omega \times (0, T)$ :

$$\partial_t \theta(p) - \operatorname{div} \left( k_r(\theta(p)) \mathbf{K}(\mathbf{x}) \nabla(p+z) \right) = 0, \quad (1.1)$$

where  $p = p(\mathbf{x}, t)$  is the piezometric head. The space coordinates are defined by  $\mathbf{x} = (x, z)$  in the case of space dimension 2 and  $\mathbf{x} = (x, y, z)$  in the case of space dimension 3. The quantity  $\theta(p)$  is the water storage capacity, also known as the saturation,  $\mathbf{K}(\mathbf{x})$  is the absolute

---

Konstantin Brenner  
LJAD University Nice Sophia-Antipolis & Coffee team Inria Sophia-Antipolis - Méditerranée, France  
E-mail: konstantin.brenner@unice.fr

Danielle Hilhorst  
Laboratoire de Mathématiques, CNRS et Université de Paris-Sud, Orsay, France  
E-mail: Danielle.Hilhorst@math.u-psud.fr

Huy-Cuong Vu-Do  
Laboratoire de Mathématiques, Université de Paris-Sud, Orsay, France  
E-mail: vdhuycuong@math.u-psud.fr

permeability tensor and the scalar function  $k_r$  is the relative permeability.

Next we perform Kirchhoff's transformation. We define

$$F(s) := \int_0^s k_r(\theta(\tau)) d\tau,$$

and suppose that the function  $F$  is invertible. Then we set  $u = F(p)$  in  $Q_T$  and  $c(u) = c(F(p)) = \theta(p)$ ; it turns out that the function  $c$  is either qualitatively similar to the function  $\theta$  or has a support which is bounded from the left as in Figure 1-(b). We remark that Kirchhoff's transformation leads to  $\nabla u = k_r(\theta(p))\nabla p$ . The equation (1.1) becomes

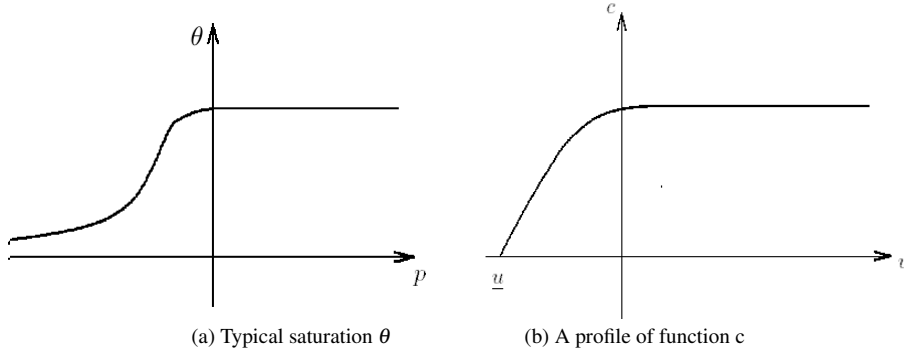
$$\partial_t c(u) - \operatorname{div}(\mathbf{K}(\mathbf{x})\nabla u) - \operatorname{div}(k_r(c(u))\mathbf{K}(\mathbf{x})\nabla z) = 0. \quad (1.2)$$

We suppose that  $\hat{u} \in W^{1,\infty}(\Omega)$  is a given function and consider the equation (1.2) together with the inhomogeneous Dirichlet boundary

$$u(\mathbf{x}, t) = \hat{u}(\mathbf{x}) \quad \text{a.e. on } \partial\Omega \times (0, T), \quad (1.3)$$

and the initial condition

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}) \quad \text{a.e. in } \Omega. \quad (1.4)$$



**Fig. 1** Typical saturation and its Kirchhoff's transformation

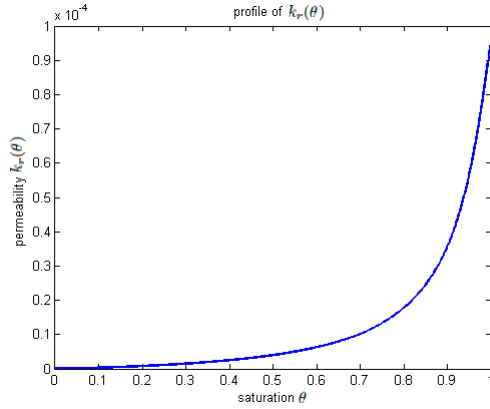
We denote by  $(P)$  the problem given by the equations (1.2), (1.3) and (1.4). We shall make the following hypotheses:

$(\mathcal{H}_1)$   $c \in W^{1,\infty}(\mathbb{R})$ ,  $0 \leq c \leq c_{max}$ , is a nondecreasing Lipschitz continuous function with Lipschitz constant  $L_c$ ;

$(\mathcal{H}_2)$   $k_r : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  is a nondecreasing Lipschitz continuous function with Lipschitz constant  $L_{k_r}$ , and  $0 \leq k_r \leq \bar{k}_r$  where  $\bar{k}_r$  is a positive constant,  $k_r(s) > 0$  for all  $s > 0$ ;

$(\mathcal{H}_3)$   $\mathbf{K}$  is a bounded function from  $\Omega$  to  $\mathbb{M}_d(\mathbb{R})$ , where  $\mathbb{M}_d(\mathbb{R})$  denotes the set of real  $d \times d$  matrices. Moreover for a.e.  $\mathbf{x}$  in  $\Omega$ ,  $\mathbf{K}(\mathbf{x})$  is a positive definite matrix and there exist two positive constants  $\bar{\mathbf{K}}$  and  $\underline{\mathbf{K}}$  such that the eigenvalues of  $\mathbf{K}(\mathbf{x})$  are included in  $[\bar{\mathbf{K}}, \underline{\mathbf{K}}]$ ;

$(\mathcal{H}_4)$   $u_0 \in L^\infty(\Omega)$  and  $\hat{u} \in W^{1,\infty}(\Omega)$ .



**Fig. 2** Typical permeability

**Definition 1** [Weak solution] A function  $u = u(\mathbf{x}, t)$  is said to be a weak solution of Problem  $(\mathcal{P})$  if:

$$\begin{aligned}
 (i) \quad & u - \hat{u} \in L^2(0, T; H_0^1(\Omega)); \\
 (ii) \quad & c(u) \in L^\infty(0, T; L^2(\Omega)); \\
 (iii) \quad & - \int_0^T \int_\Omega c(u) \partial_t \psi \, d\mathbf{x} dt - \int_\Omega c(u_0) \psi(\cdot, 0) \, d\mathbf{x} \\
 & + \int_0^T \int_\Omega \mathbf{K} \nabla u \cdot \nabla \psi \, d\mathbf{x} dt + \int_0^T \int_\Omega k_r(c(u)) \mathbf{K} \nabla z \cdot \nabla \psi \, d\mathbf{x} dt = 0,
 \end{aligned} \tag{1.5}$$

for all  $\psi \in L^2(0, T; H_0^1(\Omega))$  such that  $\psi(\cdot, T) = 0$  and  $\partial_t \psi \in L^\infty(Q_T)$ .

The discretization of the Richards equation was performed by means of the finite difference method by Hornung [17] and by means of the finite element method by Knabner [19]; Kelanemer [18] and Chounet et. al. [4] implemented a mixed finite element method and Frolkovic et. al. [15] applied a finite volume scheme on the dual mesh of a finite element mesh. We refer to Eymard, Gutnic and Hilhorst [13] and to Eymard, Gallouët, Gutnic, Herbin and Hilhorst [8] for the study of the convergence of two slightly different numerical schemes based upon the standard finite volume method.

In section 2.2, we introduce the SUSHI scheme, a finite volume scheme using stabilization and hybrid interfaces which has been proposed by Eymard et. al. [9] and define the approximate Problem  $(P_{\mathcal{D}, \delta t})$ . We also present some relevant results which will be useful in the sequel. In section 2.3, we prove an a priori estimate on the approximate solution in a discrete norm corresponding to a norm in  $L^2(0, T; H_0^1(\Omega))$ . Using these estimates and arguments based on the topological degree, we prove the existence of a solution of Problem  $(P_{\mathcal{D}, \delta t})$  in section 2.4. In section 2.5, we prove estimates on differences of time and space translates. These estimates imply the relative compactness of sequences of approximate solutions by the Fréchet-Kolmogorov theorem. We deduce the convergence in  $L^2$  of a sequence of approximate solutions to a solution of the continuous problem  $(P)$  in section 2.6. For the proofs, we apply methods inspired upon those of [9] and [10]. In the last section we describe

effective computations in the case of some well-known numerical tests which are often used in literature. We also perform simulations for a realistic model at the end of section 2.7.

The discretization of Richards equation by means of gradient schemes, which include the SUSHI method, has already been proposed by Eymard, Guichard, Herbin and Masson [12], where they consider Richards equation as a special case of two phase flow; however, they make the extra hypothesis that the relative permeability  $k_r$  is bounded away from zero, which is not satisfied in most geological contexts. Here we avoid this extra hypothesis by performing Kirchhoff's transformation.

## 2 The hybrid finite volume scheme SUSHI

In this section, we construct an approximate solution of Problem (P) corresponding to a time implicit discretization and a hybrid finite volume scheme. We follow the idea of Eymard et al. [9] to construct the fluxes using a stabilised discrete gradient.

### 2.1 Space and Time Discretization

Let us first define the notion of admissible finite volume mesh of  $\Omega$  and some notations associated with it.

**Definition 2 (Space discretization)** Let  $\Omega$  be a polyhedral open bounded connected subset of  $\mathbb{R}^d$  and  $\partial\Omega = \overline{\Omega} \setminus \Omega$  its boundary. A discretization of  $\Omega$ , denoted by  $\mathcal{D}$ , is defined as the triplet  $\mathcal{D} = (\mathcal{M}, \mathcal{E}, \mathcal{P})$ , where:

1.  $\mathcal{M}$  is a finite family of non empty convex open disjoint subsets of  $\Omega$  (the "control volumes") such that  $\overline{\Omega} = \bigcup_{K \in \mathcal{M}} \overline{K}$ . For any  $K \in \mathcal{M}$ , let  $\partial K = \overline{K} \setminus K$  be the boundary of  $K$ ; we denote by  $|K|$  the measure of  $K$  and  $d(K)$  the diameter of  $K$ .
2.  $\mathcal{E}$  is a finite family of disjoint subsets of  $\overline{\Omega}$  (the "interfaces"), such that, for all  $\sigma \in \mathcal{E}$ ,  $\sigma$  is a nonempty open subset of a hyperplane of  $\mathbb{R}^d$  and denote by  $|\sigma|$  its measure. We assume that, for all  $K \in \mathcal{M}$ , there exists a subset  $\mathcal{E}_K$  of  $\mathcal{E}$  such that  $\partial K = \bigcup_{\sigma \in \mathcal{E}_K} \sigma$ .
3.  $\mathcal{P}$  is a family of points of  $\Omega$  indexed by  $\mathcal{M}$ , denoted by  $\mathcal{P} = (\mathbf{x}_K)_{K \in \mathcal{M}}$ , such that for all  $K \in \mathcal{M}$ ,  $\mathbf{x}_K \in K$  and  $K$  is assumed to be  $\mathbf{x}_K$ -star-shaped, which means that for all  $\mathbf{x} \in K$ , the inclusion  $[\mathbf{x}_K, \mathbf{x}] \subset K$  holds.

For all  $\sigma \in \mathcal{E}$ , we denote by  $\mathbf{x}_\sigma$  the barycenter of  $\sigma$ . For all  $K \in \mathcal{M}$  and  $\sigma \in \mathcal{E}_K$ , we denote by  $D_{K,\sigma}$  the cone with vertex  $\mathbf{x}_K$  and basis  $\sigma$ , by  $\mathbf{n}_{K,\sigma}$  the unit vector normal to  $\sigma$  outward to  $K$  and by  $d_{K,\sigma}$  the Euclidean distance between  $\mathbf{x}_K$  and the hyperplane including  $\sigma$ . For any  $\sigma \in \mathcal{E}$ , we define  $\mathcal{M}_\sigma = \{K \in \mathcal{M} : \sigma \in \mathcal{E}_K\}$ . The set of boundary interfaces is denoted by  $\mathcal{E}_{ext}$  and the set of interior interfaces is denoted by  $\mathcal{E}_{int}$ .

We express the finite volume scheme in a weak form. For that purpose, let us first associate with the mesh the following spaces of discrete unknowns

$$X_{\mathcal{D}} = \{v = ((v_K)_{K \in \mathcal{M}}, (v_\sigma)_{\sigma \in \mathcal{E}}) : v_K \in \mathbb{R}, v_\sigma \in \mathbb{R}\},$$

$$X_{\mathcal{D},0} = \{v \in X_{\mathcal{D}} : v_\sigma = 0 \forall \sigma \in \mathcal{E}_{ext}\}.$$

**Definition 3 (Time discretization)** We divide the time interval  $(0, T)$  into  $N$  uniform time steps of length  $\delta t = T/N$ , and we define by  $t_n = n\delta t$  where  $n \in \{0, \dots, N\}$ .

Taking into account the time discretization leads us to define the following discrete spaces

$$\begin{aligned} X_{\mathcal{D}}^{\delta t} &= X_{\mathcal{D}}^N = \{h = (h^n)_{n \in \{1, \dots, N\}}, h^n \in X_{\mathcal{D}}\}, \\ X_{\mathcal{D},0}^{\delta t} &= X_{\mathcal{D},0}^N = \{h = (h^n)_{n \in \{1, \dots, N\}}, h^n \in X_{\mathcal{D},0}\}. \end{aligned}$$

For the sake of simplicity, we restrict our study to the case of constant time steps. Nevertheless all results presented below can be easily extended to the case of a non uniform time discretization.

## 2.2 Discrete weak formulation

We propose here a discrete scheme which is based upon the hybrid finite volume scheme SUSHI. It has been initially proposed for uniformly elliptic problems [9]. Schemes of this type in the case of a parabolic degenerate equation have been recently analyzed in [1]. Remark that this method can also be viewed as a mimetic finite difference or a mixed finite volume method (see [1], [2], [7]).

After formally integrating the equation (1.2) on the cell  $K \times (t_{n-1}, t_n)$  for each  $K \in \mathcal{M}$  and for each  $n \in \{1, \dots, N\}$ , we obtain

$$\begin{aligned} \int_K \left( c(u(\mathbf{x}, t_n)) - c(u(\mathbf{x}, t_{n-1})) \right) d\mathbf{x} - \sum_{\sigma \in \mathcal{E}_K} \int_{t_{n-1}}^{t_n} \int_{\sigma} \mathbf{K} \nabla u \cdot \mathbf{n}_{K,\sigma} d\gamma dt \\ - \sum_{\sigma \in \mathcal{E}_K} \int_{t_{n-1}}^{t_n} \int_{\sigma} k_r(c(u)) \mathbf{K} \nabla z \cdot \mathbf{n}_{K,\sigma} d\gamma dt = 0. \end{aligned} \quad (2.1)$$

For all  $K \in \mathcal{M}$  and  $\sigma \in \mathcal{E}_K$ , the diffusive flux  $-\int_{\sigma} \mathbf{K} \nabla u \cdot \mathbf{n}_{K,\sigma} d\gamma$  and the convective flux  $-\int_{\sigma} k_r(c(u)) \mathbf{K} \nabla z \cdot \mathbf{n}_{K,\sigma} d\gamma$  are approximated by  $F_{K,\sigma}(u)$  and  $Q_{K,\sigma}(u)$ , which are defined below by (2.7) and by (2.16), respectively.

Before introducing the numerical scheme, we define the following projection operator: let  $\phi \in C(\overline{Q_T})$ . We denote by  $P_{\mathcal{D}}\phi$  the element of  $X_{\mathcal{D}}$  defined by  $\left\{ \{\phi(\mathbf{x}_K, \cdot)\}, \{\phi(\mathbf{x}_{\sigma}, \cdot)\} \right\}$  for all  $K \in \mathcal{M}$  and  $\sigma \in \mathcal{E}$ .

Let  $P_{\mathcal{D}}\hat{u}$  and  $P_{\mathcal{D}}u_0$  be the projections of boundary and initial functions in (1.3) and (1.4), respectively; we present below a discrete weak problem ( $P_{\mathcal{D},\delta t}$ ):

The initial condition is discretized by

$$u^0 = \frac{1}{|K|} \int_K u_0(\mathbf{x}) d\mathbf{x} \quad \forall K \in \mathcal{M}. \quad (2.2)$$

For each  $n \in \{1, \dots, N\}$  and for all  $K \in \mathcal{M}$ , find  $u^n$  such that  $u^n - P_{\mathcal{D}}\hat{u} \in X_{\mathcal{D},0}$  satisfying

$$\sum_{K \in \mathcal{M}} |K| (c(u_K^n) - c(u_K^{n-1})) v_K + \delta t \langle u^n, v \rangle_F + \delta t \langle u^n, v \rangle_Q = 0 \quad \forall v \in X_{\mathcal{D},0}, \quad (2.3)$$

where

$$\langle w, v \rangle_F := \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}(w)(v_K - v_\sigma), \quad (2.4)$$

and

$$\langle w, v \rangle_Q := \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} Q_{K,\sigma}(w)(v_K - v_\sigma). \quad (2.5)$$

Let  $\tilde{u}_{\mathcal{D}}^{\delta t} = u_{\mathcal{D}}^{\delta t} - P_{\mathcal{D}}\hat{u} \in X_{\mathcal{D},0}^{\delta t}$ ; we rewrite the discrete equation (2.3) as

$$\sum_{K \in \mathcal{M}} |K|(c(u_K^n) - c(u_K^{n-1}))v_K + \delta t \langle \tilde{u}^n, v \rangle_F + \delta t \langle P_{\mathcal{D}}\hat{u}, v \rangle_F + \delta t \langle u^n, v \rangle_Q = 0. \quad (2.6)$$

The discrete Problem  $(P_{\mathcal{D},\delta t})$  is given by initial condition (2.2) and either the discrete equation (2.3) or the discrete equation (2.6).

### 2.3 The approximate flux

The discrete flux  $F_{K,\sigma}$  is expressed in terms of the discrete unknowns. For this purpose we apply the SUSHI scheme proposed in [9]. The idea is based upon the identification of the numerical fluxes through the mesh dependent bilinear form, using the expression of the discrete gradient

$$\sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} F_{K,\sigma}(w)(v_K - v_\sigma) = \int_{\Omega} \nabla_{\mathcal{D}} w(\mathbf{x}) \cdot \mathbf{K}(\mathbf{x}) \nabla_{\mathcal{D}} v(\mathbf{x}) \, d\mathbf{x} \quad \forall v, w \in X_{\mathcal{D},0}. \quad (2.7)$$

To this purpose, we first define

$$\nabla_K w = \frac{1}{|K|} \sum_{\sigma \in \mathcal{E}_K} |\sigma| (w_\sigma - w_K) \mathbf{n}_{K,\sigma} \quad \forall K \in \mathcal{M}, \forall w \in X_{\mathcal{D}}. \quad (2.8)$$

The consistency of formula (2.8) stems from the following geometrical relation:

$$\sum_{\sigma \in \mathcal{E}_K} |\sigma| \mathbf{n}_{K,\sigma} (\mathbf{x}_\sigma - \mathbf{x}_K)^T = |K| \mathbf{Id} \quad \forall K \in \mathcal{M}, \quad (2.9)$$

where  $(\mathbf{x}_\sigma - \mathbf{x}_K)^T$  is the transpose of  $\mathbf{x}_\sigma - \mathbf{x}_K \in \mathbb{R}^d$  and  $\mathbf{Id}$  is the  $d \times d$  identity matrix.

*Remark 1* The approximation formula (2.8) is exact for linear functions. Indeed, for any linear function defined on  $\Omega$  by  $\varphi(\mathbf{x}) = \mathbf{G} \cdot \mathbf{x}$  with  $\mathbf{G} \in \mathbb{R}^d$ , assuming that  $w_\sigma = \varphi(\mathbf{x}_\sigma)$  and  $w_K = \varphi(\mathbf{x}_K)$ , we obtain  $w_\sigma - w_K = (\mathbf{x}_\sigma - \mathbf{x}_K)^T \mathbf{G} = (\mathbf{x}_\sigma - \mathbf{x}_K)^T \nabla \varphi$ ; hence (2.8) leads to  $\nabla_K w = \nabla \varphi$ .

We also remark that

$$\sum_{\sigma \in \mathcal{E}_K} |\sigma| \mathbf{n}_{K,\sigma} = \sum_{\sigma \in \mathcal{E}_K} \int_{\sigma} \mathbf{n}_{K,\sigma} \, d\gamma = \int_K (\nabla 1) \, d\mathbf{x} = 0,$$

which means that the coefficient of  $w_K$  in (2.8) is equal to zero. Thus, a reconstruction of the discrete gradient solely based on (2.8) cannot lead to a definite discrete bilinear form in the general case. Therefore we introduce the stabilized gradient

$$\nabla_{K,\sigma} w = \nabla_K w + R_{K,\sigma} w \mathbf{n}_{K,\sigma}, \quad (2.10)$$

where

$$R_{K,\sigma}w = \frac{\sqrt{d}}{d_{K,\sigma}} \left( w_\sigma - w_K - \nabla_K w \cdot (\mathbf{x}_\sigma - \mathbf{x}_K) \right). \quad (2.11)$$

We may then define  $\nabla_{\mathcal{D}}w$  as the piecewise constant function equal to  $\nabla_{K\sigma}w$  a.e. in the cone  $D_{K\sigma}$

$$\nabla_{\mathcal{D}}w(\mathbf{x}) = \nabla_{K,\sigma}w \quad \text{for a.e. } \mathbf{x} \in D_{K,\sigma}. \quad (2.12)$$

Note that, from the definition (2.11), in view of (2.9) and (2.8), we deduce that

$$\sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} R_{K,\sigma}w \mathbf{n}_{K,\sigma} = 0 \quad \forall K \in \mathcal{M}. \quad (2.13)$$

In order to identify the numerical fluxes  $F_{K,\sigma}(w)$  through relation (2.7), we put the discrete gradient in the form

$$\nabla_{K,\sigma}w = \sum_{\sigma' \in \mathcal{E}_K} (w_{\sigma'} - w_K) \mathbf{y}_K^{\sigma\sigma'}, \quad (2.14)$$

with

$$\mathbf{y}_K^{\sigma\sigma'} = \begin{cases} \frac{|\sigma|}{|K|} \mathbf{n}_{K,\sigma} + \frac{\sqrt{d}}{d_{K,\sigma}} \left( 1 - \frac{|\sigma|}{|K|} \mathbf{n}_{K,\sigma} \cdot (\mathbf{x}_\sigma - \mathbf{x}_K) \right) \mathbf{n}_{K,\sigma} & \text{if } \sigma = \sigma', \\ \frac{|\sigma'|}{|K|} \mathbf{n}_{K,\sigma'} - \frac{\sqrt{d}}{d_{K,\sigma}} \frac{|\sigma'|}{|K|} \mathbf{n}_{K,\sigma'} \cdot (\mathbf{x}_\sigma - \mathbf{x}_K) \mathbf{n}_{K,\sigma} & \text{otherwise.} \end{cases}$$

Thus

$$\int_K \nabla_{\mathcal{D}}w(\mathbf{x}) \cdot \mathbf{K}(\mathbf{x}) \nabla_{\mathcal{D}}v(\mathbf{x}) d\mathbf{x} = \sum_{\sigma \in \mathcal{E}_K} \sum_{\sigma' \in \mathcal{E}_K} A_K^{\sigma\sigma'} (w_\sigma - w_K) (v_{\sigma'} - v_K^n),$$

with  $\sigma, \sigma' \in \mathcal{E}_K$  and

$$A_K^{\sigma\sigma'} = \sum_{\sigma'' \in \mathcal{E}_K} \mathbf{y}_K^{\sigma''\sigma} \cdot \Lambda_K^{\sigma''} \mathbf{y}_K^{\sigma''\sigma'}, \quad \Lambda_K^{\sigma''} = \int_{D_{K,\sigma''}} \mathbf{K}(\mathbf{x}) d\mathbf{x}.$$

The local matrices  $A_K^{\sigma\sigma'}$  are symmetric and positive, and the identification of the numerical fluxes using relation (2.7) leads to the expression:

$$F_{K,\sigma}(w) = \sum_{\sigma' \in \mathcal{E}_K} A_K^{\sigma\sigma'} (w_K - w_{\sigma'}). \quad (2.15)$$

Next we consider the convective flux. To this purpose we first define  $g_{K,\sigma} = \int_{\sigma} \mathbf{K} \nabla z \cdot \mathbf{n}_{K,\sigma} d\gamma$ . Then the convective flux is defined as

$$Q_{K,\sigma}(w) = -k_r(c(w_{K,\sigma})) g_{K,\sigma} \quad \forall K \in \mathcal{M}, \sigma \in \mathcal{E}_K, \quad (2.16)$$

where  $w_{K,\sigma}$  satisfies the upwind-sort formula

$$w_{K,\sigma} = \begin{cases} w_K & \text{if } g_{K,\sigma} < 0, \\ w_\sigma & \text{otherwise.} \end{cases} \quad (2.17)$$

Moreover, in view of the definition of  $g_{K,\sigma}$ , we remark that

$$g_{K,\sigma} = -g_{L,\sigma} \quad \forall \sigma \in \mathcal{E}_{int}, \mathcal{M}_\sigma = \{K, L\}. \quad (2.18)$$



## 2.4 The properties of the scheme

Next, we introduce some extra notations related to the mesh. Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2. The size of the discretization  $\mathcal{D}$  is defined by

$$l_{\mathcal{D}} = \sup_{K \in \mathcal{M}} d(K),$$

and the regularity of the mesh by:

$$\mu_{\mathcal{D}} = \max\left(\max_{\sigma_{K,L} \in \mathcal{E}_{int}} \frac{d_{K,\sigma}}{d_{L,\sigma}}, \max_{K \in \mathcal{M}, \sigma \in \mathcal{E}_K} \frac{d(K)}{d_{K,\sigma}}\right).$$

We will suppose in the sequel that  $l_{\mathcal{D}} \leq 1$ .

**Definition 4** Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2, and let  $\delta t$  be the time step defined in Definition 3. For  $v \in X_{\mathcal{D}}$ , we define the semi-norm

$$|v|_{X_{\mathcal{D}}}^2 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (v_{\sigma} - v_K)^2,$$

For all  $h = \{h^n\}_{n \in \{1, \dots, N\}} \in X_{\mathcal{D}}^{\delta t}$ , we define the semi-norm

$$|h|_{X_{\mathcal{D}}^{\delta t}}^2 = \sum_{n=1}^N \delta t |h^n|_{X_{\mathcal{D}}}^2.$$

Let  $H_{\mathcal{M}}(\Omega) \subset L^2(\Omega)$  be the set of piecewise constant functions on the control volumes of the mesh  $\mathcal{M}$ . For all  $v \in X_{\mathcal{D}}$  we denote by  $\Pi_{\mathcal{M}} v \in H_{\mathcal{M}}(\Omega)$  the piecewise function from  $\Omega$  to  $\mathbb{R}$  defined by  $\Pi_{\mathcal{M}} v(\mathbf{x}) = v_K$  for almost every  $\mathbf{x} \in K$ , for all  $K \in \mathcal{M}$ .

Let  $H_{\mathcal{M}}^{\delta t}(\Omega \times (0, T)) \subset L^2(\Omega \times (0, T))$  be the set of piecewise constant functions on the space-time control volumes. We denote by  $\Pi_{\mathcal{M}}^{\delta t} : X_{\mathcal{D}}^{\delta t} \rightarrow L^2(Q_T)$  the mapping

$$\Pi_{\mathcal{M}}^{\delta t} v(\mathbf{x}, t) = v_K^n \quad \text{for all } (\mathbf{x}, t) \in K \times (t_{n-1}, t_n]. \quad (2.19)$$

We also define  $\nabla_{\mathcal{D}}^{\delta t} : X_{\mathcal{D}}^{\delta t} \rightarrow L^2(Q_T)^d$  by

$$\nabla_{\mathcal{D}}^{\delta t} v(\mathbf{x}, t) = \nabla_{\mathcal{D}} v^n \quad \text{for all } (\mathbf{x}, t) \in K \times (t_{n-1}, t_n]. \quad (2.20)$$

Next, following [9], for all  $v \in X_{\mathcal{D}}$  we define the following related norm

$$\|\Pi_{\mathcal{M}} v\|_{1,2,\mathcal{M}}^2 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} \left( \frac{D_{\sigma} v}{d_{\sigma}} \right)^2, \quad (2.21)$$

with  $d_{\sigma} = |d_{K,\sigma} + d_{L,\sigma}|$ ,  $D_{\sigma} v = |v_K - v_L|$  if  $\mathcal{M}_{\sigma} = \{K, L\}$ , and  $d_{\sigma} = d_{K,\sigma}$ ,  $D_{\sigma} v = |v_K|$  if  $\mathcal{M}_{\sigma} = K$ . A result stated in [9] gives the relation

$$\|\Pi_{\mathcal{M}} v\|_{1,2,\mathcal{M}}^2 \leq |v|_{X_{\mathcal{D}}}^2 \quad \forall v \in X_{\mathcal{D},0}. \quad (2.22)$$

**Lemma 1 (Poincaré like inequality)** *Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2. Let  $\eta > 0$  be such that  $\eta \leq d_{K,\sigma}/d_{L,\sigma} \leq 1/\eta$  for all  $\sigma \in \mathcal{E}_{int}$ , where  $\mathcal{M}_{\sigma} = \{K, L\}$ . Then there exists  $C_1$  only depending on  $d$ ,  $\Omega$  and  $\eta$  such that*

$$\|\Pi_{\mathcal{M}} v\|_{L^2(\Omega)} \leq C_1 \|\Pi_{\mathcal{M}} v\|_{1,2,\mathcal{M}} \quad \forall v \in X_{\mathcal{D}}, \quad (2.23)$$

where  $\|\Pi_{\mathcal{M}} v\|_{1,2,\mathcal{M}}$  is defined by (2.21).

*Proof:* In view of Lemma 5.4 in [9], for each  $p \geq 1$  there exists  $q > p$  only depending on  $p$  and there exists a positive constant  $C$  only depending on  $d$  and  $\eta$  such that  $\|\Pi_{\mathcal{M}} v\|_{L^q(\Omega)} \leq C \|\Pi_{\mathcal{M}} v\|_{1,p,\mathcal{M}}$  for all  $v \in X_D$ . We remark that for all  $q > p$ , then

$$\|\Pi_{\mathcal{M}} v\|_{L^p(\Omega)} \leq |\Omega|^{(q-p)/pq} \|\Pi_{\mathcal{M}} v\|_{L^q(\Omega)}.$$

We set  $p = 2$  and  $C_1 = |\Omega|^{(q-2)/2q} C$  to conclude the proof.  $\square$

**Definition 5** Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2, and let  $\delta t$  be the time step defined in Definition 3. We define the  $L^2$ -norm of the discrete gradient by

$$\|\nabla_{\mathcal{D}} v(\mathbf{x})\|_{L^2(\Omega)^d}^2 = \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} |\nabla_{K,\sigma} v|^2 \quad \forall v \in X_{\mathcal{D}},$$

and

$$\|\nabla_{\mathcal{D}}^{\delta t} h(\mathbf{x}, t)\|_{L^2(Q_T)^d}^2 = \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} |\nabla_{K,\sigma} h^n|^2 \quad \forall h \in X_{\mathcal{D}}^{\delta t},$$

where  $\nabla_{K,\sigma}$  and  $\nabla_{\mathcal{D}}$  is defined by (2.8)-(2.12).

**Lemma 2** Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2 and suppose that there exists a positive constant  $\mu$  such that  $\mu_{\mathcal{D}} \leq \mu$  for all  $\mathcal{D}$ ; let  $\delta t$  be the time step defined in Definition 3.

(i) Then there exist positive constants  $C_2$  and  $C_3$  only depending on  $\mu$  and  $d$  such that

$$C_2 |v|_{X_{\mathcal{D}}}^2 \leq \|\nabla_{\mathcal{D}} v(\mathbf{x})\|_{L^2(\Omega)^d}^2 \leq C_3 |v|_{X_{\mathcal{D}}}^2 \quad \forall v \in X_{\mathcal{D}}.$$

(ii) Moreover, we have

$$C_2 |h|_{X_{\mathcal{D}}^{\delta t}}^2 \leq \|\nabla_{\mathcal{D}}^{\delta t} h(\mathbf{x}, t)\|_{L^2(Q_T)^d}^2 \leq C_3 |h|_{X_{\mathcal{D}}^{\delta t}}^2 \quad \forall h \in X_{\mathcal{D}}^{\delta t}.$$

*Proof:* We refer to Lemma 4.2 in [9] for the proof of (i). In view of the definition of the semi-norm in the space  $X_{\mathcal{D}}^{\delta t}$  and the  $L^2$ -norm of the discrete gradient, we deduce (ii).  $\square$

**Lemma 3** Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2 and suppose that there exists a positive constant  $\mu$  such that  $\mu_{\mathcal{D}} \leq \mu$  for all  $\mathcal{D}$ ; there exists a positive constant  $\alpha$  such that

$$\langle v, v \rangle_F \geq \alpha |v|_{X_{\mathcal{D}}}^2. \quad (2.24)$$

*Proof:* In view of Hypothesis ( $\mathcal{H}_3$ ) and Lemma 2, we also obtain

$$\begin{aligned} \langle v, v \rangle_F &= \int_{\Omega} \mathbf{K}(\mathbf{x}) (\nabla_{\mathcal{D}} v(\mathbf{x}))^2 d\mathbf{x} \\ &\geq \underline{K} \|\nabla_{\mathcal{D}} v(\mathbf{x})\|_{L^2(\Omega)}^2 \\ &\geq \underline{K} C_2 |v|_{X_{\mathcal{D}}}^2. \end{aligned}$$

Setting  $\alpha = \underline{K} C_2$  permits to complete the proof.  $\square$

**Definition 6** Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2 and let  $\delta t$  be the time step defined in Definition 3. Let  $u_{\mathcal{D}}^{\delta t} \in X_{\mathcal{D}}^{\delta t}$  be a solution of Problem  $(P_{\mathcal{D}, \delta t})$ . We say that  $\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)$  is an approximate solution of Problem  $(P)$ .

We now state a weak compactness result for the discrete gradient.

**Lemma 4** *Let  $\mathcal{F}$  be a family of discretizations of  $\Omega$  in the sense of Definition 2 and suppose that there exists a positive constant  $\mu$  such that  $\mu_{\mathcal{D}} \leq \mu$  for all  $\mathcal{D} \in \mathcal{F}$ . Let  $(h_{\mathcal{D}}^{\delta t})_{\mathcal{D} \in \mathcal{F}}$  be a family of unknowns such that*

(i)  $h_{\mathcal{D}}^{\delta t} \in X_{\mathcal{D},0}^{\delta t}$  for all  $\mathcal{D} \in \mathcal{F}$  and for all  $\delta t \in (0, 1)$ ;

(ii) there exists  $C > 0$  such that  $|h_{\mathcal{D}}^{\delta t}|_{X_{\mathcal{D}}^{\delta t}} \leq C$  for all  $\mathcal{D} \in \mathcal{F}$  and for all  $\delta t \in (0, 1)$ ;

(iii) there exists  $h \in L^2(Q_T)$  such that  $\Pi_{\mathcal{M}}^{\delta t} h_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)$  converges to  $h$  weakly in  $L^2(Q_T)$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ .

Then  $h \in L^2(0, T, H_0^1(\Omega))$  and  $\nabla_{\mathcal{D}}^{\delta t} h_{\mathcal{D}}^{\delta t}$  converges to  $\nabla h$  weakly in  $L^2(Q_T)^d$  as  $l_{\mathcal{D}}$  and  $\delta t \rightarrow 0$ .

*Proof:* We extend the functions  $\Pi_{\mathcal{M}}^{\delta t} h_{\mathcal{D}}^{\delta t}$  and  $\nabla_{\mathcal{D}}^{\delta t} h_{\mathcal{D}}^{\delta t}$  by zero outside of  $\Omega$ . In view of (ii) of Lemma 2, there exists a function  $\mathcal{H} \in L^2(\mathbb{R}^d \times (0, T))^d$  such that up to a subsequence  $\nabla_{\mathcal{D}}^{\delta t} h_{\mathcal{D}}^{\delta t}$  weakly converges to  $\mathcal{H}$  in  $L^2(\mathbb{R}^d \times (0, T))^d$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ . We show below that  $\mathcal{H} = \nabla h$ . Let  $\varphi \in C_c^\infty(\mathbb{R}^d \times (0, T))^d$  be given; we consider the term defined by

$$T_1^{\mathcal{D}} = \int_0^T \int_{\mathbb{R}^d} \nabla_{\mathcal{D}}^{\delta t} h_{\mathcal{D}}^{\delta t}(\mathbf{x}, t) \cdot \varphi(\mathbf{x}, t) \, d\mathbf{x} dt.$$

In view of the definition of  $\nabla_{\mathcal{D}}$  in (2.8)-(2.12) we infer that  $T_1^{\mathcal{D}} = T_2^{\mathcal{D}} + T_3^{\mathcal{D}}$ , where

$$T_2^{\mathcal{D}} = \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| (h_{\sigma}^n - h_K^n) \mathbf{n}_{K,\sigma} \cdot \varphi_K^n \quad \text{with } \varphi_K^n = \frac{1}{\delta t |K|} \int_{t_{n-1}}^{t_n} \int_K \varphi(\mathbf{x}, t) \, d\mathbf{x} dt,$$

and

$$T_3^{\mathcal{D}} = \sum_{n=1}^N \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} R_{K,\sigma} h^n \mathbf{n}_{K,\sigma} \cdot \int_{t_{n-1}}^{t_n} \int_{D_{K,\sigma}} \varphi(\mathbf{x}, t) \, d\mathbf{x} dt,$$

which by (2.13) yields

$$T_3^{\mathcal{D}} = \sum_{n=1}^N \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} R_{K,\sigma} h^n \mathbf{n}_{K,\sigma} \cdot \int_{t_{n-1}}^{t_n} \int_{D_{K,\sigma}} (\varphi(\mathbf{x}, t) - \varphi_K^n) \, d\mathbf{x} dt.$$

Applying Cauchy-Schwarz inequality, we deduce that

$$(T_3^{\mathcal{D}})^2 \leq \left( \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} (R_{K,\sigma} h^n)^2 \right) \cdot \left( \sum_{n=1}^N \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{d}{|\sigma| d_{K,\sigma} \delta t} \left| \int_{t_{n-1}}^{t_n} \int_{D_{K,\sigma}} (\varphi - \varphi_K^n) \, d\mathbf{x} dt \right|^2 \right). \quad (2.25)$$

We deduce from formulas (4.11) and (4.12) in [9] the inequality

$$\begin{aligned} (R_{K,\sigma} h^n)^2 &\leq 2d \left( \left( \frac{h_{\sigma}^n - h_K^n}{d_{K,\sigma}} \right)^2 + \mu^2 |\nabla_K h^n|^2 \right) \\ &\leq 2d \left( \left( \frac{h_{\sigma}^n - h_K^n}{d_{K,\sigma}} \right)^2 + \mu^2 \frac{d}{|K|} \sum_{\sigma' \in \mathcal{E}_K} \frac{|\sigma'|}{d_{K,\sigma'}} (h_{\sigma'}^n - h_K^n)^2 \right), \end{aligned}$$

which in turn implies that

$$\frac{|\sigma| d_{K,\sigma}}{d} (R_{K,\sigma} h^n)^2 \leq 2 \frac{|\sigma|}{d_{K,\sigma}} (h_{\sigma}^n - h_K^n)^2 + 2\mu^2 \frac{|\sigma| d_{K,\sigma} d}{|K|} \sum_{\sigma' \in \mathcal{E}_K} \frac{|\sigma'|}{d_{K,\sigma'}} (h_{\sigma'}^n - h_K^n)^2.$$

We remark that  $\sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d|K|} = 1$  and that  $\sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (h_\sigma^n - h_K^n)^2 = |h_{\mathcal{D}}^{\delta t}|_{X_{\mathcal{D}}^{\delta t}}|^2 \leq C$ , which yields

$$\sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} (R_{K,\sigma} h^n)^2 \leq 2(1 + \mu^2 d^2) C. \quad (2.26)$$

By the regularity properties of the function  $\varphi$ , there exists  $C_\varphi$  only depending on  $\varphi$  such that  $|\int_{t_{n-1}}^{t_n} \int_{D_{K,\sigma}} (\varphi(\mathbf{x}, t) - \varphi_K^n) d\mathbf{x} dt| \leq C_\varphi \delta t l_{\mathcal{D}} \frac{|\sigma| d_{K,\sigma}}{d}$ , which implies that

$$\frac{d}{|\sigma| d_{K,\sigma} \delta t} \left| \int_{t_{n-1}}^{t_n} \int_{D_{K,\sigma}} (\varphi(\mathbf{x}, t) - \varphi_K^n) d\mathbf{x} dt \right|^2 \leq \delta t \frac{|\sigma| d_{K,\sigma}}{d} C_\varphi^2 l_{\mathcal{D}}^2.$$

Since  $\sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma| d_{K,\sigma}}{d} = |K|$ , it follows that

$$\sum_{n=1}^N \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{d}{|\sigma| d_{K,\sigma} \delta t} \left( \int_{t_{n-1}}^{t_n} \int_{D_{K,\sigma}} (\varphi - \varphi_K^n) d\mathbf{x} dt \right)^2 \leq T |\Omega| C_\varphi^2 l_{\mathcal{D}}^2. \quad (2.27)$$

From (2.25), (2.26) and (2.27), we deduce that  $\lim_{l_{\mathcal{D}}, \delta t \rightarrow 0} T_3^{\mathcal{D}} = 0$ . Next, we compare  $T_2^{\mathcal{D}}$  to  $T_4^{\mathcal{D}}$  defined by

$$T_4^{\mathcal{D}} = \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| (h_\sigma^n - h_K^n) \mathbf{n}_{K,\sigma} \cdot \varphi_\sigma^n \text{ with } \varphi_\sigma^n = \frac{1}{\delta t |\sigma|} \int_{t_{n-1}}^{t_n} \int_\sigma \varphi d\gamma dt.$$

We have that

$$\begin{aligned} (T_2^{\mathcal{D}} - T_4^{\mathcal{D}})^2 &\leq \left( \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K,\sigma}} (h_\sigma^n - h_K^n)^2 \right) \left( \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K,\sigma} |\varphi_K^n - \varphi_\sigma^n|^2 \right) \\ &\leq |h_{X_{\mathcal{D}}}^{\delta t}|^2 T d |\Omega| C_\varphi^2 l_{\mathcal{D}}^2, \end{aligned}$$

which leads to  $\lim_{l_{\mathcal{D}}, \delta t \rightarrow 0} \{T_2^{\mathcal{D}} - T_4^{\mathcal{D}}\} = 0$ . On the other hand, since

$$T_4^{\mathcal{D}} = - \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| h_K^n \mathbf{n}_{K,\sigma} \cdot \varphi_\sigma^n = - \int_0^T \int_{\mathbb{R}^d} \Pi_{\mathcal{M}}^{\delta t} h_{\mathcal{D}}^{\delta t}(\mathbf{x}, t) \operatorname{div} \varphi(\mathbf{x}, t) d\mathbf{x} dt,$$

it follows that  $\lim_{l_{\mathcal{D}}, \delta t \rightarrow 0} T_2^{\mathcal{D}} = - \int_0^T \int_{\mathbb{R}^d} h(\mathbf{x}, t) \operatorname{div} \varphi(\mathbf{x}, t) d\mathbf{x} dt$ . Thus the function  $\mathcal{H} \in L^2(\mathbb{R}^d \times (0, T))^d$  is a.e. equal to  $\nabla h$  in  $\mathbb{R}^d \times (0, T)$ . Since  $h = 0$  outside of  $\Omega$ , we deduce that  $h \in L^2(0, T, H_0^1(\Omega))$ , and the uniqueness of the limit implies that the whole family  $\nabla_{\mathcal{D}}^{\delta t} h_{\mathcal{D}}^{\delta t}$  weakly converges in  $L^2(\mathbb{R}^d \times (0, T))^d$  to  $\nabla h$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ .  $\square$

### 3 A priori estimates

**Lemma 5** *Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2, and let  $\delta t$  be a time step in the interval  $(0, T)$  in the sense of Definition 3. Let  $u_{\mathcal{D}}^{\delta t} \in X_{\mathcal{D}}^{\delta t}$  be the solution of Problem  $(P_{\mathcal{D}, \delta t})$ . Let  $\tilde{u}_{\mathcal{D}}^{\delta t} = u_{\mathcal{D}}^{\delta t} - P_{\mathcal{D}}\hat{u}$ . There exists a positive constant  $C_5$  only depending on  $\bar{K}, \bar{k}_r, T, \Omega, \alpha$  as well as on  $\|c\|_{L^\infty(\mathbb{R})}, \|u_0\|_{L^\infty(\Omega)}, \|\hat{u}\|_{L^\infty(\Omega)}$  and  $\|\hat{u}\|_{W^{1,\infty}(\Omega)}$  such that*

$$|\tilde{u}_{\mathcal{D}}^{\delta t}|_{X_{\mathcal{D}}^{\delta t}}^2 \leq C_5, \quad (3.1)$$

and

$$|u_{\mathcal{D}}^{\delta t}|_{X_{\mathcal{D}}^{\delta t}}^2 \leq C_5. \quad (3.2)$$

*Proof:* Setting  $v = \tilde{u}^n$  in the scheme (2.6) and summing over  $n \in \{1, \dots, N\}$  implies

$$\begin{aligned} & \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| \left( c(u_K^n) - c(u_K^{n-1}) \right) \left( u_K^n - (P_{\mathcal{D}}\hat{u})_K \right) \\ & + \sum_{n=1}^N \delta t \langle \tilde{u}^n, \tilde{u}^n \rangle_F + \sum_{n=1}^N \delta t \langle P_{\mathcal{D}}\hat{u}, \tilde{u}^n \rangle_F + \sum_{n=1}^N \delta t \langle u^n, \tilde{u}^n \rangle_Q = 0, \end{aligned} \quad (3.3)$$

which can be rewritten as

$$\bar{A}_1 - \bar{A}_2 + \bar{B}_1 + \bar{B}_2 + \bar{C} = 0, \quad (3.4)$$

where

$$\begin{aligned} \bar{A}_1 &= \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| \left( c(u_K^n) - c(u_K^{n-1}) \right) u_K^n, \\ \bar{A}_2 &= \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| \left( c(u_K^n) - c(u_K^{n-1}) \right) (P_{\mathcal{D}}\hat{u})_K, \\ \bar{B}_1 &= \sum_{n=1}^N \delta t \langle \tilde{u}^n, \tilde{u}^n \rangle_F, \\ \bar{B}_2 &= \sum_{n=1}^N \delta t \langle P_{\mathcal{D}}\hat{u}, \tilde{u}^n \rangle_F, \\ \bar{C} &= \sum_{n=1}^N \delta t \langle u^n, \tilde{u}^n \rangle_Q. \end{aligned} \quad (3.5)$$

Next we define

$$\Theta_K^n = c(u_K^n) u_K^n - \int_0^{u_K^n} c(\tau) d\tau.$$

Since  $c$  is nondecreasing, it follows that  $\Theta_K^n \geq 0$  for all  $n \in \{1, \dots, N\}$  and  $K \in \mathcal{M}$ . Moreover, we have that

$$\Theta_K^n - \Theta_K^{n-1} = \left( c(u_K^n) - c(u_K^{n-1}) \right) u_K^n + \int_{u_K^{n-1}}^{u_K^n} \left( c(u_K^{n-1}) - c(\tau) \right) d\tau,$$

where the last term is negative. Thus

$$\bar{A}_1 \geq \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| (\Theta_K^n - \Theta_K^{n-1}) = \sum_{K \in \mathcal{M}} |K| \Theta_K^N - \sum_{K \in \mathcal{M}} |K| \Theta_K^0.$$

Note that  $\Theta_K^N$  is positive and  $\Theta_K^0$  can be written

$$\Theta_K^0 = \int_0^{u_K^0} (c(u_K^0) - c(\tau)) d\tau \leq 2 \|u_0\|_{L^\infty(\Omega)} \|c\|_{L^\infty(\mathbb{R})}.$$

It implies that

$$-\bar{A}_1 \leq 2|\Omega| \|u_0\|_{L^\infty(\Omega)} \|c\|_{L^\infty(\mathbb{R})}. \quad (3.6)$$

In view of the hypotheses  $(\mathcal{H}_1)$  and  $(\mathcal{H}_4)$ , we obtain

$$|\bar{A}_2| \leq 2|\Omega| \|\hat{u}\|_{L^\infty(\Omega)} \|c\|_{L^\infty(\mathbb{R})}. \quad (3.7)$$

We deduce from the coercivity property in Lemma 3 that

$$\bar{B}_1 \geq \alpha \sum_{n=1}^N \delta t |\bar{u}^n|_{X_{\mathcal{D}}}^2. \quad (3.8)$$

Applying first Hölder's inequality and then Young's inequality, we deduce that there exists a positive constant  $C_{\hat{u}}$  such that for all  $\varepsilon_1 > 0$

$$\begin{aligned} |\bar{B}_2| &\leq \sum_{n=1}^N \delta t \bar{K} \|\nabla_{\mathcal{D}} P_{\mathcal{D}} \hat{u}\|_{L^2(\Omega)^d} \|\nabla_{\mathcal{D}} \bar{u}^n\|_{L^2(\Omega)^d} \\ &\leq \frac{1}{2\varepsilon_1} C_{\hat{u}} T \bar{K} \|\hat{u}\|_{W^{1,\infty}(\Omega)}^2 + \frac{\varepsilon_1}{2} \sum_{n=1}^N \delta t \bar{K} \|\nabla_{\mathcal{D}} \bar{u}^n\|_{L^2(\Omega)^d}^2 \\ &\leq \frac{1}{2\varepsilon_1} C_{\hat{u}} T \bar{K} \|\hat{u}\|_{W^{1,\infty}(\Omega)}^2 + \frac{\varepsilon_1}{2} C_3 \bar{K} \sum_{n=1}^N \delta t |\bar{u}^n|_{X_{\mathcal{D}}}^2. \end{aligned} \quad (3.9)$$

As for the term  $\bar{C}$ , we deduce from  $|g_{K\sigma}| \leq |\sigma| \bar{K}$  and Young's inequality that for all  $\varepsilon_2 > 0$

$$\begin{aligned} |\bar{C}| &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |k_r(c(u_{K\sigma}^n))| g_{K\sigma} (\bar{u}_K^n - \bar{u}_\sigma^n) \\ &\leq \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} (\bar{k}_r \sqrt{\bar{K}} |\sigma| d_{K\sigma}) \left( \sqrt{\bar{K}} \frac{|\sigma|}{d_{K\sigma}} |\bar{u}_K^n - \bar{u}_\sigma^n| \right) \\ &\leq \frac{1}{2\varepsilon_2} \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K\sigma} \bar{k}_r^2 \bar{K} + \frac{\varepsilon_2}{2} \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \bar{K} \frac{|\sigma|}{d_{K\sigma}} (\bar{u}_K^n - \bar{u}_\sigma^n)^2 \\ &\leq \frac{d}{2\varepsilon_2} T |\Omega| \bar{k}_r^2 \bar{K} + \frac{\varepsilon_2}{2} \bar{K} \sum_{n=1}^N \delta t |\bar{u}^n|_{X_{\mathcal{D}}}^2. \end{aligned} \quad (3.10)$$

We deduce from (3.4) that

$$\bar{B}_1 = -\bar{A}_1 + \bar{A}_2 - \bar{B}_2 - \bar{C},$$

so that

$$\bar{B}_1 \leq -\bar{A}_1 + |\bar{A}_2| + |\bar{B}_2| + |\bar{C}|. \quad (3.11)$$

We gather the inequalities (3.6)-(3.11). Then in view of Definition 4 of the space-time norm

$|\tilde{u}_{\mathcal{D}}^{\delta t}|_{X_{\mathcal{D}}^{\delta t}}^2 = \sum_{n=1}^N \delta t |\tilde{u}^n|_{X_{\mathcal{D}}}^2$ , we deduce that

$$\begin{aligned} \left(\alpha - \frac{\varepsilon_1 C_3 + \varepsilon_2 \bar{K}}{2}\right) |\tilde{u}_{\mathcal{D}}^{\delta t}|_{X_{\mathcal{D}}^{\delta t}}^2 &\leq 2|\Omega| \|c\|_{L^\infty(\mathbb{R})} (\|u_0\|_{L^\infty(\Omega)} + \|\hat{u}\|_{L^\infty(\Omega)}) \\ &\quad + \frac{1}{2\varepsilon_1} C_{\hat{u}} T \bar{K} \|\hat{u}\|_{W^{1,\infty}(\Omega)}^2 + \frac{d}{2\varepsilon_2} T |\Omega| \bar{k}_r^2 \bar{K}. \end{aligned}$$

Choosing  $\varepsilon_1 = \alpha/(2C_3\bar{K})$  and  $\varepsilon_2 = \alpha/(2\bar{K})$  permits to complete the proof of Lemma 5.  $\square$

#### 4 Existence of a discrete solution

Let  $\mu \in [0, 1]$  and  $u^{n-1} \in X_{\mathcal{D}}$ ; we consider the following extended problem. Find  $u^{n,\mu} \in X_{\mathcal{D}}$  such that for all  $v \in X_{\mathcal{D},0}$

$$\mu \sum_{K \in \mathcal{M}} |K| \left( c(u_K^{n,\mu}) - c(u_K^{n-1}) \right) v_K + \delta t \langle u^{n,\mu}, v \rangle_F + \mu \delta t \langle u^{n,\mu}, v \rangle_Q = 0. \quad (4.1)$$

It can be shown by a similar proof as that of Lemma 5 that the solution of the extended problem (4.1) satisfies

$$\delta t |u^{n,\mu}|_{X_{\mathcal{D}}}^2 \leq \mu C_6 \leq C_6, \quad (4.2)$$

where  $C_6$  only depends on  $\bar{K}, \bar{k}_r, T, \Omega, \alpha$  as well as on  $\|c\|_{L^\infty(\mathbb{R})}, \|u_0\|_{L^\infty(\Omega)}, \|\hat{u}\|_{L^\infty(\Omega)}$  and  $\|\hat{u}\|_{W^{1,\infty}(\Omega)}$ .

**Theorem 1 (Existence of a discrete solution)** *The discrete problem  $(P_{\mathcal{D},\delta t})$  possesses at least one solution.*

*Proof* The extended problem (4.1) can be written as the abstract system of nonlinear equations

$$H(u^{n,\mu}, u^{n-1}, \mu) = 0, \quad (4.3)$$

where  $H$  is a continuous mapping from  $X_{\mathcal{D}} \times X_{\mathcal{D}} \times [0, 1]$  to  $X_{\mathcal{D}}$ . Indeed, setting  $v_K = 1, v_L = 0$ , for all  $L \neq K$ ,  $v_\sigma = 0$  for all  $\sigma \in \mathcal{E}$ , we obtain the equation

$$H_K \left( c(u_K^{n,\mu}), c(u_K^{n-1}), u_K^{n,\mu}, (u_\sigma^{n,\mu})_{\sigma \in \mathcal{E}_K}, \mu \right) = 0 \quad \text{for all } K \in \mathcal{M},$$

and setting  $v_K = 0$  for all  $K \in \mathcal{M}$ ,  $v_\sigma = 1$  and  $v_{\sigma'} = 0$  for all  $\sigma' \neq \sigma$ , we deduce the equation

$$H_\sigma \left( (u_K^{n,\mu})_{K \in \mathcal{M}_\sigma}, ((u_\sigma^{n,\mu})_{\sigma \in \mathcal{E}_K})_{K \in \mathcal{M}_\sigma}, \mu \right) = 0 \quad \text{for all } \sigma \in \mathcal{E}_{int}.$$

Setting  $r = 2\sqrt{\frac{C_6}{\delta t}}$ , we deduce from (4.2) that the system (4.3) has no solution on the boundary of the ball  $B_r$  of radius  $r$  for  $\mu \in [0, 1]$ .

Before pursuing the proof, we recall results due to [5, Theorem 3.1].

**Proposition 1** Let  $M = \{(f, \Omega, y) \text{ with } \Omega \text{ an open bounded set of } \mathbb{R}^n, f \in C(\overline{\Omega}) \text{ and } y \notin f(\partial\Omega)\}$  and let  $d : M \rightarrow \mathbb{Z}$  be the topological degree. Then  $d$  has the following properties.

(d1)  $d(\text{id}, \Omega, y) = 1$  for  $y \in \Omega$ .

(d2)  $d(f, \Omega, y) = d(f, \Omega_1, y) + d(f, \Omega_2, y)$  whenever  $\Omega_1$  and  $\Omega_2$  are disjoint open subsets of  $\Omega$  such that  $y \notin f(\overline{\Omega} \setminus (\Omega_1 \cup \Omega_2))$ .

(d3)  $d(h(t, \cdot), \Omega, y(t))$  is independent of  $t$  whenever  $h : [0, 1] \times \overline{\Omega} \rightarrow \mathbb{R}^n$  and  $y : [0, 1] \rightarrow \mathbb{R}^n$  are continuous and  $y(t) \notin f(t, \partial\Omega)$  for every  $t \in [0, 1]$ .

(d4)  $d(f, \Omega, y) \neq 0$  implies  $f^{-1}(y) \neq \emptyset$ .

Next we denote by  $d(H(\cdot, u^{n-1}, \mu), B_r, 0)$  the topological degree of the application  $H(\cdot, u^{n-1}, \mu)$  with respect to the ball  $B_r$  and the right-hand side 0. For  $\mu = 0$  the system  $H(\cdot, u^{n-1}, 0) = 0$  reduces to a linear system with a positive definite matrix. Applying property (d1) in Proposition 1, we obtain

$$d(H(\cdot, u^{n-1}, 0), B_r, 0) = 1.$$

Then, in view of the homotopy invariance of the topological degree (property (d3) in Proposition 1) we have that

$$d(H(\cdot, u^{n-1}, \mu), B_r, \mu) = d(H(\cdot, u^{n-1}, 0), B_r, 0) \quad \text{for all } \mu \in [0, 1].$$

As a result, in the case where  $\mu = 1$  we have that

$$d(H(\cdot, u^{n-1}, 1), B_r, 1) = 1.$$

Thus, by the property (d4) in Proposition 1, the system  $H(\cdot, u^{n-1}, 1)$  is invertible. Then there exists  $u^n$  such that  $H(\cdot, u^{n-1}, 1) = 0$  so that  $u_{\mathcal{D}}^{\delta t} = (u^n)_{n \in \{1, \dots, N\}}$  is a solution of the discrete problem  $(P_{\mathcal{D}, \delta t})$ .

## 5 Estimates on space and time translates

In this section, we perform estimates on time and space translates of the discrete saturation.

### 5.1 Estimates on time translates

Let  $\lceil s \rceil$  denote the smallest integer larger or equal to  $s$ . We state without proof two technical lemmas deduced from [13], which will be useful for proving the estimate on time translates.

**Lemma 6** Let  $T > 0, \tau > 0, \delta t > 0$  and  $N$  be a positive integer such that  $\tau \in (0, T)$  as well as  $\delta t = T/N$ . Let  $(\gamma^n)_{n \in \mathbb{N}}$  be a family of non-negative real values. Then

$$\int_0^{T-\tau} \sum_{n=\lceil t/\delta t \rceil+1}^{\lceil (t+\tau)/\delta t \rceil} \gamma^n dt \leq \tau \sum_{n=1}^N \gamma^n.$$

**Lemma 7** Let  $T > 0, \tau > 0, \zeta > 0, \delta t > 0$  and  $N$  be a positive integer such that  $\zeta \in [0, \tau], \tau \in (0, T)$  as well as  $\delta t = T/N$ . Let  $(\gamma^n)_{n \in \mathbb{N}}$  be a family of non-negative real values. Then

$$\int_0^{T-\tau} \sum_{n=\lceil t/\delta t \rceil+1}^{\lceil (t+\tau)/\delta t \rceil} \gamma^{\lceil (t+\zeta)/\delta t \rceil} dt \leq \tau \sum_{n=1}^N \gamma^n.$$



**Lemma 8** Let  $u_{\mathcal{D}}^{\delta t}$  be a solution of Problem  $(P_{\mathcal{D},\delta t})$ . There exists a positive constant  $C_7$  only depending on  $\mu$  such that for all  $\tau \in (0, T)$ , there holds

$$\int_0^{T-\tau} \int_{\Omega} \left( c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t + \tau)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)) \right)^2 d\mathbf{x} dt \leq C_7 \tau.$$

*Proof:* Let  $p_t = \lceil \frac{t+\tau}{\delta t} \rceil$  and  $q_t = \lceil \frac{t}{\delta t} \rceil$ , we obtain

$$\int_0^{T-\tau} \int_{\Omega} \left( c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t + \tau)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)) \right)^2 d\mathbf{x} dt = \int_0^{T-\tau} \sum_{K \in \mathcal{M}} |K| \left( c(u_K^{p_t}) - c(u_K^{q_t}) \right)^2 dt. \quad (5.1)$$

Since  $c$  is monotone and Lipschitz continuous, we deduce that

$$\int_0^{T-\tau} \sum_{K \in \mathcal{M}} |K| \left( c(u_K^{p_t}) - c(u_K^{q_t}) \right)^2 dt \leq \int_0^{T-\tau} L_c \sum_{K \in \mathcal{M}} |K| \left( c(u_K^{p_t}) - c(u_K^{q_t}) \right) (u_K^{p_t} - u_K^{q_t}) dt. \quad (5.2)$$

We substitute  $v = u^{p_t} - u^{q_t} \in X_{\mathcal{D},0}$  in the scheme (2.3) to obtain

$$\begin{aligned} & \sum_{K \in \mathcal{M}} |K| \left( c(u_K^n) - c(u_K^{n-1}) \right) (u_K^{p_t} - u_K^{q_t}) \\ &= -\delta t \langle u^n, u^{p_t} - u^{q_t} \rangle_F - \delta t \langle u^n, u^{p_t} - u^{q_t} \rangle_Q. \end{aligned} \quad (5.3)$$

At first, we consider the first term on the right-hand side of (5.3). Applying Hölder's inequality yields

$$\begin{aligned} |\langle u^n, u^{p_t} - u^{q_t} \rangle_F| &\leq |\langle u^n, u^{p_t} \rangle_F| + |\langle u^n, u^{q_t} \rangle_F| \\ &\leq \bar{K} \|\nabla_{\mathcal{D}} u^n\|_{L^2(\Omega)^d} \|\nabla_{\mathcal{D}} u^{p_t}\|_{L^2(\Omega)^d} + \bar{K} \|\nabla_{\mathcal{D}} u^n\|_{L^2(\Omega)^d} \|\nabla_{\mathcal{D}} u^{q_t}\|_{L^2(\Omega)^d}. \end{aligned} \quad (5.4)$$

Since  $2ab \leq a^2 + b^2$  and in view of Lemma 2, one has

$$\begin{aligned} & \|\nabla_{\mathcal{D}} u^{p_t}\|_{L^2(\Omega)^d} \|\nabla_{\mathcal{D}} u^n\|_{L^2(\Omega)^d} + \|\nabla_{\mathcal{D}} u^{q_t}\|_{L^2(\Omega)^d} \|\nabla_{\mathcal{D}} u^n\|_{L^2(\Omega)^d} \\ &\leq \frac{1}{2} \|\nabla_{\mathcal{D}} u^{p_t}\|_{L^2(\Omega)^d}^2 + \frac{1}{2} \|\nabla_{\mathcal{D}} u^{q_t}\|_{L^2(\Omega)^d}^2 + \|\nabla_{\mathcal{D}} u^n\|_{L^2(\Omega)^d}^2 \\ &\leq \frac{C_3}{2} |u^{p_t}|_{X_{\mathcal{D}}}^2 + \frac{C_3}{2} |u^{q_t}|_{X_{\mathcal{D}}}^2 + C_3 |u^n|_{X_{\mathcal{D}}}^2. \end{aligned}$$

Next, we consider the second term in the right-hand side of (5.3); in view of (3.10) with  $\varepsilon_2 = 1$

$$\begin{aligned} |\langle u^n, u^{p_t} - u^{q_t} \rangle_Q| &\leq |\langle u^n, u^{p_t} \rangle_Q| + |\langle u^n, u^{q_t} \rangle_Q| \\ &\leq d |\Omega| \bar{k}_r^{-2} \bar{K} + \frac{1}{2} \bar{K} |u^{p_t}|_{X_{\mathcal{D}}}^2 + \frac{1}{2} \bar{K} |u^{q_t}|_{X_{\mathcal{D}}}^2. \end{aligned} \quad (5.5)$$

Taking the sum of (5.3) with respect to  $n$  from  $q_t + 1$  to  $p_t$  and substituting (5.4) - (5.5) yields

$$\begin{aligned} & \sum_{K \in \mathcal{M}} |K| \left( c(u_K^{p_t}) - c(u_K^{q_t}) \right) (u_K^{p_t} - u_K^{q_t}) \\ &\leq \delta t \bar{K} \left( \sum_{n=q_t+1}^{p_t} \frac{C_3+1}{2} |u^{p_t}|_{X_{\mathcal{D}}}^2 + \sum_{n=q_t+1}^{p_t} \frac{C_3+1}{2} |u^{q_t}|_{X_{\mathcal{D}}}^2 + \sum_{n=q_t+1}^{p_t} (C_3 |u^n|_{X_{\mathcal{D}}}^2 + d |\Omega| \bar{k}_r^{-2}) \right). \end{aligned} \quad (5.6)$$

In view of estimates (5.2)-(5.6) and lemmas 6, 7 and 5, (5.2) becomes

$$\begin{aligned} & \int_0^{T-\tau} \sum_{K \in \mathcal{M}} |K| \left( c(u_K^{p_t}) - c(u_K^{q_t}) \right)^2 dt \\ & \leq \tau L_c \bar{K} \sum_{n=1}^N \delta t \left( \frac{C_3+1}{2} |u^n|_{X_{\mathcal{D}}}^2 + \frac{C_3+1}{2} |u^n|_{X_{\mathcal{D}}}^2 + C_3 |u^n|_{X_{\mathcal{D}}}^2 + d |\Omega| \bar{k}_r^{-2} \right) \\ & \leq \tau L_c \bar{K} \left( (2C_3+1)C_5 + dT |\Omega| \bar{k}_r^{-2} \right). \end{aligned}$$

We obtain

$$\int_0^{T-\tau} \int_{\Omega} \left( c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t + \tau)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)) \right)^2 d\mathbf{x} dt \leq C_7 \tau.$$

where  $C_7 = L_c \bar{K} \left( (2C_3+1)C_5 + dT |\Omega| \bar{k}_r^{-2} \right)$ .  $\square$

We remark that there also holds that

$$\int_{T-\tau}^T \sum_{K \in \mathcal{M}} |K| \left( c(u_K^{p_t}) - c(u_K^{q_t}) \right)^2 dt \leq c_{\max}^2 |\Omega| \tau. \quad (5.7)$$

## 5.2 Estimates on space translates

In this section we prove an estimate in the  $L^2$  norm of differences of space translates of the discrete saturation. At first we state without proof the following result from [1].

**Lemma 9** *Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2 and let  $\eta > 0$  be such that  $\eta \leq d_{K,\sigma}/d_{L,\sigma} \leq 1/\eta$  for all  $\sigma \in \mathcal{E}_{\text{int}}$ , where  $\mathcal{M}_{\sigma} = \{K, L\}$ . There exist  $q > 2$  and  $C_8 > 0$  only depending on  $d, \Omega$  and  $\eta$  such that*

$$\|\Pi_{\mathcal{M}} w(\mathbf{x} + \mathbf{y}) - \Pi_{\mathcal{M}} w(\mathbf{x})\|_{L^2(\mathbb{R}^d)} \leq C_8 |\mathbf{y}|^{\rho} \|\Pi_{\mathcal{M}} w\|_{1,2,\mathcal{M}},$$

where  $\rho = \frac{1}{2} \frac{q-2}{q-1}$ ,  $w \in X_{\mathcal{D}}^0$ ,  $w = 0$  outside  $Q_T$  and  $\|\cdot\|_{1,2,\mathcal{M}}$  is defined by (2.21).

We will apply the result of Lemma 9 to  $\Pi_{\mathcal{M}}^{\delta t} \tilde{u}_{\mathcal{D}}^{\delta t}$ . We first extend  $\Pi_{\mathcal{M}}^{\delta t} \tilde{u}_{\mathcal{D}}^{\delta t}$  by 0 outside  $Q_T$  and extend  $\Pi_{\mathcal{M}} P_{\mathcal{D}} \hat{u}$  by the boundary value outside  $Q_T$ . Figure 3 shows how we proceed.

**Lemma 10** *Let  $\mathcal{D}$  be a discretization of  $\Omega$  in the sense of Definition 2 and let  $\eta > 0$  be such that  $\eta \leq d_{K,\sigma}/d_{L,\sigma} \leq 1/\eta$  for all  $\sigma \in \mathcal{E}_{\text{int}}$ , where  $\mathcal{M}_{\sigma} = \{K, L\}$ . There exist  $q > 2$  and  $C_9$  only depending on  $d, \Omega$  and  $\eta$  such that*

$$\|c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t))\|_{L^2(\mathbb{R}^d)} \leq C_9 |\mathbf{y}|^{\rho} \|\Pi_{\mathcal{M}} \tilde{u}^n\|_{1,2,\mathcal{M}} \quad \forall t \in (t_{n-1}, t_n], \forall \mathbf{y} \in \mathbb{R}^d. \quad (5.8)$$

*Proof:* From the Lipschitz continuity of  $c$  we have

$$\|c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t))\|_{L^2(\mathbb{R}^d)} \leq L_c \|\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t) - \Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)\|_{L^2(\mathbb{R}^d)}.$$

We remark that  $\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t) - \Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t) = \Pi_{\mathcal{M}}^{\delta t} \tilde{u}_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t) - \Pi_{\mathcal{M}}^{\delta t} \tilde{u}_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)$ . Applying Lemma 9, we deduce that there exists  $\rho > 0$  such that

$$\|\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t) - \Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)\|_{L^2(\mathbb{R}^d)} \leq C_8 |\mathbf{y}|^{\rho} \|\Pi_{\mathcal{M}} \tilde{u}^n\|_{1,2,\mathcal{M}}. \quad (5.9)$$

The inequality (5.8) follows from (5.9) by setting  $C_9 = C_8 L_c$ .  $\square$

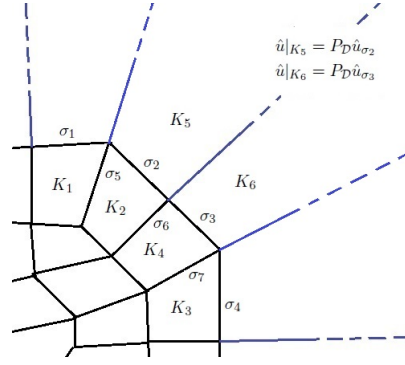


Fig. 3 Extension of function  $\hat{u}$ .

**Theorem 2** Let  $\mathcal{F}$  be a family of discretizations of  $\Omega$  in the sense of Definition 2 such that there exists  $\mu \geq \mu_{\mathcal{D}}$  for all  $\mathcal{D} \in \mathcal{F}$  and let  $\delta t \in (0, 1)$ . The family  $(c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}))_{\mathcal{D} \in \mathcal{F}}$  of approximate saturations is relatively compact in  $L^2(Q_T)$ . In particular, there exists a subsequence of  $\{c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\}$ , which we denote again by  $\{c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\}$ , and a function  $\vartheta \in L^2(Q_T)$  such that  $\{c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\}$  converges strongly to  $\vartheta$  in  $L^2(Q_T)$  as  $l_{\mathcal{D}}$  and  $\delta t$  tend to zero.

*Proof:* In view of Lemma 10, integrating in time we obtain

$$\|c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t))\|_{L^2(\mathbb{R}^d \times (0, T))}^2 \leq C_9^2 |\mathbf{y}|^{2p} \sum_{n=1}^N \delta t \|\Pi_{\mathcal{M}} \bar{u}^n\|_{1,2,\mathcal{M}}^2,$$

which by the inequality (2.22) and Lemma 5 yields

$$\|c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t))\|_{L^2(\mathbb{R}^d \times (0, T))} \leq \sqrt{C_5} C_9 |\mathbf{y}|^p.$$

We combine this result with Lemma 8 to obtain

$$\begin{aligned} & \|c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t + \tau)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t))\|_{L^2(\mathbb{R}^d \times (0, T))} \\ & \leq \|c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t + \tau)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t))\|_{L^2(\mathbb{R}^d \times (0, T))} \\ & \quad + \|c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x} + \mathbf{y}, t)) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t))\|_{L^2(\mathbb{R}^d \times (0, T))} \\ & \leq C_{10} (|\tau|^{1/2} + |\mathbf{y}|^p), \end{aligned}$$

where  $C_{10} = \max(\sqrt{C_5} C_9, \sqrt{C_7})$ .

Moreover we recall that  $\|c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\|_{L^2(\mathbb{R}^d \times (0, T))}^2 \leq |\Omega| T \|c\|_{L^\infty(\mathbb{R})}^2$ . Applying the Fréchet-Kolmogorov compactness theorem we deduce that the sequence  $\{c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\}$  is relatively compact in  $L^2(Q_T)$ . Thus, there exists a  $\vartheta \in L^2(Q_T)$  and a subsequence of  $\{c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\}$  which converges to  $\vartheta$  strongly in  $L^2(Q_T)$  as  $l_{\mathcal{D}}$  and  $\delta t$  tend to zero.  $\square$

**Lemma 11** Let  $\mathcal{F}$  be a family of discretizations of  $\Omega$  in the sense of Definition 2 and let  $\delta t \in (0, 1)$ . Let  $(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})_{\mathcal{D} \in \mathcal{F}}$  be a sequence of approximate solutions of Problem  $(P_{\mathcal{D}, \delta t})$  such that  $\{\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}\}$  converges to  $\bar{u}$  weakly in  $L^2(Q_T)$  and  $\{c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\}$  be a sequence of approximate saturations which converges to a limit  $\vartheta$  strongly in  $L^2(Q_T)$  and a.e. in  $Q_T$  as  $l_{\mathcal{D}}$  and  $\delta t$  tend to zero. Then  $\vartheta = c(\bar{u})$ .

*Proof:* We deduce from the monotonicity of  $c$  that for all  $\phi \in L^2(Q_T)$

$$\int_0^T \int_{\Omega} \left( c(\phi) - c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}) \right) \left( \phi - (\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}) \right) dxdt \geq 0.$$

Because of the weak convergence of  $\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}$  and the strong convergence of  $c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})$  respectively, the expression above tends to  $\int_0^T \int_{\Omega} (c(\phi) - \vartheta)(\phi - \bar{u}) dxdt$ . Let  $\delta > 0$  and set  $\phi = \bar{u} + \delta(\vartheta - c(\bar{u}))$ ; we obtain

$$\delta \int_0^T \int_{\Omega} \left( c(\bar{u} + \delta(\vartheta - c(\bar{u}))) - \vartheta \right) (\vartheta - c(\bar{u})) dxdt \geq 0.$$

We divide the inequality above by  $\delta$  and let  $\delta \rightarrow 0$ . This implies that

$$- \int_0^T \int_{\Omega} \left( c(\bar{u}) - \vartheta \right)^2 dxdt \geq 0,$$

so that  $c(\bar{u}) = \vartheta$  a.e. in  $Q_T$ .  $\square$

## 6 Convergence

**Theorem 3** *Let  $\mathcal{F}$  be a family of discretizations of  $\Omega$  in the sense of Definition 2 such that there exists  $\mu \geq \mu_{\mathcal{D}}$  for all  $\mathcal{D} \in \mathcal{F}$ . Let  $\delta t \in (0, 1)$  and let  $(u_{\mathcal{D}}^{\delta t})_{\mathcal{D} \in \mathcal{F}}$  be a family of solution of Problem  $(P_{\mathcal{D}, \delta t})$ . Then*

(i) *There exists a function  $\bar{u} \in L^2(\Omega)$  and a subsequence of  $\{\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}\}$ , which we denote again by  $\{\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}\}$ , which converges to  $\bar{u}$  weakly in  $L^2(Q_T)$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ ;*

(ii) *There exists a subsequence of  $\{c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\}$  which converges to  $c(\bar{u})$  strongly in  $L^2(Q_T)$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ .*

*Moreover  $\bar{u}$  is a weak solution of Problem (P),  $\bar{u} - \hat{u} \in L^2(0, T; H_0^1(\Omega))$  and  $\nabla_{\mathcal{D}}^{\delta t} u_{\mathcal{D}}^{\delta t}$  converges to  $\nabla \bar{u}$  weakly in  $L^2(Q_T)^d$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ .*

*Proof:* Estimate (3.1) together with the discrete Poincaré inequality (2.23) imply that the sequences  $\{\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}\}$  and  $\{\nabla_{\mathcal{D}}^{\delta t} u_{\mathcal{D}}^{\delta t}\}$  are bounded. Thus there exists  $\bar{u}$  in  $L^2(Q_T)$  and a subsequence of  $\{\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}\}$ , which we denote again by  $\{\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t}\}$ , which converges weakly to  $\bar{u}$  in  $L^2(Q_T)$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ .

Let  $\tilde{u}_{\mathcal{D}}^{\delta t} = u_{\mathcal{D}}^{\delta t} - P_{\mathcal{D}} \hat{u}$ . It is easy to see that  $\Pi_{\mathcal{M}} P_{\mathcal{D}} \hat{u}$  converges to  $\hat{u}$  strongly in  $L^2(\Omega)$  as  $l_{\mathcal{D}} \rightarrow 0$ . We deduce that

$$\Pi_{\mathcal{M}}^{\delta t} \tilde{u}_{\mathcal{D}}^{\delta t} \rightharpoonup \bar{u} = \bar{u} - \hat{u} \quad \text{in } L^2(Q_T) \text{ as } l_{\mathcal{D}}, \delta t \rightarrow 0. \quad (6.1)$$

It follows from Lemma 4 that  $\tilde{u} \in L^2(0, T; H_0^1(\Omega))$  and that  $\nabla_{\mathcal{D}}^{\delta t} \tilde{u}_{\mathcal{D}}^{\delta t}$  converges weakly to  $\nabla \tilde{u}$  in  $L^2(Q_T)^d$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ . Moreover  $\nabla_{\mathcal{D}}^{\delta t} P_{\mathcal{D}} \hat{u}$  converges to  $\nabla \hat{u}$  strongly in  $L^2(\Omega)$  as  $l_{\mathcal{D}} \rightarrow 0$  [13, Lemma 4.4]. Thus we deduce that

$$\nabla_{\mathcal{D}}^{\delta t} \tilde{u}_{\mathcal{D}}^{\delta t} + \nabla_{\mathcal{D}}^{\delta t} P_{\mathcal{D}} \hat{u} \rightharpoonup \nabla \tilde{u} + \nabla \hat{u} \quad \text{in } L^2(Q_T) \text{ as } l_{\mathcal{D}}, \delta t \rightarrow 0, \quad (6.2)$$

or else

$$\nabla_{\mathcal{D}}^{\delta t} \tilde{u}_{\mathcal{D}}^{\delta t} \rightharpoonup \nabla \bar{u} \quad \text{in } L^2(Q_T) \text{ as } l_{\mathcal{D}}, \delta t \rightarrow 0, \quad (6.3)$$

and that  $\bar{u} - \hat{u} \in L^2(0, T; H_0^1(\Omega))$ .

By Theorem 2, there exists a function  $\vartheta \in L^2(Q_T)$  and a subsequence of  $\{c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})\}$  such that  $c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})$  converges to  $\vartheta$  strongly in  $L^2(Q_T)$  as  $l_{\mathcal{D}}$  and  $\delta t$  tend to zero. Also applying Lemma 11 we deduce that  $\vartheta = c(\bar{u})$ .

Next we prove that  $\bar{u}$  is a weak solution of Problem (P). We first introduce the function space

$$\Psi = \left\{ \psi \in C^2(\bar{\Omega} \times [0, T]), \psi = 0 \text{ on } \partial\Omega \times [0, T], \psi(\cdot, T) = 0 \right\}. \quad (6.4)$$

Let  $\psi \in \Psi$  and set  $v = P_{\mathcal{D}}\psi(\mathbf{x}, t_{n-1})$  in (2.3). Taking the sum on  $n = \{1, \dots, N\}$ , we obtain  $T_T + T_F + T_Q = 0$ , with

$$\begin{aligned} T_T &= \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| \left( c(u_K^n) - c(u_K^{n-1}) \right) \psi(\mathbf{x}_K, t_{n-1}), \\ T_F &= \sum_{n=1}^N \delta t \int_{\Omega} \mathbf{K} \nabla_{\mathcal{D}} u^n \cdot \nabla_{\mathcal{D}} P_{\mathcal{D}} \psi(\mathbf{x}, t_{n-1}) \, d\mathbf{x}, \\ T_Q &= - \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} k_r(c(u_{K\sigma}^n)) g_{K\sigma} \left( \psi(\mathbf{x}_K, t_{n-1}) - \psi(\mathbf{x}_{\sigma}, t_{n-1}) \right). \end{aligned} \quad (6.5)$$

#### Time evolution term

Let  $p^n = c(u_K^n)$  and  $q^n = \psi(\mathbf{x}_K, t_n)$ . Adding and subtracting  $\sum_{K \in \mathcal{M}} |K| p^N q^N$  in the expression of  $T_T$ , we deduce that

$$\begin{aligned} T_T &= \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| p^n q^{n-1} - \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| p^{n-1} q^{n-1} - \sum_{K \in \mathcal{M}} |K| p^N q^N + \sum_{K \in \mathcal{M}} |K| p^N q^N \\ &= \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| p^n q^{n-1} - \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| p^n q^n - \sum_{K \in \mathcal{M}} |K| p^0 q^0 + \sum_{K \in \mathcal{M}} |K| p^N q^N \\ &= - \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| p^n (q^n - q^{n-1}) - \sum_{K \in \mathcal{M}} |K| p^0 q^0 + \sum_{K \in \mathcal{M}} |K| p^N q^N. \end{aligned} \quad (6.6)$$

As a result, we deduce that  $T_T = A_1 - A_2 - A_3$  where

$$\begin{aligned} A_1 &= \sum_{K \in \mathcal{M}} |K| c(u_K^N) \psi(\mathbf{x}_K, t_N), \\ A_2 &= \sum_{K \in \mathcal{M}} |K| c(u_K^0) \psi(\mathbf{x}_K, 0), \\ A_3 &= \sum_{n=1}^N \sum_{K \in \mathcal{M}} |K| c(u_K^n) \left( \psi(\mathbf{x}_K, t_n) - \psi(\mathbf{x}_K, t_{n-1}) \right). \end{aligned}$$

Since  $\psi(\mathbf{x}, t_N) = \psi(\mathbf{x}, T) = 0$ , the first term  $A_1$  vanishes.

Next we add and subtract  $\sum_{K \in \mathcal{M}} \int_K c(u_K^0) \psi(\mathbf{x}, 0) d\mathbf{x}$  to the term  $A_2$  and compare  $A_2$  with  $\int_{\Omega} c(u_0(\mathbf{x})) \psi(\mathbf{x}, 0) d\mathbf{x}$ . This yields

$$\begin{aligned} A_2 - \int_{\Omega} c(u_0(\mathbf{x})) \psi(\mathbf{x}, 0) d\mathbf{x} &= \sum_{K \in \mathcal{M}} |K| c(u_K^0) \psi(\mathbf{x}_K, 0) d\mathbf{x} - \sum_{K \in \mathcal{M}} \int_K c(u_K^0) \psi(\mathbf{x}, 0) d\mathbf{x} \\ &\quad + \sum_{K \in \mathcal{M}} \int_K c(u_K^0) \psi(\mathbf{x}, 0) d\mathbf{x} - \int_{\Omega} c(u_0(\mathbf{x})) \psi(\mathbf{x}, 0) d\mathbf{x} \\ &= \sum_{K \in \mathcal{M}} \int_K c(u_K^0) (\psi(\mathbf{x}_K, 0) - \psi(\mathbf{x}, 0)) d\mathbf{x} \\ &\quad + \sum_{K \in \mathcal{M}} \int_K (c(u_K^0) - c(u_0(\mathbf{x}))) \psi(\mathbf{x}, 0) d\mathbf{x}. \end{aligned}$$

Applying hypothesis  $(\mathcal{H}_1)$ , we deduce that

$$\begin{aligned} |A_2 - \int_{\Omega} c(u_0(\mathbf{x})) \psi(\mathbf{x}, 0) d\mathbf{x}| &\leq \|c\|_{L^\infty(\mathbb{R})} \sum_{K \in \mathcal{M}} \int_K |\psi(\mathbf{x}_K, 0) - \psi(\mathbf{x}, 0)| d\mathbf{x} \\ &\quad + L_c \sum_{K \in \mathcal{M}} \int_K |u_K^0 - u_0(\mathbf{x})| |\psi(\mathbf{x}, 0)| d\mathbf{x}. \end{aligned} \quad (6.7)$$

Since  $\psi \in C^2(\overline{\Omega} \times [0, T])$ , there exists a positive constant  $C_1^\psi$ , which only depends on  $\psi, T$  and  $\Omega$ , such that

$$\sum_{K \in \mathcal{M}} \int_K |\psi(\mathbf{x}_K, 0) - \psi(\mathbf{x}, 0)| \leq \Omega C_1^\psi l_{\mathcal{D}}.$$

We conclude that the first term on the right-hand side of (6.7) converges to zero as  $l_{\mathcal{D}}$  tends to zero. By the definition of  $u_K^0$  in (2.2), the second term on the right-hand side of (6.7) tends to zero as  $l_{\mathcal{D}}$  tends to zero. Finally  $A_2 \rightarrow \int_{\Omega} c(u_0(\mathbf{x})) \psi(\mathbf{x}, 0) d\mathbf{x}$  as  $l_{\mathcal{D}}, \delta t$  tend to zero.

Next, we add and subtract  $\sum_{n=1}^N \sum_{K \in \mathcal{M}} \int_{t_{n-1}}^{t_n} \int_K c(u_K^n) \partial_t \psi d\mathbf{x} dt$  to the difference

$$A_3 - \int_0^T \int_{\Omega} c(\bar{u}(\mathbf{x}, t)) \partial_t \psi d\mathbf{x} dt,$$

to deduce that

$$\begin{aligned}
A_3 &= \int_0^T \int_{\Omega} c(\bar{u}(\mathbf{x}, t)) \partial_t \psi \, d\mathbf{x} dt \\
&= \sum_{n=1}^N \sum_{K \in \mathcal{M}} \int_K c(u_K^n) \left( \psi(\mathbf{x}_K, t_n) - \psi(\mathbf{x}_K, t_{n-1}) \right) d\mathbf{x} - \sum_{n=1}^N \sum_{K \in \mathcal{M}} \int_{t_{n-1}}^{t_n} \int_K c(u_K^n) \partial_t \psi \, d\mathbf{x} dt \\
&\quad + \sum_{n=1}^N \sum_{K \in \mathcal{M}} \int_{t_{n-1}}^{t_n} \int_K c(u_K^n) \partial_t \psi \, d\mathbf{x} dt - \int_0^T \int_{\Omega} c(\bar{u}(\mathbf{x}, t)) \partial_t \psi \, d\mathbf{x} dt \\
&= \sum_{n=1}^N \int_{t_{n-1}}^{t_n} \int_{\Omega} c(u_K^n) \left( \partial_t \psi(\mathbf{x}_K, t) - \partial_t \psi(\mathbf{x}, t) \right) d\mathbf{x} dt \\
&\quad + \sum_{n=1}^N \sum_{K \in \mathcal{M}} \int_{t_{n-1}}^{t_n} \int_K \left( c(u_K^n) - c(\bar{u}(\mathbf{x}, t)) \right) \partial_t \psi(\mathbf{x}, t) \, d\mathbf{x} dt.
\end{aligned}$$

Thus

$$\begin{aligned}
|A_3 - \int_0^T \int_{\Omega} c(\bar{u}(\mathbf{x}, t)) \partial_t \psi \, d\mathbf{x} dt| & \\
&\leq \sum_{n=1}^N \sum_{K \in \mathcal{M}} \int_{t_{n-1}}^{t_n} \int_K |c(u_K^n)| |\partial_t \psi(\mathbf{x}_K, t) - \partial_t \psi(\mathbf{x}, t)| \, d\mathbf{x} dt \quad (6.8) \\
&\quad + \sum_{n=1}^N \sum_{K \in \mathcal{M}} \int_{t_{n-1}}^{t_n} \int_K |c(u_K^n) - c(\bar{u}(\mathbf{x}, t))| |\partial_t \psi(\mathbf{x}, t)| \, d\mathbf{x} dt.
\end{aligned}$$

For all  $\mathbf{x} \in K$  and for all  $K \in \mathcal{M}$ , we have

$$|\partial_t \psi(\mathbf{x}_K, t) - \partial_t \psi(\mathbf{x}, t)| \leq C_2^{\psi} l_{\mathcal{Q}},$$

where  $C_2^{\psi}$  is a positive constant. Since  $c$  is bounded, the first term on the right-hand side of (6.8) tends to zero as  $l_{\mathcal{Q}}, \delta t$  tend to zero. Since  $c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{Q}}^{\delta t})$  strongly converges to  $c(\bar{u})$  in  $L^2(Q_T)$ , the second term tends to zero as  $l_{\mathcal{Q}}, \delta t$  tend to zero. Thus  $A_3 \rightarrow \int_0^T \int_{\Omega} c(\bar{u}(\mathbf{x}, t)) \partial_t \psi \, d\mathbf{x} dt$  as  $l_{\mathcal{Q}}, \delta t$  tend to zero. We conclude that

$$T_T \rightarrow - \int_0^T \int_{\Omega} c(\bar{u}(\mathbf{x}, t)) \partial_t \psi \, d\mathbf{x} dt - \int_{\Omega} c(u_0(\mathbf{x})) \psi(\mathbf{x}, 0) \, d\mathbf{x} \quad \text{as } l_{\mathcal{Q}}, \delta t \rightarrow 0. \quad (6.9)$$

### Diffusion term

Next, we consider the diffusion term  $T_F$  in (6.5). Adding and subtracting  $\int_0^T \int_{\Omega} \mathbf{K}(\mathbf{x}) \nabla_{\mathcal{Q}}^{\delta t} u_{\mathcal{Q}}^{\delta t}(\mathbf{x}, t) \cdot \nabla \psi(\mathbf{x}, t) \, d\mathbf{x} dt$  yields

$$\begin{aligned}
T_F &= \int_0^T \int_{\Omega} \mathbf{K}(\mathbf{x}) \nabla \bar{u}(\mathbf{x}, t) \cdot \nabla \psi(\mathbf{x}, t) \, d\mathbf{x} dt \\
&= \sum_{n=1}^N \delta t \int_{\Omega} \mathbf{K} \nabla_{\mathcal{Q}} u^n \cdot \nabla_{\mathcal{Q}} P_{\mathcal{Q}} \psi(\mathbf{x}, t_{n-1}) \, d\mathbf{x} - \int_0^T \int_{\Omega} \mathbf{K}(\mathbf{x}) \nabla_{\mathcal{Q}}^{\delta t} u_{\mathcal{Q}}^{\delta t}(\mathbf{x}, t) \cdot \nabla \psi(\mathbf{x}, t) \, d\mathbf{x} dt \\
&\quad + \int_0^T \int_{\Omega} \mathbf{K}(\mathbf{x}) \nabla_{\mathcal{Q}}^{\delta t} u_{\mathcal{Q}}^{\delta t}(\mathbf{x}, t) \cdot \nabla \psi(\mathbf{x}, t) \, d\mathbf{x} dt - \int_0^T \int_{\Omega} \mathbf{K}(\mathbf{x}) \nabla \bar{u}(\mathbf{x}, t) \cdot \nabla \psi(\mathbf{x}, t) \, d\mathbf{x} dt \\
&= \sum_{n=1}^N \int_{t_{n-1}}^{t_n} \int_{\Omega} \mathbf{K}(\mathbf{x}) \nabla_{\mathcal{Q}} u^n(\mathbf{x}, t) \cdot \left( \nabla_{\mathcal{Q}} P_{\mathcal{Q}} \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t) \right) \, d\mathbf{x} dt \\
&\quad + \int_0^T \int_{\Omega} \mathbf{K}(\mathbf{x}) \left( \nabla_{\mathcal{Q}}^{\delta t} u_{\mathcal{Q}}^{\delta t}(\mathbf{x}, t) - \nabla \bar{u}(\mathbf{x}, t) \right) \cdot \nabla \psi(\mathbf{x}, t) \, d\mathbf{x} dt. \quad (6.10)
\end{aligned}$$

Since  $\nabla_{\mathcal{D}}^{\delta t} u_{\mathcal{D}}^{\delta t}(\mathbf{x}, t)$  weakly converges to  $\nabla \bar{u}$  in  $L^2(Q_T)$ , the second term on the right-hand side of (6.10) tends to zero as  $l_{\mathcal{D}}, \delta t$  tend to zero.

We denote by  $\tilde{T}_F$  the first term on the right-hand side of (6.10). We have that

$$\begin{aligned} |\tilde{T}_F| &= \sum_{n=1}^N \int_{t_{n-1}}^{t_n} \int_{\Omega} \mathbf{K}(\mathbf{x}) \nabla_{\mathcal{D}} u^n(\mathbf{x}, t) \cdot \left( \nabla_{\mathcal{D}} P_{\mathcal{D}} \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t) \right) d\mathbf{x} dt \\ &\leq \sum_{n=1}^N \int_{t_{n-1}}^{t_n} \bar{K} \|\nabla_{\mathcal{D}} u^n(\mathbf{x}, t)\|_{L^2(\Omega)^d} \|\nabla_{\mathcal{D}} P_{\mathcal{D}} \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t)\|_{L^2(\Omega)^d}, \end{aligned} \quad (6.11)$$

together with

$$\begin{aligned} \|\nabla_{\mathcal{D}} P_{\mathcal{D}} \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t)\|_{L^\infty(\Omega)^d} &\leq \|\nabla_{\mathcal{D}} P_{\mathcal{D}} \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t_{n-1})\|_{L^\infty(\Omega)^d} \\ &\quad + \|\nabla \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t)\|_{L^\infty(\Omega)^d}. \end{aligned} \quad (6.12)$$

In view of the regularity of  $\psi$  there holds

$$\|\nabla \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t)\|_{L^\infty(\Omega)^d} \leq C_3^\psi \delta t,$$

where  $C_3^\psi$  is a constant.

In view of [9, Lemma 4.4], we have  $\|\nabla_{\mathcal{D}} P_{\mathcal{D}} \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t_{n-1})\|_{L^\infty(\Omega)^d} \leq C_4^\psi$  where  $C_4^\psi$  is a positive constant. As a result, the term  $\|\nabla_{\mathcal{D}} P_{\mathcal{D}} \psi(\mathbf{x}, t_{n-1}) - \nabla \psi(\mathbf{x}, t)\|_{L^\infty(\Omega)^d}$  tends to 0 as  $l_{\mathcal{D}}, \delta t$  tend to zero. In view of Lemma 2 and estimate (3.2), the first term on the right-hand side of (6.10) tends to zero as  $l_{\mathcal{D}}, \delta t$  tend to zero. We conclude that

$$T_F \rightarrow \int_0^T \int_{\Omega} \mathbf{K}(\mathbf{x}) \nabla \bar{u}(\mathbf{x}, t) \cdot \nabla \psi(\mathbf{x}, t) d\mathbf{x} dt \quad \text{as } l_{\mathcal{D}}, \delta t \rightarrow 0. \quad (6.13)$$

### Convection term

Finally we prove that the convection term  $T_Q$  tends to  $-\int_0^T \int_{\Omega} k_r(c(\bar{u})) \mathbf{K}(\mathbf{x}) \nabla z \cdot \nabla \psi(\mathbf{x}, t) d\mathbf{x} dt$  as  $l_{\mathcal{D}}, \delta t \rightarrow 0$ . For this purpose, we introduce the following two terms

$$\begin{aligned} T_Q^1 &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} k_r(c(u_K^n)) g_{K\sigma} \psi(\mathbf{x}_\sigma, t_{n-1}), \\ T_Q^2 &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} k_r(c(u_K^n)) g_{K\sigma} \psi(\mathbf{x}_K, t_{n-1}). \end{aligned} \quad (6.14)$$

We show below that  $\lim_{l_{\mathcal{D}}, \delta t \rightarrow 0} |T_Q - (T_Q^2 - T_Q^1)| = 0$ .

$$\begin{aligned} T_Q - (T_Q^2 - T_Q^1) &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} k_r(c(u_{K\sigma}^n)) g_{K\sigma} \left( \psi(\mathbf{x}_K, t_{n-1}) - \psi(\mathbf{x}_\sigma, t_{n-1}) \right) \\ &\quad - \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} k_r(c(u_K^n)) g_{K\sigma} \left( \psi(\mathbf{x}_K, t_{n-1}) - \psi(\mathbf{x}_\sigma, t_{n-1}) \right) \\ &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} g_{K\sigma} \left( \psi(\mathbf{x}_K, t_{n-1}) - \psi(\mathbf{x}_\sigma, t_{n-1}) \right) \left( k_r(c(u_{K\sigma}^n)) - k_r(c(u_K^n)) \right). \end{aligned} \quad (6.15)$$



We denote by  $\tilde{T}_Q$  the term on the right-hand side of (6.15). Since  $|g_{K,\sigma}| \leq \bar{K}|\sigma|$ , using the Cauchy-Schwarz inequality, we deduce from the Lipschitz continuity of the functions  $k_r$  and  $c$  that

$$\begin{aligned} (\tilde{T}_Q)^2 &\leq \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} |\sigma| d_{K\sigma} \bar{K} \left( \psi(\mathbf{x}_K, t_{n-1}) - \psi(\mathbf{x}_\sigma, t_{n-1}) \right)^2 \\ &\quad \cdot \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} \frac{|\sigma|}{d_{K\sigma}} \left( k_r(c(u_{K\sigma}^n)) - k_r(c(u_K^n)) \right)^2 \\ &\leq \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} d |D_{K\sigma}| \bar{K} |\psi(\mathbf{x}_K, t_{n-1}) - \psi(\mathbf{x}_\sigma, t_{n-1})|^2 \\ &\quad \cdot \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} \sum_{\sigma \in \mathcal{E}_K} L_c^2 L_{k_r}^2 \frac{|\sigma|}{d_{K\sigma}} (u_{K\sigma}^n - u_K^n)^2. \end{aligned} \quad (6.16)$$

It follows from the definition (2.17) that  $(u_{K\sigma}^n - u_K^n)^2 \leq (u_\sigma^n - u_K^n)^2$ , which together with Definition 4 and estimate (3.2) implies that the second term on the right-hand side of (6.16) is bounded. In view of the regularity properties of  $\psi$  we deduce that  $(\tilde{T}_Q)^2 \leq C_4^\psi l_{\mathcal{D}}^2$ . As a result, we have

$$\lim_{l_{\mathcal{D}}, \delta t \rightarrow 0} |T_Q - (T_Q^2 - T_Q^1)| = 0. \quad (6.17)$$

Next we consider term  $T_Q^1$ . Because of the regularity of  $\psi$ , it is easy to see that  $(T_Q^1 - \bar{T}_Q^1) \rightarrow 0$  as  $l_{\mathcal{D}} \rightarrow 0$  where

$$\begin{aligned} \bar{T}_Q^1 &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} k_r(c(u_K^n)) \sum_{\sigma \in \mathcal{E}_K} \int_{\sigma} \mathbf{K} \nabla_z \psi \mathbf{n}_{K\sigma} d\gamma \\ &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} k_r(c(u_K^n)) \int_K \operatorname{div}(\mathbf{K}(\mathbf{x}) \nabla_z \psi(\mathbf{x}, t)) d\mathbf{x} \\ &= \int_0^T \int_{\Omega} k_r(c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})) \operatorname{div}(\mathbf{K}(\mathbf{x}) \nabla_z \psi) d\mathbf{x} dt. \end{aligned}$$

Since  $c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})$  converges to  $c(\bar{u})$  strongly in  $L^2(Q_T)$ , we have that

$$T_Q^1 \rightarrow \int_0^T \int_{\Omega} k_r(c(\bar{u})) \operatorname{div}(\mathbf{K}(\mathbf{x}) \nabla_z \psi) d\mathbf{x} dt \quad \text{as } l_{\mathcal{D}}, \delta t \rightarrow 0. \quad (6.18)$$

Next we consider term  $T_Q^2$ .

$$\begin{aligned} T_Q^2 &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} k_r(c(u_K^n)) \psi(\mathbf{x}_K, t_{n-1}) \sum_{\sigma \in \mathcal{E}_K} \int_{\sigma} \mathbf{K} \nabla_z \mathbf{n}_{K\sigma} d\gamma \\ &= \sum_{n=1}^N \delta t \sum_{K \in \mathcal{M}} k_r(c(u_K^n)) \psi(\mathbf{x}_K, t_{n-1}) \int_K \operatorname{div}(\mathbf{K}(\mathbf{x}) \nabla_z) d\mathbf{x} \\ &= \int_0^T \int_{\Omega} k_r(c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})) \operatorname{div}(\mathbf{K}(\mathbf{x}) \nabla_z) \bar{\Pi}_{\mathcal{M}}^{\delta t} P_{\mathcal{D}} \psi d\mathbf{x} dt, \end{aligned}$$

where  $\bar{\Pi}_{\mathcal{M}}^{\delta t} P_{\mathcal{D}} \psi(\mathbf{x}, t) = \psi(\mathbf{x}_K, t_{n-1})$  for all  $(\mathbf{x}, t) \in K \times [t_{n-1}, t_n]$ . Since  $c(\Pi_{\mathcal{M}}^{\delta t} u_{\mathcal{D}}^{\delta t})$  converges to  $c(\bar{u})$  strongly in  $L^2(Q_T)$  and since  $\bar{\Pi}_{\mathcal{M}}^{\delta t} P_{\mathcal{D}} \psi$  converges to  $\psi$  strongly in  $L^2(Q_T)$ , we deduce that

$$T_Q^2 \rightarrow \int_0^T \int_{\Omega} k_r(c(\bar{u})) \operatorname{div}(\mathbf{K}(\mathbf{x}) \nabla_z) \psi d\mathbf{x} dt \quad \text{as } l_{\mathcal{D}}, \delta t \rightarrow 0. \quad (6.19)$$

We remark that  $\operatorname{div}(\mathbf{K}(\mathbf{x})\nabla z\psi) = \operatorname{div}(\mathbf{K}(\mathbf{x})\nabla z)\psi + \mathbf{K}(\mathbf{x})\nabla z\nabla\psi$ . It follows from (6.17) - (6.19) that

$$T_Q \rightarrow - \int_0^T \int_{\Omega} k_r(c(\bar{u}))\mathbf{K}(\mathbf{x})\nabla z \cdot \nabla\psi \, dxdt \quad \text{as } l_{\mathcal{Q}}, \delta t \rightarrow 0. \quad (6.20)$$

From (6.5), (6.9), (6.13) (6.20), we deduce that  $\bar{u}$  satisfies the weak form (1.5) for test functions  $\psi \in \Psi$ . Finally, we deduce from the density of the set  $\Psi$  in the set

$$\Psi = \left\{ \psi \in L^2(0, T; H_0^1(\Omega)), \partial_t \psi \in L^\infty(Q_T), \psi(\cdot, T) = 0 \right\}. \quad (6.21)$$

that  $\bar{u}$  is a weak solution of the continuous problem (P) in the sense of Definition 1.  $\square$

## 7 Numerical tests

### 7.1 The Hornung-Messing problem

The Hornung-Messing problem is a standard test (cf. for instance [13]). We consider a horizontal flow in a homogeneous ground  $\Omega = [0, 1]^2$  and set  $T = 1$ . Its characteristics are given by

$$\theta(\psi) = \begin{cases} \pi^2/2 - 2\arctan^2(\psi) & \text{if } \psi < 0, \\ \pi^2/2 & \text{otherwise,} \end{cases}$$

$$k_\theta(\psi) = \begin{cases} 2/(1 + \psi)^2 & \text{if } \psi < 0, \\ 2 & \text{otherwise,} \end{cases} \quad \mathbf{K}(\mathbf{x}) = \mathbf{Id}.$$

An analytical solution is given by

$$p(x, z, t) = \begin{cases} -s/2 & \text{if } s < 0, \\ -\tan\left(\frac{e^s - 1}{e^s + 1}\right) & \text{otherwise,} \end{cases} \quad (7.1)$$

where  $s = x - z - t$ . The problem after Kirchhoff's transformation is given by Problem (1.2) with

$$c(u) = \theta(p) = \begin{cases} \pi^2/2 - 2\arctan^2\left(\frac{u}{2-u}\right) & \text{if } p < 0, \\ \pi^2/2 & \text{otherwise,} \end{cases}$$

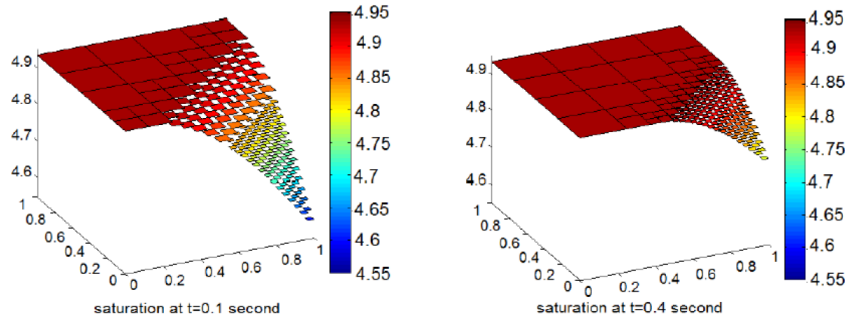
where

$$u(x, z, t) = \begin{cases} \frac{2p(x, z, t)}{1 + p(x, z, t)} & \text{if } p < 0, \\ 2p(x, z, t) & \text{otherwise.} \end{cases} \quad (7.2)$$

We apply the SUSHI scheme using an adaptive mesh driven by the variations of the saturation. We prescribe the Neumann boundary condition deduced from (7.2) on the line  $x = 0$  and an inhomogeneous Dirichlet boundary condition elsewhere. We use an initially square mesh, which is such that each square can be decomposed again into four smaller square elements. Whereas the standard finite volume scheme is not suited to handle such a non-conforming adaptive mesh, the SUSHI scheme is compatible with these non-conforming volume elements.

We introduce the relative error in  $L^2(Q_T)$  between the exact and the numerical solution

$$\operatorname{err}(u) = \frac{\|u_{\text{exact}} - u_{\mathcal{Q}, \delta t}\|_{L^2(Q_T)}}{\|u_{\text{exact}}\|_{L^2(Q_T)}}, \quad (7.3)$$



**Fig. 4** Saturation at  $t = 0.1$  second and at  $t = 0.4$  second. The medium is unsaturated on the right-hand side of the space domain where  $\theta < 4.9348$  and fully saturated elsewhere.

Mesh	$N$	$l_{\mathcal{D}}$	$N_{unk}$	$err(u)$	$err(c(u))$	$eoc(u)$
Uniform	25	0.2	85	$2.40 \cdot 10^{-2}$	$1.60 \cdot 10^{-5}$	-
Uniform	100	0.1	320	$6.09 \cdot 10^{-3}$	$4.13 \cdot 10^{-6}$	1.98
Uniform	400	0.05	1240	$1.53 \cdot 10^{-3}$	$2.90 \cdot 10^{-6}$	2.00
Uniform	1600	0.025	4880	$3.76 \cdot 10^{-3}$	$1.83 \cdot 10^{-6}$	2.02
Adaptive	200	0.143	302	$5.62 \cdot 10^{-3}$	$3.67 \cdot 10^{-6}$	-
Adaptive	800	0.071	1232	$1.32 \cdot 10^{-3}$	$2.19 \cdot 10^{-6}$	-

**Table 1** Number of time steps  $N$ , mesh size  $l_{\mathcal{D}}$ , number of unknowns  $N_{unk}$ , the error on the solution  $err(u)$ , the error on the saturation  $err(c(u))$  and the experimental order of convergence  $eoc(u)$ .

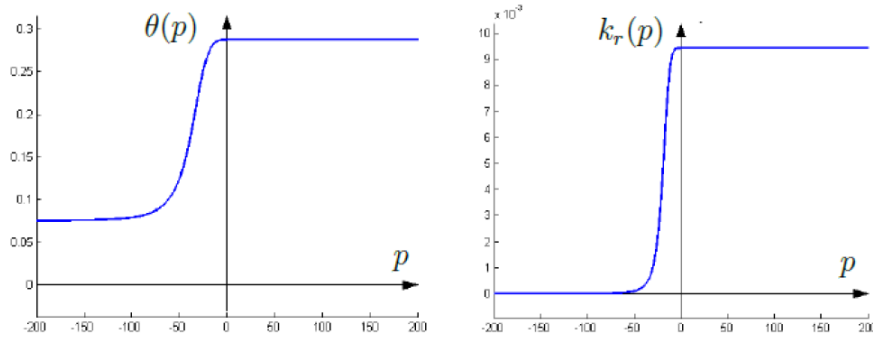
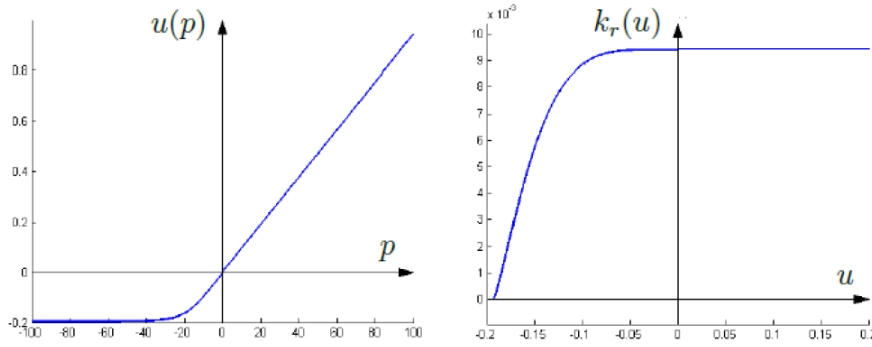
as well as the experimental order of convergence

$$eoc_{i+1}(u) = \frac{\log(err(u_i)/err(u_{i+1}))}{\log(h_{\mathcal{D}_i}/h_{\mathcal{D}_{i+1}})}, \quad (7.4)$$

where  $u_i$  is the solution corresponding to the space discretization  $\mathcal{D}_i$ . Table 1 shows the error using a uniform square mesh with various mesh sizes and time steps in the four first lines. Note that the scheme is only first order accurate with respect to time; therefore in order to obtain second order convergence we choose  $\delta t$  proportional to  $h_{\mathcal{D}}^2$ . We also compare the error for the approximate saturation using a uniform mesh and an adaptive mesh with a similar number of unknowns. In both cases: about 300 unknowns (line 2 - line 5) and 1200 unknowns (line 3 - line 6), the adaptive mesh compared to the fixed one provides slightly better results for the saturation  $c(u)$ . The observed computational gain in relative error is rather small (about 10 – 20%), which is due to the fact that the area of high gradients of  $c$  is comparatively large.

## 7.2 The Haverkamp problem

We consider the case of a sand ground represented by the space domain  $\Omega = (0, 2) \times (0, 40)$  on the time interval  $[0, 600]$ . The parameters are given by [16]

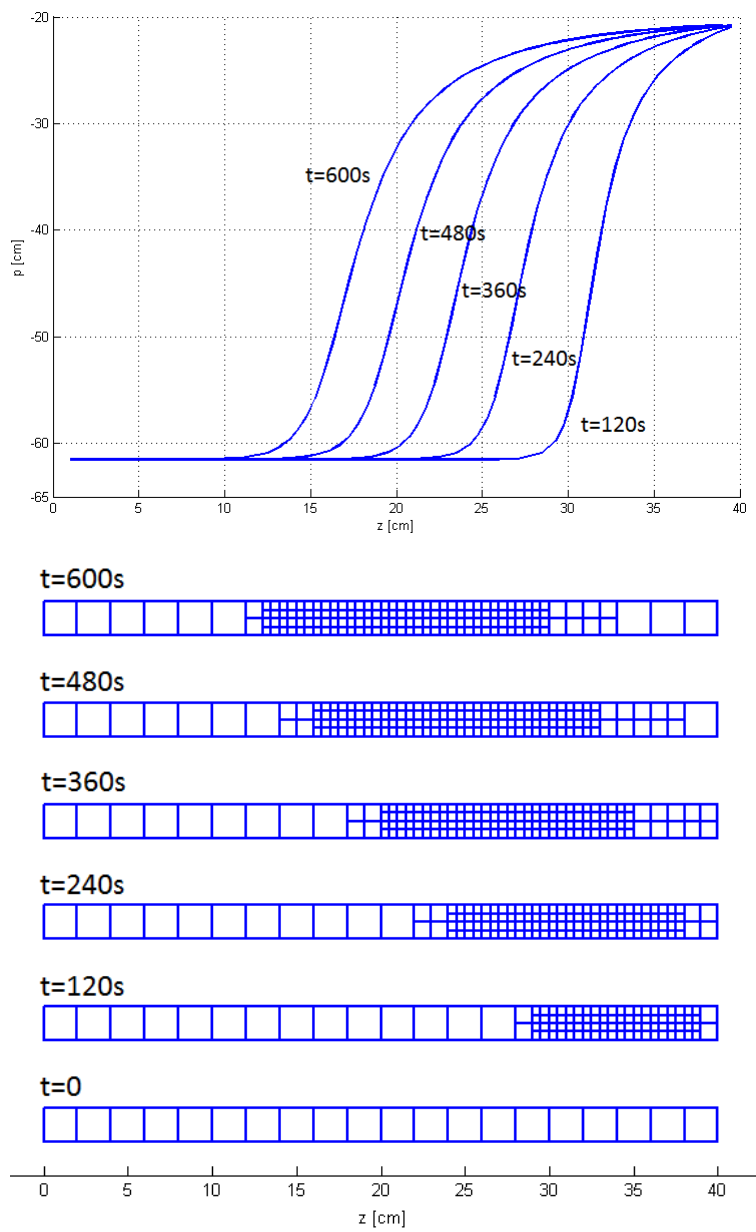
Profiles of saturation  $\theta(p)$  and permeability  $k_r(p)$  in the Haverkamp problem.The functions  $u(p)$  and  $k_r(u)$  in the Haverkamp problem.**Fig. 5** Parameters in the Haverkamp problem.

$$\theta(p) = \begin{cases} \frac{\theta_s - \theta_r}{1 + |\alpha p|^\beta} + \theta_r & \text{if } p < 0, \\ \theta_s & \text{otherwise,} \end{cases}$$

$$k_r(\theta(p)) = \begin{cases} \frac{K_s}{1 + |A p|^\gamma} & \text{if } p < 0, \\ K_s & \text{otherwise,} \end{cases}$$

where  $\theta_s = 0.287$ ,  $\theta_r = 0.075$ ,  $\alpha = 0.0271$ ,  $\beta = 3.96$ ,  $K_s = 9.44e - 3$ ,  $A = 0.0524$  and  $\gamma = 4.74$ . From  $\theta$  and  $K$ , we have tabulated suitable values for the functions  $c$  and  $K_c$ . We have taken here the initial condition  $p = -61.5$ , a homogeneous Neumann boundary condition for  $x = 0$  and  $x = 2$ , the Dirichlet boundary condition  $p = -61.5$  for  $z = 0$  and  $p = -20.7$  for  $z = 40$ .

We use an adaptive mesh and the time step  $\delta t = 1$  to perform the test. Figure 6-(a) represents the pressure profile at various times. In this test, no analytical solution is known. Therefore we compare our numerical solution with that of Pierre Sochala [22, Fig. 2.6, p. 35] which is obtained by means of a finite element method. Our results are quite similar to his. Figure 6-(b) shows the time evolution of the mesh at different times corresponding to the pressure profiles in Figure 6-(a).



**Fig. 6** Time evolution of the pressure  $p$  and the adaptive mesh in the Haverkamp problem.

**Acknowledgements** We thank Professor Pascal Omnès as well as the referees for a careful rereading of our manuscript which has led to many improvements. This work was supported by the ITN Marie Curie Network FIRST and Fondation Jacques Hadamard.

## References

1. Ophélie Angelini, Konstantin Brenner, Danielle Hilhorst. *A finite volume method on general meshes for a degenerate parabolic convection-reaction-diffusion equation*, Numerische Mathematik 123.2, 219-257, 2013.
2. Lourenço Beiro Da Veiga, Jérôme Droniou, Gianmarco Manzini. *A unified approach for handling convection terms in finite volumes and mimetic discretization methods for elliptic problems*, IMA J.Numer. Anal, 31, 4, 1357-1401, 2011.
3. Konstantin Brenner, *Méthodes de volumes finis sur maillages quelconques pour des systèmes d'évolution non linéaires*, Thèse de doctorat, Université Paris-Sud, 2011.
4. L. M. Chounet, Danielle Hilhorst, Claude Jouron, Youcef Kelanemer, P. Nicolas. *Saturated-unsaturated simulation for coupled heat and mass transfer in the ground by means of a mixed finite element method*, Adv. Water Resources, 22.5, 445-460, 1999.
5. Klaus Deimling. *Nonlinear functional analysis*, Springer-Verlag, 1985.
6. Jérôme Droniou. *Finite volume schemes for diffusion equations: introduction to and review of modern methods*, Math. Models Methods Appl. Sci. 24, 8, 1575-1619, 2014.
7. Jérôme Droniou, Robert Eymard, Thierry Gallouët, Raphaèle Herbin. *A unified approach to mimetic finite difference, hybrid finite volume and mixed finite volume methods*, Math. Models Methods Appl. Sci. 20, 2, 265-295, 2010.
8. Robert Eymard, Thierry Gallouët, Michaël Gutnic, Raphaèle Herbin, and D. Hilhorst. *Approximation by the finite volume method of an elliptic-parabolic equation arising in environmental studies*, Mathematical Models and Methods in Applied Sciences. 11.9, 1505-1528, 2001.
9. Robert Eymard, Thierry Gallouët, Raphaèle Herbin. *Discretization of heterogeneous and anisotropic diffusion problems on general nonconforming meshes SUSHI: a scheme using stabilization and hybrid interfaces*, IMA J. Numer. Anal. 30.4, 1009-1043, 2010.
10. Robert Eymard, Thierry Gallouët, Raphaèle Herbin. *Finite Volume Methods, Handbook of Numerical Analysis*, volume 7. P. G. Ciarlet and J.L.Lions eds, Elsevier Science B.V., 2000.
11. Robert Eymard, Thierry Gallouët, Raphaèle Herbin, Anthony Michel. *Convergence of a finite volume scheme for nonlinear degenerate parabolic equations*, Numer. Math., 92, 41-82, 2002.
12. Robert Eymard, Cindy Guichard, R. Herbin, R. Masson. *Gradient schemes for two-phase flow in heterogeneous porous media and Richards equation ZAMM* 94.7-8, 560, 2014.
13. Robert Eymard, Michaël Gutnic, Danielle Hilhorst. *The finite volume method for Richards equation*, Comput. Geosci. 3.3-4, 259-294, 2000.
14. Robert Eymard, Danielle Hilhorst and Martin Vohralik. *A combined finite volume scheme nonconforming/ mixed-hybrid finite element scheme for degenerate parabolic problems*. Numer. Math., 105, 73-131, 2006.
15. Peter Frolkovic, Peter Knabner, Christoph Tapp, Kathrin Thiele. *Adaptive finite volume discretization of density driven flows in porous media*, in Transport de Contaminants en Milieux Poreux (Support de Cours), INRIA, 322-355, 1997.
16. R. Haverkamp, M. Vauclin, J. Touma, P. Wierenga, and G. Vachaud, *A comparison of numerical simulation models for one-dimensional infiltration*, Soil Sci. Soc. Am. J., 41, 285-294, 1977.
17. Ulrich Hornung, *Numerische Simulation von gesättigt-ungesättigt Wasserflüssen in porösen Medien*, in *Numerische Behandlung von Differentialgleichungen mit besonderer Berücksichtigung freier Randwertaufgaben*, eds I.S.N.M., 39, Birkhäuser Verlag, 214-232, 1978.
18. Youcef Kelanemer, *Transferts couplés de masse et de chaleur dans les milieux poreux: Modélisation et étude numérique*, Thèse de doctorat, Université Paris-Sud, 1994.
19. Peter Knabner, *Finite element simulation of saturated-unsaturated flow through porous media, Large scale scientific computing*, Prog. Sci. Comput. 7, 83-93, 1987.
20. F. Otto,  *$L^1$ -contraction and uniqueness for quasilinear elliptic-parabolic equations*, J. Differential Equations 131.1, 20-38, 1996.
21. Narisoa Ramarosy, *Application de la méthode des volumes finis à des problèmes d'environnement et de traitement d'image*, Thèse de doctorat, Université Paris-Sud, 1999.
22. Pierre Sochala, *Méthodes numériques pour les écoulements souterrains et couplage avec le ruissellement*, Thèse de doctorat, École National des Ponts et Chaussées, 2008.