

NUMERICAL HOMOGENIZATION OF A NONLINEARLY COUPLED ELLIPTIC-PARABOLIC SYSTEM, REDUCED BASIS METHOD, AND APPLICATION TO NUCLEAR WASTE STORAGE

ANTOINE GLORIA

*Département de mathématiques, Université Libre de Bruxelles (ULB), 1050 Brussels, Belgium,
& Project-team SIMPAF, Inria Lille - Nord Europe, 59650 Villeneuve d'Ascq, France
gloria@ulb.ac.be*

THIERRY GOUDON

*Project-team COFFEE, Inria Sophia Antipolis Méditerranée, France
& Labo. J. A. Dieudonné, UMR 7351, CNRS-Université Nice Sophia Antipolis
thierry.goudon@inria.fr*

STELLA KRELL

*ANDRA DRD/EAP & Project-team SIMPAF, Inria Lille-Nord Europe, France
(Current address: Labo. J. A. Dieudonné, UMR 7351, CNRS-Université Nice Sophia Antipolis)
stella.krell@unice.fr*

Received (Day Month Year)

Revised (Day Month Year)

Communicated by (xxxxxxxxxx)

We consider the homogenization of a coupled system of PDEs describing flows in heterogeneous porous media. Due to the coupling, the effective coefficients always depend on the slow variable, even in the simple case when the porosity is periodic. Therefore the most important part of the computational time for the numerical simulation of such flows is dedicated to the determination of these coefficients. We propose a new numerical algorithm based on Reduced Basis techniques, which significantly improves the computational performances.

Keywords: Porous media ; Homogenization ; Reduced basis method.

AMS Subject Classification: 35B27, 65M60, 65Z05, 76M50, 76M10

1. Introduction

This work is concerned with the numerical treatment of a nonlinearly coupled elliptic-parabolic system of equations whose coefficients vary on a small scale. Resolving the finest scales induces a prohibitive numerical cost, both in terms of computational time and memory storage. Our goal consists in finding relevant “averaged” models, combined with efficient numerical methods. It turns out that the main part of the computational effort is precisely devoted to the evaluation of the

coefficients of the effective equations which are obtained by homogenization. We shall propose methods which lead to a considerable speed-up of this crucial step.

A strong motivation comes from the modeling of radionuclide transport in nuclear waste storage devices. This yields a nonlinear system of parabolic equations, coupling the time-evolution of the radionuclide concentration $C(t, x)$ (for the sake of simplicity we consider only one single species of radionuclides) to the velocity field $U(t, x)$ of the water flow. The flow takes place in a complex porous medium made of clay, limestone and marl — so that the physical properties vary a lot from place to place. The modeling of radionuclide transport in disposal facilities of radioactive waste therefore requires to deal with PDEs whose coefficients are heterogeneous at small scales. The realization of routine simulations should however rely on fast computations, which excludes to resolve the finest scales. Homogenization is the natural tool to derive effective models, which hopefully smooth out in a consistent way the small scale features of the problem. In the case of the nonlinearly coupled system treated here, (periodic) homogenization alone is not enough to drastically reduce the computational cost, since a so-called cell-problem (which is itself an elliptic PDE) has to be solved at each Gauss point of the computational domain — this could surprise the expert: although diffusion coefficients are assumed to be periodic, and the equations are linear, the nonlinear coupling condition makes the homogenized diffusion matrix depend on the space variable. This is where the reduced basis (RB) method comes into the picture: these cell-problems can be viewed as a d -parameter (d being the space dimension) family of elliptic equations, which is an ideal setting for the RB method. A further practical issue is related to the dependence of the elliptic operator upon the parameters, which is not affine (according to the terminology of the RB approach) and therefore requires a specific treatment.

The model under investigation here has been derived for the benchmark COU-PLEX, see ⁴. This benchmark is based on a set of simplified but realistic models for the transport of radionuclides around a nuclear waste repository. It allows one to evaluate the pros and cons of several numerical strategies that can be used in this context. We wish to complete the benchmark by considering the corresponding homogenization problem. Let us recall the model. We are interested in the evolution in time and space of the concentration C of a species (pollutant, nuclear waste, etc.) in a saturated porous medium Ω , which is driven by reaction, transport and diffusion:

$$R\phi\partial_t C - \nabla \cdot (D(U)\nabla C - UC) + R\phi\lambda C = S. \quad (1.1)$$

This equation involves the following quantities: the fluid velocity U , a nonlinear diffusion matrix field $D(U)$, the porosity $\phi > 0$, the species-dependent latency retardation factor $R > 0$ and degradation coefficient $\lambda = \frac{\ln(2)}{T}$ (where T is the half-time of the species), and a source term S . The system is completed by initial and

boundary conditions, which can be of Dirichlet, Neumann, or mixed type. Both the diffusion coefficients $D(U)$ and transport term $U \cdot \nabla C$ depend on the fluid velocity, which is itself related to the hydrodynamic load (or piezoelectric height) $\Theta = \frac{p}{\rho g} + z$ ($z = x_d \in \mathbb{R}$ being the height coordinate), through the formula

$$U = -K\nabla\Theta, \quad (1.2)$$

where K is the heterogeneous permeability tensor of the porous medium. In this work, we focus on the simplest situation possible where the hydraulic regime is established and governed by a mere diffusion equation relating the charge to the stationary source flow q :

$$\nabla \cdot U = q. \quad (1.3)$$

The diffusion coefficients satisfy $D(U) = D_0 + \mathbb{D}(U)$, where D_0 is a diffusion matrix field and $\mathbb{D}(U)$ is given by

$$\mathbb{D}(U) = \alpha|U|\mathbb{I} + \beta \frac{U \otimes U}{|U|}, \quad (1.4)$$

for some $\alpha, \beta \geq 0$. We shall write the COUPLEX system (1.1)–(1.4) in dimensionless form, and identify a small parameter $0 < \varepsilon \ll 1$, which is the ratio between the typical period of the heterogeneities and the characteristic length scale of Ω . The rescaled system has the same form as (1.1)–(1.4), with however a permeability matrix of the form $K(x/\varepsilon)$ for some periodic function K . We may consider oscillating source terms q, S and diffusion coefficient D_0 as well.

This work addresses the following questions:

- (1) prove existence and uniqueness of suitable weak solutions to COUPLEX system (1.1)–(1.4),
- (2) find effective equations for this model of transport of radionuclides in porous media describing the regime $0 < \varepsilon \ll 1$,
- (3) design numerical methods to compute efficiently the coefficients of the effective models, without using the brutal and prohibitive method consisting of solving independently the cell-problems and computing the suitable averages at each Gauss point of the computational domain.

The first two points are considered as a preliminary to the design of a numerical solution method. Note that the COUPLEX system can be embedded into a much more general family of systems, which has been widely studied in the literature. The main difference with the existing literature is the coupling condition: in the COUPLEX system (1.1)–(1.4), C depends on Θ but Θ does not depend on C ; this is a weak coupling. In more general models, Θ depends on C through an additional viscosity term in the Darcy equation (1.3); this is a strong coupling. In the latter case, existence results have been obtained in ¹⁶, ¹⁷, and ¹⁰. Uniqueness of weak solutions has not been proved in general. This is a specific feature of the COUPLEX

system, which we shall address in Section 2. The problem is not trivial since the coefficients $D(U)$ are unbounded unless $\nabla\Theta$ is bounded (which does not hold in general). The system with strong coupling has been homogenized by Choquet and Sili in ¹¹. We provide a much simpler proof in the case of weak coupling, which includes in addition the uniqueness property. As will be shown in Section 2, the effective model in the regime $\varepsilon \rightarrow 0$ has the form

$$\begin{cases} \nabla \cdot U^* = q^*, \\ \partial_t C^* - \nabla \cdot (D^* \nabla C^* - U^* C^*) + \lambda C^* = S^*. \end{cases} \quad (1.5)$$

where the coefficients S^* , q^* are determined by suitable “averages” of the oscillating coefficients S , q , while the velocity field U^* and the diffusion coefficient D^* are given by relations of the form

$$U^* = \mathcal{U}(K, q), \quad D^* = \mathcal{D}(D_0, U^*), \quad (1.6)$$

for some nonlinear maps \mathcal{U} and \mathcal{D} . The core of this article is the numerical approximation of this homogenized system in Section 3. We shall see that the direct approach for the computation of the effective coefficients, which consists in solving corrector equations at each Gauss point of Ω , remains prohibitive in terms of computational cost. We then propose a numerical strategy based on the RB method (see ²¹ and the references therein for elliptic equations, and ⁵ for an application to homogenization). The guideline of the RB approach is the construction of a suitable Galerkin basis “adapted” to the parametrized set of equations. We present in detail the application of the RB method to the COUPLEX system (1.1)–(1.4). Again, the fact that the diffusion coefficients are unbounded raises some interesting questions, this time not only for the analysis but also for the practical implementation of the RB method, and more specifically for the choice of the estimator. Numerical results demonstrate the ability of the method to provide accurate results with a substantial speed-up.

We shall make use of the following notation:

- $\mathbb{R}^+ = [0, +\infty)$;
- $0 < T < \infty$ is a final time;
- $d \geq 1$ denotes the space dimension;
- $\mathcal{M}_d(\mathbb{R})$ is the set of $d \times d$ real matrices, \mathbb{I} is the identity matrix;
- Ω is an open bounded Lipschitz domain of \mathbb{R}^d ;
- For all $p \in [1, \infty]$ and $s \in \mathbb{N}$, $L^p(\Omega)$ denotes the space of p -integrable functions on Ω , $W^{s,p}(\Omega)$ denotes the Sobolev space of p -integrable functions whose s -first distributional derivatives are p -integrable, $W_0^{1,p}(\Omega)$ the closure in $W^{1,p}(\Omega)$ of the space $C_0^\infty(\Omega)$ of smooth functions compactly supported in Ω ;
- For $p = 2$, we denote the Hilbert spaces $W^{1,2}(\Omega)$ and $W_0^{1,2}(\Omega)$ by $H^1(\Omega)$ and $H_0^1(\Omega)$, respectively.
- $\mathbb{Y} = (0, 1)^d$ is the periodic cell, and $H_{\#}^1(\mathbb{Y})$ denotes the closure of the subspace of $C^\infty(\mathbb{R}^d)$ made of \mathbb{Y} -periodic functions with vanishing mean.

2. Well-posedness and homogenization

2.1. Main results

We consider the following weakly coupled system of PDEs:

$$\begin{cases} U = -K\nabla\Theta & \text{in } \Omega, \\ \nabla \cdot U = q & \text{in } \Omega, \\ \partial_t C - \nabla \cdot (D(U)\nabla C - UC) + \lambda C = S & \text{in }]0, T[\times \Omega. \end{cases} \quad (2.1)$$

Let $\lambda > 0$, $q \in L^\infty(\Omega)$, and $S \in L^2(0, T; H^{-1}(\Omega))$. The weak coupling condition reads

$$D(U)(x) := D_0(x) + \alpha|U(x)|\mathbb{I} + \beta \frac{U(x) \otimes U(x)}{|U(x)|}, \quad (2.2)$$

for a. e. $x \in \Omega$, where $\alpha > 0, \beta \geq 0$. The functions $x \mapsto K(x)$ and $x \mapsto D_0(x)$ are matrix-valued; they both satisfy uniform bounds and strong ellipticity conditions: there exists $\Lambda > 0$ such that for a. e. $x \in \Omega$ and all $\xi \in \mathbb{R}^d$

$$\begin{aligned} |K(x)\xi| &\leq \Lambda|\xi|, \quad \xi \cdot K(x)\xi \geq \Lambda^{-1}|\xi|^2, \\ |D_0(x)\xi| &\leq \Lambda|\xi|, \quad \xi \cdot D_0(x)\xi \geq \Lambda^{-1}|\xi|^2. \end{aligned}$$

The system (2.1) is completed by boundary conditions and an initial condition. For the mathematical analysis of the problem, we restrict ourselves to homogeneous Dirichlet boundary conditions; namely, we set

$$\begin{cases} \Theta = 0 & \text{on } \partial\Omega, \\ C(0, \cdot) = C_{\text{init}} & \text{in } \Omega, \\ C = 0 & \text{on }]0, T[\times \partial\Omega, \end{cases} \quad (2.3)$$

for some $C_{\text{init}} \in L^2(\Omega)$. The adaptation to more general (time-independent) boundary conditions, as treated in the numerical tests later on, could be considered without further difficulty.

We are interested in the case when K is an ε -periodic matrix, and $\varepsilon \rightarrow 0$. Before we turn to this problem, we first define a notion of weak solution for the coupled system (2.1)–(2.3), and give an existence and uniqueness result.

Definition 1. A weak solution of (2.1)–(2.3) is a pair $(\Theta, C) \in H_0^1(\Omega) \times L^2(0, T; H_0^1(\Omega)) \cap C^0(0, T; L^2(\Omega))$ such that $\partial_t C \in L^2(0, T; H^{-1}(\Omega))$, $\int_0^T \int_\Omega \nabla C \cdot D(U)\nabla C < \infty$ with $U = -K\nabla\Theta$, and which satisfies (2.1)–(2.3) in the following sense:

- Θ is a weak solution in $H_0^1(\Omega)$ to (2.1)_{1,2} & (2.3)₁;
- For all $v \in L^2(0, T; H_0^1(\Omega)) \cap L^2(0, T; L^\infty(\Omega))$ such that $\int_0^T \int_\Omega \nabla v \cdot D(U)\nabla v < \infty$, we have

$$\begin{aligned} \int_0^T \langle \partial_t C, v \rangle_{H^{-1}, H_0^1} + \int_0^T \int_\Omega \nabla v \cdot D(U)\nabla C + \int_0^T \int_\Omega vU \cdot \nabla C \\ + \int_0^T \int_\Omega Cv(q + \lambda) = \int_0^T \langle S, v \rangle_{H^{-1}, H_0^1}. \end{aligned}$$

The following theorem states the existence and uniqueness of such weak solutions.

Theorem 1. *For all $q \in L^\infty(\Omega)$, $S \in L^2(0, T; H^{-1}(\Omega))$, and $C_{\text{init}} \in L^2(\Omega)$, there exists a unique weak solution to (2.1)–(2.3) in the sense of Definition 1.*

We now turn to the periodic homogenization of (2.1)–(2.3). Let K be a $\mathbb{Y} = (0, 1)^d$ -periodic matrix. For all $\varepsilon > 0$, we consider the coupled system

$$\begin{cases} U_\varepsilon = -K_\varepsilon \nabla \Theta_\varepsilon & \text{in } \Omega, \\ \nabla \cdot U_\varepsilon = q & \text{in } \Omega, \\ \partial_t C_\varepsilon - \nabla \cdot (D(U_\varepsilon) \nabla C_\varepsilon - U_\varepsilon C_\varepsilon) + \lambda C_\varepsilon = S & \text{in }]0, T[\times \Omega, \\ \Theta_\varepsilon = 0 & \text{on } \partial\Omega, \\ C_\varepsilon(0, \cdot) = C_{\text{init}} & \text{in } \Omega, \\ C_\varepsilon = 0 & \text{on }]0, T[\times \partial\Omega, \end{cases} \quad (2.4)$$

where q, S, C_{init} and the function D are as above, and K_ε is defined by $K_\varepsilon(x) := K(x/\varepsilon)$ on Ω . Theorem 1 ensures the existence and uniqueness of a weak solution $(\Theta_\varepsilon, C_\varepsilon)$ of (2.4) for all $\varepsilon > 0$. In order to characterize the asymptotic behavior of $(\Theta_\varepsilon, C_\varepsilon)$ as $\varepsilon \rightarrow 0$ we need to introduce a few auxiliary quantities. For all $i \in \{1, \dots, d\}$, we let φ_i denote the unique periodic weak solution in $H_{\#}^1(\mathbb{Y})$ to the following elliptic equation

$$-\nabla \cdot K(\mathbf{e}_i + \nabla \varphi_i) = 0. \quad (2.5)$$

We define the matrix K^* by: for all $i, j \in \{1, \dots, d\}$,

$$\mathbf{e}_j \cdot K^* \mathbf{e}_i = \int_{\mathbb{Y}} (\mathbf{e}_j + \nabla \varphi_j) \cdot K(\mathbf{e}_i + \nabla \varphi_i). \quad (2.6)$$

The matrix K^* defined this way is strongly elliptic. This allows one to define the unique weak solution $\Theta_0 \in H_0^1(\Omega)$ to the elliptic equation

$$-\nabla \cdot K^* \nabla \Theta_0 = q. \quad (2.7)$$

The homogenized drift is then given by

$$U_0 = -K^* \nabla \Theta_0. \quad (2.8)$$

Next we have to consider two-scale functions \tilde{U}, \tilde{D} defined on $\Omega \times \mathbb{Y}$ by

$$\begin{aligned} \tilde{U}(x, y) &= -K(y)(\mathbb{I} + \nabla \varphi(y)) \nabla \Theta_0(x), \\ \tilde{D}(x, y) &= D_0(x) + \alpha |\tilde{U}(x, y)| \mathbb{I} + \beta \frac{\tilde{U}(x, y) \otimes \tilde{U}(x, y)}{|\tilde{U}(x, y)|} = D_0(x) + \mathbb{D}(\tilde{U}(x, y)) \end{aligned} \quad (2.9)$$

where $\varphi = (\varphi_1, \dots, \varphi_d)$. For all $i \in \{1, \dots, d\}$ and a. e. $x \in \Omega$, we let $\Phi_i(x, \cdot)$ denote the unique periodic weak solution in $H_{\#}^1(\mathbb{Y})$ to the elliptic equation parametrized by x :

$$-\nabla_y \cdot \tilde{D}(x, y)(\mathbf{e}_i + \nabla_y \Phi_i(x, y)) = 0.$$

We finally define a homogenized matrix field D^* by: for all $x \in \Omega$ and all $i, j \in \{1, \dots, d\}$,

$$\mathbf{e}_j \cdot D^*(x) \mathbf{e}_i = \int_{\mathbb{Y}} (\mathbf{e}_j + \nabla_y \Phi_j(x, y)) \cdot \tilde{D}(x, y) (\mathbf{e}_i + \nabla_y \Phi_i(x, y)) dy. \quad (2.11)$$

We point out that, while K^* is a constant matrix, D^* is not a constant matrix field, since $\nabla \Theta_0$ is in general not constant on Ω . We are now in position to describe the asymptotic behavior.

Theorem 2. *Let $q \in L^\infty(\Omega)$, $S \in L^2(0, T; H^{-1}(\Omega))$, $C_{\text{init}} \in L^2(\Omega)$, D be as in (2.2), and let K be a \mathbb{Y} -periodic bounded and strongly elliptic matrix. For all $\varepsilon > 0$, we set $K_\varepsilon := K(\cdot/\varepsilon)$. We let K^* , Θ_0 , U_0 , and D^* be as in (2.6), (2.7), (2.8), and (2.11), respectively. Then there exists a unique weak solution C_0 to*

$$\begin{cases} \partial_t C_0 - \nabla \cdot (D^* \nabla C_0 - U_0 C_0) + \lambda C_0 = S & \text{in }]0, T[\times \Omega, \\ C_0(0, \cdot) = C_{\text{init}} & \text{in } \Omega, \\ C_0 = 0 & \text{on }]0, T[\times \partial\Omega, \end{cases} \quad (2.12)$$

in the sense of Definition 1 (with D^* in place of $D(U)$), and the unique weak solution $(\Theta_\varepsilon, C_\varepsilon)$ to (2.4) converges to (Θ_0, C_0) strongly in $L^2(\Omega)$ and $L^2((0, T) \times \Omega)$, and weakly in $H^1(\Omega)$ and $L^2(0, T; H^1(\Omega))$; in addition C_ε converges in $C^0([0, T], L^2(\Omega) - \text{weak})$ to C_0 .

Although the diffusion D^* is not of the form (2.2), for all $x \in \Omega$, $D^*(x)$ only depends on $\nabla \Theta_0(x)$, and $D^* \in L^2(\Omega)$. Hence existence and uniqueness of weak solutions for the homogenized system can be proved the same way as for Theorem 1, and we leave the details to the reader. From the homogenization point of view, Theorem 2 is a rather direct application of two-scale convergence and Theorem 1. Although $D(U)$ is unbounded, it is square-integrable and the homogenized system remains elliptic-parabolic (for the homogenization of elliptic equations with unbounded coefficients which are not equi-integrable, nonlocal effects may appear, and we refer the reader to ² and ⁶). In the case of strong coupling (that is when the equation for U depends on C through a nonscantant viscosity term), homogenization has been proved in ¹¹. Yet ¹¹ is an overkill for the problem under consideration (uniqueness is not discussed in ¹¹ though), and we display the main arguments of a simpler proof in Appendix Appendix A.

Remark 1. In this statement we have considered only the case of a periodic oscillating matrix K . Note that, even in this simple case, the matrix \tilde{D} depends on both the slow variable $x \in \Omega$ and the fast variable $y \in \mathbb{Y}$. The result generalizes readily to the case of a locally periodic fields of the form $K(x, x/\varepsilon)$. Similarly, oscillating source terms and diffusion coefficients D_0 , depending on x and x/ε , can be considered.

Since the specific feature of the coupled model under investigation is the uniqueness of weak solutions, we prove Theorem 1 in detail in the following subsection.

2.2. Proof of Theorem 1

The difference with respect to previous contributions on the strongly coupled system is the fact that weak solutions can be proved to be unique for the weakly coupled system. The only subtle feature of the system is the integrability condition on D and U , which are square-integrable but not necessarily essentially bounded. The proof is based on standard regularization and compactness arguments. We shall only prove the uniqueness of weak solutions in detail.

For the reader convenience we quickly sketch the proof of existence of weak solutions as well. We divide the proof into six steps, and proceed by regularization. In the first step we recall a classical result essentially due to J.-L. Lions. In the second step we introduce the regularizations for U and $D(U)$. In the third step we apply Step 1 and derive a priori estimates. In Step 4 we deduce from these a priori estimates, by compactness and Aubin-Simon's arguments, that the weakly coupled system admits a distributional solution. We then show in Step 5 that this solution satisfies a weak formulation of the equation. We prove uniqueness of weak solutions in Step 6.

Step 1. Case of bounded coefficients and drifts.

Let $D : \Omega \rightarrow \mathcal{M}_d(\mathbb{R})$ be uniformly bounded and strongly elliptic, and $U \in L^\infty(\Omega, \mathbb{R}^d)$. We consider the equation

$$\begin{cases} \partial_t C - \nabla \cdot (D \nabla C - UC) + \lambda C = S & \text{in }]0, T[\times \Omega, \\ C(0, \cdot) = C_{\text{init}} & \text{in } \Omega, \\ C = 0 & \text{on }]0, T[\times \partial\Omega. \end{cases} \quad (2.13)$$

Then for all $S \in L^2(0, T; H^{-1}(\Omega))$ and $C_{\text{init}} \in L^2(\Omega)$ there exists a unique function $C \in L^2(0, T; H_0^1(\Omega)) \cap C^0(0, T; L^2(\Omega))$ such that $\partial_t C \in L^2(0, T; H^{-1}(\Omega))$, which satisfies the weak form of (2.13): for all $v \in L^2(0, T; H_0^1(\Omega))$ and all $T > 0$,

$$\begin{aligned} \int_0^T \langle \partial_t C, v \rangle_{H^{-1}, H_0^1} + \int_0^T \int_\Omega \nabla v \cdot D \nabla C - \int_0^T \int_\Omega C U \cdot \nabla v \\ + \lambda \int_0^T \int_\Omega C v = \int_0^T \langle S, v \rangle_{H^{-1}, H_0^1}. \end{aligned}$$

Note that by an exponential change of time, one may assume λ to be as large as desired (which then ensures the coercivity of the bilinear form). We refer the reader to ²³ for details. In the sequel we shall use the following equivalent weak formulation: for all $v \in L^2(0, T; H_0^1(\Omega))$ and all $t > 0$,

$$\begin{aligned} \int_0^t \langle \partial_t C, v \rangle_{H^{-1}, H_0^1} + \int_0^t \int_\Omega \nabla v \cdot D \nabla C + \int_0^t \int_\Omega \nabla C \cdot U v \\ + \int_0^t \int_\Omega (q + \lambda) C v = \int_0^t \langle S, v \rangle_{H^{-1}, H_0^1}, \end{aligned} \quad (2.14)$$

which we obtain by using the divergence theorem and the identity $\nabla \cdot U = q$.

Step 2. Regularizations.

First note that the elliptic part of the system

$$\begin{cases} U = -K\nabla\Theta & \text{in } \Omega, \\ \nabla \cdot U = q & \text{in } \Omega, \\ \Theta = 0 & \text{on } \partial\Omega, \end{cases} \quad (2.15)$$

admits a unique weak solution $\Theta \in H_0^1(\Omega)$. The associated drift $U = -K\nabla\Theta$ is not essentially bounded, but square-integrable. Likewise the associated diffusion coefficients $D(U)$ are not essentially bounded, but square-integrable. In particular, the advection-diffusion equation

$$\begin{cases} \partial_t C - \nabla \cdot (D(U)\nabla C - UC) + \lambda C = S & \text{in }]0, T[\times \Omega, \\ C(0, \cdot) = C_{\text{init}} & \text{in } \Omega, \\ C = 0 & \text{on }]0, T[\times \partial\Omega, \end{cases} \quad (2.16)$$

does not satisfy the assumptions of Step 1. We regularize $D(U)$ and U , and begin with the diffusion coefficients. Since $D(U)$ is a symmetric matrix, for a. e. $x \in \Omega$ there exist $\alpha_1(x), \dots, \alpha_d(x) \geq 0$, and a unitary matrix $P(x)$ such that $D(U)(x) = P^T(x)\text{diag}(\alpha_1(x), \dots, \alpha_d(x))P(x)$. For all $M > 0$ and a. e. $x \in \Omega$, we define $D^M(x)$ as follows:

$$D^M(x) := P^T(x)\text{diag}(\min\{\alpha_1(x), M\}, \dots, \min\{\alpha_d(x), M\})P(x).$$

In particular, D^M converges monotonically to $D(U)$ in $L^2(\Omega)$ as $M \rightarrow \infty$. For the regularization of U , we prefer to regularize the defining equation $\nabla \cdot U = q$ rather than using truncations. We consider K^M and q^M two sequences of smooth functions such that K^M and q^M converge to K and q in $L^r(\Omega)$ for all $r < \infty$, respectively. We define Θ^M as the unique weak solution in $H_0^1(\Omega)$ to

$$-\nabla \cdot (K^M \nabla \Theta^M) = q^M,$$

and set $U^M := -K^M \nabla \Theta^M$. By elliptic regularity, U^M belongs to $L^\infty(\Omega)$. Furthermore, U^M converges to U in $L^2(\Omega)$ (the argument relies on Meyers' estimate, which implies that $\nabla \Theta^M$ converges to $\nabla \Theta$ in $L^p(\Omega)$ for some $p > 2$ depending only on the constant Λ , whereas K^M converges in $L^{p'}(\Omega)$ to K , with $1/p + 1/p' = 1$).

Hence, for all $M > 0$, Step 1 implies there exists a unique weak solution $C^M \in L^2(0, T; H_0^1(\Omega)) \cap C^0(0, T; L^2(\Omega))$ such that $\partial_t C^M \in L^2(0, T; H^{-1}(\Omega))$ to the regularized equation

$$\begin{cases} \partial_t C^M - \nabla \cdot (D^M \nabla C^M - U^M C^M) + \lambda C^M = S & \text{in }]0, T[\times \Omega, \\ C^M(0, \cdot) = C_{\text{init}} & \text{in } \Omega, \\ C^M = 0 & \text{on }]0, T[\times \partial\Omega. \end{cases} \quad (2.17)$$

It remains to pass to the limit as $M \rightarrow \infty$.

Step 3. A priori estimates.

The weak form of (2.17) with test-function C^M itself yields for all $0 < t \leq T$,

$$\int_0^t \langle \partial_t C^M, C^M \rangle_{H^{-1}, H_0^1} + \int_0^t \int_{\Omega} \nabla C^M \cdot D^M \nabla C^M - \int_0^t \int_{\Omega} C^M U^M \cdot \nabla C^M + \lambda \int_0^t \int_{\Omega} (C^M)^2 = \int_0^t \langle S, C^M \rangle_{H^{-1}, H_0^1}.$$

Since $\nabla \cdot U^M = q^M$, we have by the divergence theorem

$$- \int_0^t \int_{\Omega} C^M U^M \cdot \nabla C^M = -\frac{1}{2} \int_0^t \int_{\Omega} U^M \cdot \nabla (C^M)^2 = \frac{1}{2} \int_0^t \int_{\Omega} q (C^M)^2$$

so that the weak form turns into

$$\begin{aligned} \frac{1}{2} \int_{\Omega} (C^M(\cdot, t))^2 + \int_0^t \int_{\Omega} \nabla C^M \cdot D^M \nabla C^M \\ + \int_0^t \int_{\Omega} (C^M)^2 \left(\frac{1}{2} q^M + \lambda \right) = \frac{1}{2} \int_{\Omega} C_{\text{init}}^2 + \int_0^t \langle S, C^M \rangle_{H^{-1}, H_0^1}. \end{aligned}$$

Recalling that one may take λ such that $\frac{1}{2} \inf q^M + \lambda \geq \frac{1}{2} \inf q + \lambda = \lambda^* > 0$, we finally deduce by coercivity of D^M (and arbitrariness of t):

$$\begin{aligned} \frac{1}{2} \sup_{0 < t \leq T} \int_{\Omega} (C^M(\cdot, t))^2 + \Lambda^{-1} \|\nabla C^M\|_{L^2(0, T; L^2(\Omega))}^2 + \lambda^* \|C^M\|_{L^2(0, T; L^2(\Omega))}^2 \\ \leq \|S\|_{L^2(0, T; H^{-1}(\Omega))} \|C^M\|_{L^2(0, T; H_0^1(\Omega))} + \frac{1}{2} \|C_{\text{init}}\|_{L^2(\Omega)}^2. \end{aligned}$$

Using this estimate and the equation again, we finally obtain that C^M is bounded in $L^2(0, T; H_0^1(\Omega)) \cap L^\infty(0, T; L^2(\Omega))$ and that $\partial_t C^M$ is bounded in $L^2(0, T; H^{-1}(\Omega))$, uniformly in M .

Step 4. Compactness and existence of distributional solutions.

By weak compactness and Aubin-Simon's theorem (see ²⁵), there exists a function $C \in L^2(0, T; H_0^1(\Omega)) \cap C^0(0, T; L^2(\Omega))$ with $\partial_t C \in L^2(0, T; H^{-1}(\Omega))$ such that C^M converges weakly to C in $L^2(0, T; H_0^1(\Omega))$, strongly in $L^2(0, T; L^2(\Omega))$, and such that $\partial_t C^M$ converges weakly to $\partial_t C$ in $L^2(0, T; H^{-1}(\Omega))$.

It is easy matter to check that C solves (2.16) in the sense of distributions, and satisfies the initial condition as a continuous function in time taking values in $L^2(\Omega)$.

Step 5. Weak formulation of the system.

In this step, we shall prove that (Θ, C) is a weak solution of (2.1)–(2.3) in the sense of Definition 1. We start by showing that $\int_0^T \int_{\Omega} \nabla C \cdot D(U) \nabla C < \infty$, which is not obvious a priori since $\nabla C \in L^2(\Omega)$ and $D(U) \in L^2(\Omega)$. Let $M' > 0$ be fixed. By weak lower-semicontinuity, since ∇C^M converges weakly to ∇C in $L^2(0, T; L^2(\Omega))$,

$$\int_0^T \int_{\Omega} \nabla C \cdot D^{M'} \nabla C \leq \liminf_{M \rightarrow \infty} \int_0^T \int_{\Omega} \nabla C^M \cdot D^{M'} \nabla C^M.$$

Since $M \mapsto D^M$ is an increasing function in the sense of symmetric matrices, the a priori estimate of Step 3 implies for all $M \geq M'$,

$$\int_0^T \int_{\Omega} \nabla C^M \cdot D^{M'} \nabla C^M \leq \int_0^T \int_{\Omega} \nabla C^M \cdot D^M \nabla C^M \leq \|C_{\text{init}}\|_{L^2(\Omega)}^2 + \|S\|_{L^2(0,T;H^{-1}(\Omega))}^2.$$

Hence,

$$\int_0^T \int_{\Omega} \nabla C \cdot D^{M'} \nabla C \leq \|C_{\text{init}}\|_{L^2(\Omega)}^2 + \|S\|_{L^2(0,T;H^{-1}(\Omega))}^2,$$

and the desired estimate

$$\int_0^T \int_{\Omega} \nabla C \cdot D(U) \nabla C \leq \|C_{\text{init}}\|_{L^2(\Omega)}^2 + \|S\|_{L^2(0,T;H^{-1}(\Omega))}^2$$

follows from the monotone convergence theorem as $M' \rightarrow \infty$, using again the monotonicity of $M \mapsto D^M$.

Let $v \in L^2(0, T; H_0^1(\Omega)) \cap L^2(0, T; L^\infty(\Omega))$ be such that $\int_0^T \int_{\Omega} \nabla v \cdot D(U) \nabla v < \infty$. In order to prove that C is a weak solution, we need to pass to the limit as $M \rightarrow \infty$ in the weak formulation (2.14) for C^M , that is,

$$\begin{aligned} \int_0^t \langle \partial_t C^M, v \rangle_{H^{-1}, H_0^1} + \int_0^t \int_{\Omega} \nabla v \cdot D^M \nabla C^M + \int_0^t \int_{\Omega} \nabla C^M \cdot U^M v \\ + \int_0^t \int_{\Omega} (q + \lambda) C^M v = \int_0^t \langle S, v \rangle_{H^{-1}, H_0^1}. \end{aligned}$$

In view of the regularity of v , the convergence of U^M to U in $L^2(\Omega)$, and the weak compactness of C_M obtained in Step 4, the only nontrivial term to treat is the second-order term, and we shall prove that

$$\lim_{M \rightarrow \infty} \int_0^T \int_{\Omega} \nabla v \cdot D^M \nabla C^M = \int_0^T \int_{\Omega} \nabla v \cdot D(U) \nabla C.$$

Let $M' > 0$ be fixed. We rewrite the above term as

$$\int_0^T \int_{\Omega} \nabla v \cdot D^M \nabla C^M = \int_0^T \int_{\Omega} \nabla v \cdot (D^M - D^{M'}) \nabla C^M + \int_0^T \int_{\Omega} \nabla v \cdot D^{M'} \nabla C^M. \quad (2.18)$$

We focus on the second term of the r. h. s. and first take the limit as $M \rightarrow \infty$. Since $D^{M'}$ is bounded, this yields

$$\lim_{M \rightarrow \infty} \int_0^T \int_{\Omega} \nabla v \cdot D^{M'} \nabla C^M = \int_0^T \int_{\Omega} \nabla v \cdot D^{M'} \nabla C.$$

We conclude by the dominated convergence theorem that

$$\lim_{M' \rightarrow \infty} \lim_{M \rightarrow \infty} \int_0^T \int_{\Omega} \nabla v \cdot D^{M'} \nabla C^M = \int_0^T \int_{\Omega} \nabla v \cdot D(U) \nabla C$$

since by Young's inequality and monotonicity of $M \mapsto D^M$,

$$|\nabla v \cdot D^{M'} \nabla C| \leq \frac{1}{2} (\nabla v \cdot D(U) \nabla v + \nabla C \cdot D(U) \nabla C),$$

which is integrable.

It remains to prove that the first term of the r. h. s. of (2.18) vanishes as M' and M go to infinity. By Cauchy-Schwarz inequality and monotonicity of $M \mapsto D^M$, we have

$$\begin{aligned} & \left| \int_0^T \int_{\Omega} \nabla v \cdot (D^M - D^{M'}) \nabla C^M \right| \\ & \leq \left(\int_0^T \int_{\Omega} \nabla v \cdot (D - D^{M'}) \nabla v \right)^{1/2} \left(\int_0^T \int_{\Omega} \nabla C^M \cdot D^M \nabla C^M \right)^{1/2}. \end{aligned}$$

The second factor of the r. h. s. is bounded by Step 3 uniformly in M . We therefore focus on the first factor. Since $D - D^{M'} \leq D$ in the sense of symmetric matrices and $\nabla v \cdot D \nabla v \in L^1(\Omega)$, the dominated convergence theorem yields

$$\lim_{M' \rightarrow \infty} \int_0^T \int_{\Omega} \nabla v \cdot (D - D^{M'}) \nabla v = 0,$$

and therefore

$$\lim_{M' \rightarrow \infty} \limsup_{M \rightarrow \infty} \left| \int_0^T \int_{\Omega} \nabla v \cdot (D^M - D^{M'}) \nabla C^M \right| = 0,$$

which concludes the proof of this step.

Step 6. Uniqueness of weak solutions.

Since equation (2.16) is linear with respect to C , uniqueness follows formally from the weak formulation tested with the solution C itself. However, we cannot directly proceed this way since $C \notin L^2(0, T; L^\infty(\Omega))$ a priori and it is not clear whether it can be used as an admissible test function. Instead we use a standard truncation argument: for all $N > 0$ we define a function $\varphi_N : \mathbb{R} \rightarrow \mathbb{R}$ as

$$\varphi_N(x) := \begin{cases} -N & \text{for } x < -N, \\ x & \text{for } |x| \leq N, \\ N & \text{for } x > N, \end{cases}$$

and we test the weak formulation of (2.16) with $C_N := \varphi_N(C) \in L^2(0, T; H_0^1(\Omega)) \cap L^\infty((0, T) \times \Omega)$. This yields

$$\begin{aligned} & \int_0^T \langle \partial_t C, C_N \rangle_{H^{-1}, H_0^1} + \int_0^T \int_{\Omega} \nabla C_N \cdot D(U) \nabla C + \int_0^T \int_{\Omega} C_N U \cdot \nabla C \\ & \quad + \int_0^T \int_{\Omega} C C_N (q + \lambda) = \int_0^T \langle S, C_N \rangle_{H^{-1}, H_0^1}. \end{aligned} \quad (2.19)$$

It is easy to prove that $C_N \rightarrow C$ in $L^2(0, T; H^1(\Omega))$ as $N \rightarrow \infty$ so that we can pass to the limit in the first and last terms of the l. h. s. and in the r. h. s. of (2.19). It remains to treat the last two terms. We begin with the Dirichlet form. By definition of φ_N and C_N ,

$$\nabla C_N \cdot D(U) \nabla C = \nabla C \cdot D(U) \nabla C 1_{|C| \leq N} \leq \nabla C \cdot D(U) \nabla C.$$

Hence, the dominated convergence theorem yields

$$\lim_{N \rightarrow \infty} \int_0^T \int_{\Omega} \nabla C_N \cdot D(U) \nabla C = \int_0^T \int_{\Omega} \nabla C \cdot D(U) \nabla C.$$

We now turn to the third term of the l. h. s. of (2.19), which we treat together with the term involving q . In particular since $\nabla \cdot U = q$, the divergence theorem yields

$$\int_0^T \int_{\Omega} C_N U \cdot \nabla C + \int_0^T \int_{\Omega} C C_N q = - \int_0^T \int_{\Omega} C U \cdot \nabla C_N.$$

Note that by definition of φ_N and C_N we can rewrite this identity as

$$\int_0^T \int_{\Omega} C_N U \cdot \nabla C + \int_0^T \int_{\Omega} C C_N q = - \int_0^T \int_{\Omega} C_N U \cdot \nabla C_N.$$

Using that $\nabla \cdot U = q$ and the divergence theorem again, this turns into

$$\begin{aligned} \int_0^T \int_{\Omega} C_N U \cdot \nabla C + \int_0^T \int_{\Omega} C C_N q &= - \int_0^T \int_{\Omega} C_N U \cdot \nabla C_N \\ &= - \frac{1}{2} \int_0^T \int_{\Omega} U \cdot \nabla C_N^2 \\ &= \frac{1}{2} \int_0^T \int_{\Omega} q C_N^2. \end{aligned}$$

Passing to the limit in the last identity yields

$$\lim_{N \rightarrow \infty} \left(\int_0^T \int_{\Omega} C_N U \cdot \nabla C + \int_0^T \int_{\Omega} C C_N q \right) = \frac{1}{2} \int_0^T \int_{\Omega} q C^2.$$

Gathering the results of this step, we obtain the following identity:

$$\begin{aligned} \frac{1}{2} \int_{\Omega} C^2(T, \cdot) + \int_0^T \int_{\Omega} \nabla C \cdot D(U) \nabla C \\ + \int_0^T \int_{\Omega} C^2 \left(\frac{1}{2} q + \lambda \right) &= \int_0^T \langle S, C \rangle_{H^{-1}, H_0^1} + \frac{1}{2} \int_{\Omega} C_{\text{init}}^2, \end{aligned}$$

since $\int_0^T \langle \partial_t C, C \rangle_{H^{-1}, H_0^1} = \frac{1}{2} \int_{\Omega} C^2(T, \cdot) - \frac{1}{2} \int_{\Omega} C^2(0, \cdot)$. This implies uniqueness of weak solutions, and concludes the proof of Theorem 1.

3. Numerical approximation of the homogenized system

In this section we propose a numerical strategy to approximate the weak solution to the homogenized system (2.7)–(2.12). There are essentially three steps to solve (2.7)–(2.12):

- (1) the computation of K^* and the approximation of Θ_0 . The latter is solution of a standard elliptic equation once K^* is known, see (2.7).
- (2) the approximation of $D^*(x)$ at every Gauss point x of Ω . This requires to solve a family of elliptic equations on the periodic cell \mathbb{Y} , parametrized by the Gauss points x via $\nabla \Theta_0(x)$.

(3) the numerical solution of the advection-diffusion equation (2.12).

As we shall see, the bottleneck of the numerical approximation of (2.7)–(2.12) in terms of computational cost is the approximation of D^* in the second step. A large part of this section is dedicated to this problem, and we shall use a reduced basis approach to drastically reduce this computational cost. We have chosen not to focus on the numerical strategy to solve the advection-diffusion equation (2.12) since the equation is rather “standard” once D^* is known. In particular, for the numerical tests of the coupled system we use a naive \mathbb{P}_1 -finite element method in space combined with the implicit Euler method in time. For more efficient and modern methods, we refer the reader to ^{24,14,15,26,27,8,9}. The main contribution of this section (a numerical method for the computation of D^*) can indeed be combined with any strategy to solve the advection-diffusion equation (2.12).

In the first subsection we present a direct approach to solve (2.7)–(2.12), and complement the homogenization result of Theorem 2 by numerical tests showing the rate of convergence of $(\Theta_\varepsilon, C_\varepsilon)$ towards (Θ_0, C_0) . As expected, the computational time to approximate D^* becomes rapidly prohibitive as the number of discretization points increases. In the second subsection we turn to the RB method. We first quickly recall the rationale of the approach, and discuss what can be expected in terms of convergence. We then turn to the practical implementation of the method, propose an a posteriori estimator adapted to homogenization problems (but not limited to the specific one treated here), and present an original and effective way of fast-assembling of the RB matrix, which is the major difficulty encountered in the RB method when the dependence of the diffusion matrix upon the parameter is not affine — as it is the case here.

Before we turn to the core of this section let us point out that, as the attentive reader may have already noticed, it is not clear a priori that the finite element method converges since the diffusion matrix in (2.12) is unbounded. The method does indeed converge to the expected solution. This property can be proved along the lines of the existence-uniqueness theory developed in Section 2.

3.1. Direct approach

3.1.1. Space and time discretizations

We discretize the homogenized problem (2.7)–(2.12) with a finite element method in space and the implicit Euler scheme in time. Let T_{Ω, h_0} and $T_{\mathbb{Y}, h_1}, T_{\mathbb{Y}, \bar{h}_1}$ be regular tessellations of Ω and of \mathbb{Y} , respectively, into d -simplices of meshsizes $h_0, h_1, \bar{h}_1 > 0$. We denote by $\mathcal{V}_{\Omega, h_0}^k$ the space of \mathbb{P}_k finite elements associated with T_{Ω, h_0} for $k = 0$ and 1 (for $k = 1$, we only consider functions which vanish on the boundary), and by $\mathcal{V}_{\mathbb{Y}, h_1}^1, \mathcal{V}_{\mathbb{Y}, \bar{h}_1}^1$ (resp. $\mathcal{V}_{\mathbb{Y}, h_1}^0$) the subspaces of $H_{\#}^1(\mathbb{Y})$ (resp. $L^2(\mathbb{Y})$) made of \mathbb{P}_1 -periodic (resp. \mathbb{P}_0) finite elements associated with $T_{\mathbb{Y}, h_1}, T_{\mathbb{Y}, \bar{h}_1}$. As quickly mentioned above, a natural strategy to solve (2.7)–(2.12) is as follows:

Algorithm 1.

- (1) Numerical approximation $K_{\bar{h}_1}^*$ of K^* : compute for all $k \in \{1, \dots, d\}$ Galerkin approximations $\varphi_k^{\bar{h}_1}$ of φ_k in $\mathcal{V}_{\mathbb{Y}, \bar{h}_1}^1$ defined by: For all $\psi \in \mathcal{V}_{\mathbb{Y}, \bar{h}_1}^1$

$$\int_{\mathbb{Y}} \nabla \psi \cdot K(\mathbf{e}_k + \nabla \varphi_k^{\bar{h}_1}) = 0. \quad (3.1)$$

Define then for all $k, l \in \{1, \dots, d\}$,

$$\mathbf{e}_l \cdot K_{\bar{h}_1}^* \mathbf{e}_k = \int_{\mathbb{Y}} (\mathbf{e}_l + \nabla \varphi_l^{\bar{h}_1}) \cdot K(\mathbf{e}_k + \nabla \varphi_k^{\bar{h}_1}).$$

For future reference, we set $\varphi^{\bar{h}_1} = (\varphi_1^{\bar{h}_1}, \dots, \varphi_d^{\bar{h}_1}) \in H_{\#}^1(\mathbb{Y}, \mathbb{R}^d)$.

- (2) Compute the Galerkin approximation $\Theta_0^{\text{h}_0} \in \mathcal{V}_{\Omega, \text{h}_0}^1$ of the solution to (2.7) with $K_{\bar{h}_1}^*$ in place of K^* , unique solution in $\mathcal{V}_{\Omega, \text{h}_0}^1$ to: for all $w \in \mathcal{V}_{\Omega, \text{h}_0}^1$,

$$\int_{\Omega} \nabla w \cdot K_{\bar{h}_1}^* \nabla \Theta_0^{\text{h}_0} = \int_{\Omega} qw.$$

It defines $U_0^{\text{h}_0} = -K_{\bar{h}_1}^* \nabla \Theta_0^{\text{h}_0}$, too.

- (3) Approximation $D_{\text{h}_0}^* \in \mathcal{V}_{\Omega, \text{h}_0}^0$ (each entry of the matrix is piecewise constant on T_{Ω, h_0}) of the homogenized diffusion D^* . Let $\Pi_{\mathcal{V}_{\Omega, \text{h}_0}^0}$ denote the L^2 -projection onto $\mathcal{V}_{\Omega, \text{h}_0}^0$. For every element T of the tessellation T_{Ω, h_0} , $\nabla \Theta_0^{\text{h}_0}|_T$ is constant, and we define $D_{\text{h}_0}^*|_T$ as follows: for all $k, l \in \{1, \dots, d\}$,

$$\mathbf{e}_l \cdot D_{\text{h}_0}^*|_T \mathbf{e}_k = \int_{\mathbb{Y}} (\mathbf{e}_l + \nabla \Phi_l^{\text{h}_1}|_T) \cdot \tilde{D}^{\bar{h}_1}|_T (\mathbf{e}_k + \nabla \Phi_k^{\text{h}_1}|_T),$$

where

$$\tilde{D}^{\bar{h}_1}|_T(y) := \Pi_{\mathcal{V}_{\Omega, \text{h}_0}^0} D_0|_T + \mathbb{D}(\tilde{U}^{\bar{h}_1}|_T(y)),$$

$$\tilde{U}^{\bar{h}_1}|_T(y) := -K(y)(\mathbb{I} + \nabla \varphi^{\bar{h}_1}(y)) \nabla \Theta_0^{\text{h}_0}|_T,$$

and $\Phi_k^{\text{h}_1}|_T \in \mathcal{V}_{\mathbb{Y}, \text{h}_1}^1$ is the unique periodic weak solution to: for all $\Psi \in \mathcal{V}_{\mathbb{Y}, \text{h}_1}^1$,

$$\int_{\mathbb{Y}} \nabla \Psi \cdot \tilde{D}^{\bar{h}_1}|_T (\mathbf{e}_k + \nabla \Phi_k^{\text{h}_1}|_T) = 0. \quad (3.2)$$

- (4) Approximation of C_0 . Let $N \in \mathbb{N}^*$. The time interval $[0, T]$ is uniformly discretized with a fixed time step $\Delta t = \frac{T}{N}$. For all $n \in \{0, \dots, N\}$, we set $t_n = n\Delta t$, and define the approximation $C_0^{\text{h}_0, n} \in \mathcal{V}_{\Omega, \text{h}_0}^1$ of $C_0(t_n, \cdot)$ by induction as the unique solution to: for all $v \in \mathcal{V}_{\Omega, \text{h}_0}^1$,

$$\int_{\Omega} \frac{C_0^{\text{h}_0, n+1} - C_0^{\text{h}_0, n}}{\Delta t} + \int_{\Omega} \nabla v \cdot (D_{\text{h}_0}^* \nabla C_0^{\text{h}_0, n+1} - U_0^{\text{h}_0} C_0^{\text{h}_0, n+1}) + \int_{\Omega} \lambda C_0^{\text{h}_0, n+1} v = \int_{\Omega} S^{n+1} v.$$

Since we use an implicit time discretization, there is no CFL condition — note that we could have used a semi-implicit scheme as well (see for instance ²³).

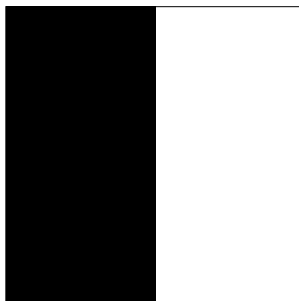


Fig. 1. Laminate structure

In this algorithm we have used two different discretizations $\mathcal{V}_{\mathbb{Y}, h_1}^1$ and $\mathcal{V}_{\mathbb{Y}, \bar{h}_1}^1$ of $H_{\#}^1(\mathbb{Y})$. Indeed, equation (3.2) is an elliptic equation whose diffusion coefficients vary a priori at scale \bar{h}_1 , and it is reasonable to approximate its solutions with a finer discretization parameter $h_1 \leq \bar{h}_1$.

3.1.2. Numerical tests

To illustrate Theorem 2 when the homogenized system is solved using the direct approach of Algorithm 1, we consider a numerical test suggested by ANDRA^a. We take $d = 2$ and let $\Omega =]0, 2[^2$ be a square domain, and $[0, T]$ be the time interval with $T = 1$. The permeability is defined on the domain $\mathbb{Y} =]0, 1[^2$ by:

$$\forall y = (y_1, y_2) \in \mathbb{Y}, \forall y_1 \in]0, 1[, K(y_1, y_2) = \begin{cases} 4.94064, & \text{if } y_2 \geq 0.5, \\ 0.57816, & \text{if } y_2 < 0.5. \end{cases}$$

It has a laminate structure (see Figure 1). We consider boundary conditions which are slightly different than in Theorem 2 and Algorithm 1 — note that the adaptations are straightforward in both cases — :

$$\left\{ \begin{array}{l} \text{Dirichlet boundary conditions: Let } x = (x_1, x_2) \in \partial\Omega \\ \text{For } x_1 \in (0, 2), h_0(x_1, x_2) = \begin{cases} \frac{5}{3}, & \text{if } x_2 = 2, \\ \frac{5}{3} + 0.5, & \text{if } x_2 = 0, \end{cases} \\ \text{For } x_2 \in (0, 2), C_0(x_1, x_2) = \begin{cases} 1, & \text{if } x_1 = 0, \\ 0, & \text{if } x_1 = 2. \end{cases} \\ \text{Homogeneous Neumann boundary conditions elsewhere.} \end{array} \right. \quad (3.3)$$

The parameters used in the numerical tests are gathered in Table 1. As a consequence of the laminate structure, the correctors φ_1 and φ_2 , and Φ_1 and Φ_2 belong

^a Agence nationale pour la gestion des déchets radioactifs — <http://www.andra.fr>

$T = 1$	$D_0 = 4.38 \mathbb{I}$	$\alpha = \frac{2}{10}$	$\beta = \frac{2}{100}$	$\lambda = \frac{\ln(2)}{1.57}$	$\Delta t = 10^{-3}$
---------	-------------------------	-------------------------	-------------------------	---------------------------------	----------------------

Table 1. Parameters

to the finite element space $\mathcal{V}_{\mathbb{Y}, h_1}^1$ provided the geometry of $T_{\mathbb{Y}, h_1}$ matches the laminate structure of Figure 1. In this case, one can therefore take $h_1 = \bar{h}_1$. In Table 2 we compare the approximations $(\Theta_\varepsilon^{\varepsilon h_1}, C_\varepsilon^{\varepsilon h_1})$ of the solutions to the heterogeneous system (2.4) to the approximation $(\Theta_0^{h_0}, C_0^{h_0})$ of the solution to the homogenized system (2.7)–(2.12), for several values of ε (the discretization parameters h_1 and h_0 being fixed). The periodic cell \mathbb{Y} is discretized with 8 elements per dimension, and the macroscopic domain Ω is discretized using $2 \times 8/\varepsilon$ elements per dimension to compute $(\Theta_\varepsilon^{\varepsilon h_1}, C_\varepsilon^{\varepsilon h_1})$. For the approximation of $(\Theta_0^{h_0}, C_0^{h_0})$, we take $h_0 = 1/100$. This yields

- $\mathcal{V}_{\mathbb{Y}, h_1}^1$ has dimension 81;
- $\mathcal{V}_{\Omega, h_0}^1$ has dimension ~ 40000 ;
- $\mathcal{V}_{\Omega, h_0, \varepsilon}^1$ has dimension $\sim 256\varepsilon^{-2}$.

We display in Table 2 the $L^2(\Omega)$ norm of the error $\Theta_0^{h_0} - \Theta_\varepsilon^{\varepsilon h_1}$ and the $L^2(\Omega \times]0, T[)$ -norm of the error $C_0^{h_0} - C_\varepsilon^{\varepsilon h_1}$ for $\varepsilon \in \{0.2, 0.1, 0.05, 0.025\}$, that is, we compare solutions to the homogenized problems to approximations of the solution resolving the ε -scale (the case $\varepsilon = 0.025$ is already borderline in terms of computational cost). These results have been obtained using **FreeFem++** (see ¹⁹). The linear systems are solved with a direct solver. We obtain a first order of convergence for both errors.

ε	$\frac{\ h_0^{h_0} - h_\varepsilon^{\varepsilon h_1}\ _{L^2}}{\ h_0^{h_0}\ _{L^2}}$	Rate	$\frac{\ C_0^{h_0} - C_\varepsilon^{\varepsilon h_1}\ _{L^2(L^2)}}{\ C_0^{h_0}\ _{L^2(L^2)}}$	Rate
0.2	1.667e-3	-	5.528e-4	-
0.1	8.095e-4	1.04	2.525e-4	1.13
0.05	3.992e-4	1.02	1.270e-4	0.99
0.025	1.983e-4	1.01	6.704e-5	0.92

Table 2. Error in function of ε

As can be seen on Table 2 the apparent convergence rates are of order 1, which is consistent with a formal two-scale expansion, and shows the interest of replacing $(\Theta_\varepsilon, C_\varepsilon)$ by its homogenized counterpart (h_0, C_0) . Although the computational time for the approximation of (Θ_0, C_0) is much smaller than the computational time for the approximation of $(\Theta_\varepsilon, C_\varepsilon)$ when ε is small, this method rapidly becomes

prohibitive when the tessellation of \mathbb{Y} gets finer since the approximation of D^* then becomes quite expensive.

The last part of this article is devoted to the speed up of the approximation of D^* , with a numerical cost which should ideally be independent of the meshsize h_1 of $\mathbb{T}_{\mathbb{Y}, h_1}$. From now on we assume \tilde{D} in (2.10) to be a symmetric matrix (that is, we assume D^0 to be symmetric).

3.2. Reduced basis method for homogenization problems

In this section we describe how to apply the reduced basis method to the homogenized problem under consideration, assuming in addition that D_0 in (2.2) is a constant matrix.

3.2.1. General presentation

The reduced basis method was introduced for the accurate online evaluation of (outputs of) solutions to a parameter-dependent family of elliptic PDEs. The basis of the method and further references can be found in ²¹. The application to the homogenization of elliptic equations is discussed in ⁵. Abstractly, it can be viewed as a method to determine a “good” N -dimensional space \mathcal{S}_N to be used in approximating the elements of a set $\mathcal{F} = \{(\bar{\Phi}_1(\xi), \dots, \bar{\Phi}_d(\xi)), \xi \in \mathcal{P}\}$ of parametrized elements lying in a Hilbert space \mathcal{S} , the parameter ξ ranging a certain subset $\mathcal{P} \subset \mathbb{R}^n$.

Let us describe how the computation of the effective coefficients we are concerned with belongs to such a framework. First of all, the auxiliary function Θ_0 is simply determined by solving the problem (2.7), with effective coefficients obtained by solving the cell equations (2.5). There is no difficulty in this step and $\nabla_x \Theta_0$ can be considered as given in this discussion. Then, we write the effective coefficient (2.10) for the concentration equation (2.12) as follows

$$\tilde{D}(x, y) = \hat{\mathcal{D}}(\nabla_x \Theta_0(x))(y)$$

where $\xi \in \mathbb{R}^d \mapsto \hat{\mathcal{D}}(\xi) \in L^\infty(\mathbb{Y}, \mathcal{M}_d(\mathbb{R}))$ is defined by

$$\begin{aligned} \hat{\mathcal{D}}(\xi)(y) &= D_0 + \alpha |M(y)\xi| \mathbb{I} + \beta \frac{M(y)\xi \otimes M(y)\xi}{|M(y)\xi|} = D_0 + \mathbb{D}(M(y)\xi), \\ M(y) &= K(y)(\mathbb{I} + \nabla \varphi(y)), \\ \varphi &= (\varphi_1, \dots, \varphi_d) \quad \text{solutions of (2.5)}. \end{aligned} \tag{3.4}$$

We recall that $\alpha, \beta \geq 0$, and D_0 is a positive-definite symmetric matrix while $M : \mathbb{Y} \rightarrow \mathcal{M}_d(\mathbb{R})$ is a square-integrable function. We are interested in the solution $\bar{\Phi}_k(\xi) \in H_{\#}^1(\mathbb{Y})$ to the problem: for all $\bar{\Psi} \in H_{\#}^1(\mathbb{Y})$,

$$\int_{\mathbb{Y}} \nabla \bar{\Psi}(y) \cdot \hat{\mathcal{D}}(\xi)(y)(\mathbf{e}_k + \nabla \bar{\Phi}_k(\xi)(y)) \, dy = 0.$$

In the present context, $S = \mathcal{H}_{\#}^{\infty}(\mathbb{Y})$ and we wish to find a convenient finite dimensional approximation space \mathcal{S}_N which allows to describe the set \mathcal{F} of solutions. In the rest of this paragraph we particularize the standard RB method to homogenization problems, by choosing a specific error estimator and orthogonalization procedure. To avoid further heavy notation, we do not display the variable y in what follows.

Let $n \geq 1$ and let $\mathcal{D} : \mathbb{R}^n \rightarrow L^{\infty}(\mathbb{Y}, \mathcal{M}_d(\mathbb{R}))$ be a function taking values in a set of $(y \in \mathbb{Y})$ -dependent $d \times d$ symmetric real matrices, satisfying uniform bounds and elliptic estimates. We suppose that $\mathcal{D}(\xi)$ depends continuously on the parameter $\xi \in \mathbb{R}^n$. Given a compact subset $\bar{\mathcal{K}}$ of \mathbb{R}^n , we set $\mathcal{F}_{\bar{\mathcal{K}}} = \{(\bar{\Phi}_1(\xi), \dots, \bar{\Phi}_d(\xi)), \xi \in \bar{\mathcal{K}}\}$, where $\bar{\Phi}_k(\xi) \in H_{\#}^1(\mathbb{Y})$ denotes the unique periodic weak solution to the problem: for all $\bar{\Psi} \in H_{\#}^1(\mathbb{Y})$, and for all $k \in \{1, \dots, d\}$,

$$\int_{\mathbb{Y}} \nabla \bar{\Psi}(\xi) \cdot \mathcal{D}(\xi)(\mathbf{e}_k + \nabla \bar{\Phi}_k(\xi)) = 0.$$

The set $\mathcal{F}_{\bar{\mathcal{K}}}$ is therefore compact in $H_{\#}^1(\mathbb{Y})$.

To construct the N -finite dimensional space \mathcal{S}_N intended to approximate $\mathcal{F}_{\bar{\mathcal{K}}}$, we proceed by induction using a greedy algorithm. To this aim we need a reliable estimator which measures the error between $\bar{\Phi}_k(\xi)$ for some $\xi \in \bar{\mathcal{K}}$ and its approximation $\bar{\Phi}_k^j(\xi)$ in \mathcal{S}_j for $0 \leq j \leq N$, which is defined as the unique weak solution $\bar{\Phi}_k^j(\xi) \in \mathcal{S}_j$ to: for all $\bar{\Psi}^j \in \mathcal{S}_j$,

$$\int_{\mathbb{Y}} \nabla \bar{\Psi}^j(\xi) \cdot \mathcal{D}(\xi)(\mathbf{e}_k + \nabla \bar{\Phi}_k^j(\xi)) = 0. \quad (3.5)$$

Recalling that we are dealing with a homogenization problem, the quantity of interest is the symmetric homogenized matrix $\bar{D}^*(\xi)$ defined for all $k, l \in \{1, \dots, d\}$ by (see (2.11))

$$\mathbf{e}_l \cdot \bar{D}^*(\xi) \mathbf{e}_k = \int_{\mathbb{Y}} (\mathbf{e}_l + \nabla \bar{\Phi}_l(\xi)) \cdot \mathcal{D}(\xi)(\mathbf{e}_k + \nabla \bar{\Phi}_k(\xi)).$$

We denote by $\bar{D}^{*,j}(\xi)$ the approximation of $\bar{D}^*(\xi)$ using \mathcal{S}_j , that is for all $k, l \in \{1, \dots, d\}$

$$\mathbf{e}_l \cdot \bar{D}^{*,j}(\xi) \mathbf{e}_k = \int_{\mathbb{Y}} (\mathbf{e}_l + \nabla \bar{\Phi}_l^j(\xi)) \cdot \mathcal{D}(\xi)(\mathbf{e}_k + \nabla \bar{\Phi}_k^j(\xi)).$$

A standard calculation using (3.5) and the symmetry of \mathcal{D} yields

$$\mathbf{e}_l \cdot (\bar{D}^*(\xi) - \bar{D}^{*,j}(\xi)) \mathbf{e}_k = \int_{\mathbb{Y}} (\nabla \bar{\Phi}_l(\xi) - \nabla \bar{\Phi}_l^j(\xi)) \cdot \mathcal{D}(\xi)(\nabla \bar{\Phi}_k(\xi) - \nabla \bar{\Phi}_k^j(\xi)).$$

This shows that the error on the homogenized matrix is a suitable estimator of the error at the level of the $\bar{\Phi}_k$. We thus define the estimator $\bar{\mathcal{E}}^j : \bar{\mathcal{K}} \times \{1, \dots, d\} \rightarrow \mathbb{R}^+$ by

$$\bar{\mathcal{E}}^j(\xi, k) = \sqrt{\frac{|\mathbf{e}_k \cdot (\bar{D}^*(\xi) - \bar{D}^{*,j}(\xi)) \mathbf{e}_k|}{\mathbf{e}_k \cdot \bar{D}^*(\xi) \mathbf{e}_k}}.$$

So defined, and recalling that \mathcal{D} is assumed to take values in the set of uniformly elliptic symmetric matrices (say with ellipticity constants $0 < \underline{\nu} \leq \bar{\nu} < \infty$), the estimator is such that there exist $C_1, C_2 > 0$ verifying for all suitable j, k, ξ the inequality

$$C_1 \bar{\mathcal{E}}^j(\xi, k) \leq \|\nabla \bar{\Phi}_k(\xi) - \nabla \bar{\Phi}_k^j(\xi)\|_{L^2(\mathbb{Y})} \leq C_2 \bar{\mathcal{E}}^j(\xi, k). \quad (3.6)$$

The induction procedure is then as follows. For all $j \in \{0, \dots, N-1\}$, choose $\xi_{j+1} \in \bar{\mathcal{K}}$ and $k_{j+1} \in \{1, \dots, d\}$ such that

$$(\xi_{j+1}, k_{j+1}) = \operatorname{argmax}_{\bar{\mathcal{K}}, \{1, \dots, d\}} \bar{\mathcal{E}}^j,$$

define

$$\bar{\Psi}_{j+1} = \frac{\bar{\Phi}_{k_{j+1}}(\xi_{j+1}) - \bar{\Phi}_{k_{j+1}}^j(\xi_{j+1})}{\|\nabla \bar{\Phi}_{k_{j+1}}(\xi_{j+1}) - \nabla \bar{\Phi}_{k_{j+1}}^j(\xi_{j+1})\|_{L^2(\mathbb{Y})}}$$

and set

$$\mathcal{S}_{j+1} := \operatorname{span} \{\bar{\Psi}_1, \dots, \bar{\Psi}_{j+1}\}.$$

By induction, for all $j \in \{0, \dots, N\}$, $\dim \mathcal{S}_{j+1} = j+1$, since by construction $\bar{\Psi}_{j+1}$ is orthogonal to \mathcal{S}_j for the following scalar product of $H_{\#}^1(\mathbb{Y})$

$$(\bar{\Psi}_1, \bar{\Psi}_2) \mapsto \int_{\mathbb{Y}} \nabla \bar{\Psi}_1 \cdot \mathcal{D}(\xi_{j+1}) \nabla \bar{\Psi}_2. \quad (3.7)$$

Note that usually, in the RB literature, the vectors $\bar{\Psi}_j$ are orthogonalized using the same scalar product for all j (whereas here, the scalar product depends on j). The choice made here makes the computation of the reduced basis simpler (and the generated space \mathcal{S}_{j+1} is the same). The influence of this choice in practice is investigated numerically in Paragraph 3.3.3.

What convergence rate can be expected in terms of N ? In the case when $n = 1$ and \mathcal{D} has a dependence of the form $\mathcal{D}(\xi) = D_0 + \xi D_1$ (that is \mathcal{D} is an affine function on the real line, and $\bar{\mathcal{K}}$ is just a segment), the combination of results from ²¹ (see also the more general case treated in ¹²) and ³ (see also ⁷) shows that there exist $c, C > 0$ such that for all $N \geq 1$,

$$\sup_{\bar{\mathcal{K}}, \{1, \dots, d\}} \bar{\mathcal{E}}^N \leq C \exp(-cN).$$

The convergence being exponential in N , the reduced basis method is expected to yield accurate results for moderate N (say for N which is much smaller than the dimension of the finite element space $\mathcal{V}_{\mathbb{Y}, h_1}^1$ for instance). Note that this yields a complete control of the error on the homogenized coefficients since for all $k, l \in$

$\{1, \dots, d\}$ by Cauchy-Schwarz' inequality and definition of the estimator,

$$\begin{aligned} |\mathbf{e}_l \cdot (\overline{D}^*(\xi) - \overline{D}^{*,N}(\xi))\mathbf{e}_k| &= \left| \int_{\mathbb{Y}} (\nabla \overline{\Phi}_l(\xi) - \nabla \overline{\Phi}_l^N(\xi)) \cdot \mathcal{D}(\xi) (\nabla \overline{\Phi}_k(\xi) - \nabla \overline{\Phi}_k^N(\xi)) \right| \\ &\leq \sqrt{(\mathbf{e}_l \cdot \overline{D}^*(\xi)\mathbf{e}_l)(\mathbf{e}_k \cdot \overline{D}^*(\xi)\mathbf{e}_k)} \overline{\mathfrak{E}}^N(\xi, k) \overline{\mathfrak{E}}^N(\xi, l) \\ &\leq \tilde{C} \exp(-2cN) \end{aligned}$$

for some constant \tilde{C} depending only on C and $d, \underline{\nu}, \overline{\nu}$.

In the case under consideration here, \mathcal{D} is replaced by $\widehat{\mathcal{D}}$ defined in (3.4). Things are more complex than in ²¹ and ³ for the following three reasons:

- the parameter ξ is in \mathbb{R}^d (that is $n = d > 1$ in the case of interest);
- a priori $\overline{\mathcal{K}} = \mathbb{R}^d$, which is not a compact set;
- the function $\xi \mapsto \widehat{\mathcal{D}}(\xi)$ is nonlinear.

More generally, our working plan faces the following technical difficulties:

- the parameter ξ ranges over the whole \mathbb{R}^d while the method is designed to deal with parameters lying in a compact set.
- the method simplifies significantly when the dependence of \mathcal{D} upon ξ is affine. In such a case it is described and analyzed in full details, whereas here the dependence with respect to the parameter is more intricate. The implementation of the method will require additional devices.
- the matrix M arising in the definition (3.4) of $\widehat{\mathcal{D}}$ is not essentially bounded as a function of $y \in \mathbb{Y}$, but square-integrable only. Therefore the available results that could be used to analyze the method simply do not apply.

The algorithm described in this paragraph is not of any practical use yet since in order to choose ξ_{j+1} and k_{j+1} , one needs to know $\overline{\mathfrak{E}}^j(\xi, k)$ for all ξ and k . In the following paragraph we describe the standard way to proceed in practice.

3.2.2. Practical reduced basis method

In practice we do not have access to $\{\overline{\mathfrak{E}}^j(\xi, k), \xi \in \overline{\mathcal{K}}, k \in \{1, \dots, k\}\}$ since:

- the corrector $\overline{\Phi}_k(\xi)$ has to be approximated in a finite-dimensional subspace \mathcal{V} of $H_{\#}^1(\mathbb{Y})$, so that $\overline{\mathfrak{E}}^j$ is approximated by some \mathfrak{E}^j .
- the space $\overline{\mathcal{K}}$ has to be replaced by some finite set \mathcal{K} .

The construction of the reduced basis is then as follows.

Algorithm 2. Let $N \in \mathbb{N}$, $p \geq N$, $\mathcal{K} = \{\xi_m, m \in \{1, \dots, p\}\}$ be a finite subset of $\overline{\mathcal{K}}$, and \mathcal{V} be a finite-dimensional subspace of $H_{\#}^1(\mathbb{Y})$.

- (1) For all $m \in \{1, \dots, p\}$ and $k \in \{1, \dots, d\}$, let $\Phi_k(\xi_m) \in \mathcal{V}$ be an approximation of the corrector $\overline{\Phi}_k(\xi_m)$ in \mathcal{V} , that is the unique element of \mathcal{V} such that for all

22

$\Psi \in \mathcal{V}$

$$\int_{\mathbb{Y}} \nabla \Psi \cdot \mathcal{D}(\xi_m)(\mathbf{e}_k + \nabla \Phi_k(\xi_m)) = 0,$$

and let $D_{kk}^*(\xi_m)$ be the approximation of $\mathbf{e}_k \cdot \bar{D}^*(\xi_m)\mathbf{e}_k$ given by

$$D_{kk}^*(\xi_m) = \int_{\mathbb{Y}} (\mathbf{e}_k + \nabla \Phi_k(\xi_m)) \cdot \mathcal{D}(\xi_m)(\mathbf{e}_k + \nabla \Phi_k(\xi_m)).$$

(2) Set $\mathcal{V}_0 = \{0\}$.

(3) While $0 \leq j < N$

(a) For all $m \in \{1, \dots, p\}$ and $k \in \{1, \dots, d\}$, let $\Phi_k^j(\xi_m) \in \mathcal{V}$ be an approximation of the corrector $\Phi_k(\xi_m)$ in \mathcal{V}_j , that is the unique element of \mathcal{V}_j such that for all $\Psi^j \in \mathcal{V}_j$

$$\int_{\mathbb{Y}} \nabla \Psi^j \cdot \mathcal{D}(\xi_m)(\mathbf{e}_k + \nabla \Phi_k^j(\xi_m)) = 0,$$

and let $D_{kk}^{*,j}(\xi_m)$ be the approximation of $D_{kk}^*(\xi_m)$ given by

$$D_{kk}^{*,j}(\xi_m) = \int_{\mathbb{Y}} (\mathbf{e}_k + \nabla \Phi_k^j(\xi_m)) \cdot \mathcal{D}(\xi_m)(\mathbf{e}_k + \nabla \Phi_k^j(\xi_m)).$$

(b) For all $m \in \{1, \dots, p\}$ and $k \in \{1, \dots, d\}$, define the estimator $\mathfrak{E}^j(m, k)$ as

$$\mathfrak{E}^j(m, k) = \sqrt{\frac{|\mathbf{e}_k \cdot D^*(\xi_m)\mathbf{e}_k - D_{kk}^{*,j}(\xi_m)|}{\mathbf{e}_k \cdot D^*(\xi_m)\mathbf{e}_k}},$$

and set

$$(m_j, k_j) = \operatorname{argmax}_{\mathcal{K}, \{1, \dots, d\}} \mathfrak{E}^j(m, k).$$

(c) Define

$$\Psi_{j+1} := \frac{\Phi_{k_j}(\xi^m) - \Phi_{k_j}^j(\xi^m)}{\|\nabla \Phi_{k_j}(\xi^m) - \nabla \Phi_{k_j}^j(\xi^m)\|_{L^2(\mathbb{Y})}},$$

and set

$$\mathcal{V}_{j+1} = \operatorname{span} \{\Psi_i, 1 \leq i \leq j+1\}.$$

(d) $j = j + 1$.

Provided p is chosen large enough and \mathcal{V} has dimension larger than N , one has as in the previous paragraph $\dim \mathcal{V}_N = N$.

What convergence rate can be expected in terms of N ? Going back to the example mentioned in the previous paragraph, that is for $\mathcal{D}(\xi) = D_0 + \xi D_1$ and $\bar{\mathcal{K}}$ a segment, the answer is given in ³. In particular it is proved that the exponential estimate is stable in the sense that if the reduced basis \mathcal{V}_N is constructed starting

from approximations of the correctors $\{\bar{\Phi}_k(\xi), \xi \in \bar{\mathcal{K}}, k \in \{1, \dots, d\}\}$ within an error e , then the error estimate is of the form

$$\sup_{\bar{\mathcal{K}}, \{1, \dots, d\}} \bar{\mathcal{E}}^N \leq C \exp(-cN) + Ce.$$

In Algorithm 2 there are two origins for the error e :

- The fact that $H_{\#}^1(\mathbb{Y})$ is replaced by a finite-dimensional space \mathcal{V} , so that for all $\xi \in \mathcal{K}$ and $k \in \{1, \dots, d\}$, $\Phi_k(\xi)$ is a finite-dimensional approximation of $\bar{\Phi}_k$;
- The fact that for the greedy algorithm, the argmax of the estimator is taken in \mathcal{K} and not in $\bar{\mathcal{K}}$.

The first source of error is standard and can be controlled by a priori or a posteriori estimates. In the affine case above, the second source of error can also be estimated. Indeed, as proved in ¹³, the maps $\bar{\Phi}_k : \bar{\mathcal{K}} \rightarrow H_{\#}^1(\mathbb{Y}), \xi \mapsto \bar{\Phi}_k(\xi)$ are analytic for all $k \in \{1, \dots, d\}$. In particular, if \mathcal{K} is a sampling of $\bar{\mathcal{K}}$ with “meshsize” h , for all $\xi \in \bar{\mathcal{K}}$, $\bar{\Phi}_k(\xi)$ can be approximated by interpolation in $\{\bar{\Phi}_k(\xi_m), m \in \{1, \dots, p\}\}$ within an error of order h^q for any $q \in \mathbb{N}$. Hence the practical reduced basis method remains efficient in this specific case. However, this analysis is restricted to the affine case and it does not apply in our context.

3.2.3. Fast-assembly of the matrix

Let $\bar{\mathcal{K}}, \mathcal{K}$, and $N \in \mathbb{N}$ and \mathcal{V}_N be as in Algorithm 2. For all $\xi \in \bar{\mathcal{K}}$ and $k \in \{1, \dots, d\}$, the approximation of $\bar{\Phi}_k(\xi)$ in the reduced basis \mathcal{V}_N is given by the unique function $\Phi_k^N(\xi) \in \mathcal{V}_N$ such that for all $\Psi^N \in \mathcal{V}_N$,

$$\int_{\mathbb{Y}} \nabla \Psi^N \cdot \mathcal{D}(\xi)(\mathbf{e}_k + \nabla \Phi_k^N(\xi)) = 0. \quad (3.8)$$

Expanding $\Phi_k^N(\xi)$ in the basis \mathcal{V}_N as $\Phi_k^N(\xi) = \sum_{j=1}^N u_j(\xi) \Psi_j$, the above equation is equivalent to the linear system

$$\mathbb{M}(\xi)U = B(\xi, k),$$

where for all $j \in \{1, \dots, N\}$, $U_j = u_j(\xi)$ and $B(\xi, k)_j = -\int_{\mathbb{Y}} \nabla \Psi_j \cdot \mathcal{D}(\xi) \mathbf{e}_k$, and the $N \times N$ matrix $\mathbb{M}(\xi)$ is given by its entries $\mathbb{M}(\xi)_{j_1 j_2} = \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot \mathcal{D}(\xi) \nabla \Psi_{j_2}$ for all $1 \leq j_1, j_2 \leq N$. In particular, in order to compute $\Phi_k^N(\xi)$, one needs to solve (3.8), and therefore construct the matrix $\mathbb{M}(\xi)$ and the r. h. s. $B(\xi, k)$.

Without further assumption on the function $\xi \mapsto \mathcal{D}(\xi)$, the exact calculation of $\mathbb{M}(\xi)$ and $B(\xi, k)$ requires:

- the storage of the coordinates of each vector Ψ_j of the reduced basis \mathcal{V}_N in the finite dimensional space \mathcal{V} ,
- the computation of integrals over \mathbb{Y} .

Both the information to be stored and the computational cost to construct $\mathbb{M}(\xi)$ and $B(\xi, k)$ scale like the dimension $\dim(\mathcal{V})$ of the finite-dimensional space \mathcal{V} (which can be prohibitively large). Yet, if $\xi \mapsto \mathcal{D}(\xi)$ has specific structural properties, the information to be stored and the computational cost can be drastically reduced. This is the case when $\xi \mapsto \mathcal{D}(\xi)$ is affine. Let us go back to the example of $\mathcal{D}(\xi) = D_0 + \xi D_1$. Then, the entries of the matrix $\mathbb{M}(\xi)$ and of the r. h. s. $B(\xi, j)$ take the form: for all $1 \leq j, j_1, j_2 \leq N$,

$$\begin{aligned}\mathbb{M}(\xi)_{j_1 j_2} &= \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot D_0 \nabla \Psi_{j_2} + \xi \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot D_1 \nabla \Psi_{j_2}, \\ B(\xi, k)_j &= \int_{\mathbb{Y}} \nabla \Psi_j \cdot D_0 \mathbf{e}_k + \xi \int_{\mathbb{Y}} \nabla \Psi_j \cdot D_1 \mathbf{e}_k.\end{aligned}$$

In particular, provided we store the following two $N \times N$ matrices M_1 and M_2 , and the following two $d \times N$ matrices B_1 and B_2 defined by: for all $1 \leq j, j_1, j_2 \leq N$ and $k \in \{1, \dots, d\}$,

$$\begin{aligned}(M_1)_{j_1 j_2} &= \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot D_0 \nabla \Psi_{j_2}, & (M_2)_{j_1 j_2} &= \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot D_1 \nabla \Psi_{j_2}, \\ (B_1)_{kj} &= \int_{\mathbb{Y}} \nabla \Psi_j \cdot D_0 \mathbf{e}_k, & (B_2)_{kj} &= \int_{\mathbb{Y}} \nabla \Psi_j \cdot D_1 \mathbf{e}_k,\end{aligned}$$

one may reconstruct $\mathbb{M}(\xi)$ and $B(\xi, k)$ by the simple formulae

$$\mathbb{M}(\xi) = M_1 + \xi M_2, \quad B(\xi, k) = (B_1)_k + \xi (B_2)_k,$$

where $(B_1)_k$ and $(B_2)_k$ are the k -th column of B_1 and B_2 , respectively.

The gain is twofold:

- the dimension of the information to store is $2N^2 + 2dN$, which is independent of $\dim(\mathcal{V})$,
- the computation of $\mathbb{M}(\xi)$ and $B(\xi, k)$ only requires $N^2 + N$ multiplications and $N^2 + N$ additions, and not the computation of $N^2 + N$ integrals on \mathbb{Y} (using an integration rule which should be exact for functions of \mathcal{V}).

The same strategy allows one to easily compute the approximation $D^{*,N}(\xi)$ of the homogenized matrix $\bar{D}^*(\xi)$, via the formula: for all $k, l \in \{1, \dots, d\}$,

$$\begin{aligned}\mathbf{e}_k \cdot D^{*,N}(\xi) \mathbf{e}_l &= \int_{\mathbb{Y}} (\mathbf{e}_k + \nabla \Phi_k^N(\xi)) \cdot \mathcal{D}(\xi) (\mathbf{e}_l + \nabla \Phi_l^N(\xi)) \\ &= \mathbf{e}_k \cdot \left(\int_{\mathbb{Y}} D_0 + \xi \int_{\mathbb{Y}} D_1 \right) \mathbf{e}_l \\ &\quad + \sum_{j=1}^N u_j(\xi) \mathbf{e}_k \cdot \left(\int_{\mathbb{Y}} D_0 \nabla \Psi_j + \xi \int_{\mathbb{Y}} D_1 \nabla \Psi_j \right),\end{aligned}$$

so that one has to store $2d^2 + 2dN$ real numbers only, to compute the approximation of the homogenized matrix.

This fast-assembly method is very convenient and efficient, but requires the diffusion matrix $\mathcal{D}(\xi)$ to be affine with respect to ξ .

3.3. Application of the reduced basis method to the homogenized system

As said above, the evaluation of the effective coefficients for the homogenized problem involves the parametrized matrices $\widehat{\mathcal{D}}$ defined in (3.4) where the parameter ξ ranges the unbounded set \mathbb{R}^d . In the following paragraph we shall rewrite the problem in an equivalent form which allows one to work with a compact set of parameters. We address the issue of fast assembly in the second paragraph, bearing in mind that the dependence with respect to the parameter is not affine. We provide with a numerical study of the method in the last paragraph.

3.3.1. Rewriting of the problem

The starting point to rewrite the problem is the following observation: for all $\xi \in \mathbb{R}^d$ and all $k \in \{1, \dots, d\}$, the corrector $\overline{\Phi}_k(\xi) \in H_{\#}^1(\mathbb{Y})$ is solution to

$$-\nabla \cdot \frac{\widehat{\mathcal{D}}(\xi)}{1 + |\xi|} (\mathbf{e}_k + \nabla \overline{\Phi}_k(\xi)) = 0. \quad (3.9)$$

Let S^{d-1} denote the unit hypersphere in dimension d . Define

$$\begin{aligned} \mathcal{D} : [0, 1] \times S^{d-1} &\longrightarrow L^2(\mathbb{Y}, \mathcal{M}_d(\mathbb{R})) \\ (\rho, X) &\longmapsto \mathcal{D}(\rho, X) \end{aligned}$$

by

$$\mathcal{D}(\rho, X) : y \mapsto (1 - \rho)D_0 + \rho \left(\alpha |M(y)X| \mathbb{I} + \beta \frac{M(y)X \otimes M(y)X}{|M(y)X|} \right), \quad (3.10)$$

For all $(\rho, X) \in [0, 1] \times S^{d-1}$ and $k \in \{1, \dots, d\}$, we let $\overline{\Phi}_k(\rho, X)$ be the unique weak solution in $H_{\#}^1(\mathbb{Y})$ to

$$-\nabla \cdot \mathcal{D}(\rho, X) (\mathbf{e}_k + \nabla \overline{\Phi}_k(\rho, X)) = 0. \quad (3.11)$$

Let $\xi \in \mathbb{R}^d$, and set

$$\rho = \frac{|\xi|}{1 + |\xi|}, \quad X = \frac{\xi}{|\xi|}$$

so that

$$\frac{\widehat{\mathcal{D}}(\xi)}{1 + |\xi|} = \mathcal{D}(\rho, X);$$

the identities (3.9) and (3.11) imply that

$$\overline{\Phi}_k(\xi) \equiv \overline{\Phi}_k(\rho, X)$$

by uniqueness of correctors. In particular, this shows that

$$\left\{ \overline{\Phi}_k(\xi), \xi \in \mathbb{R}^d, k \in \{1, \dots, d\} \right\} = \left\{ \overline{\Phi}_k(\rho, X), (\rho, X) \in [0, 1] \times S^{d-1}, k \in \{1, \dots, d\} \right\}.$$

What we gain by applying the reduced basis method on this new formulation is that the parameters now belong to closed unit ball $[0, 1] \times S^{d-1}$.

To complete the description of the RB method, we need to choose an estimator. We shall make use of the estimator defined in the previous subsection. Let $j \in \mathbb{N}$ and let \mathcal{V}_j be a subspace of $H_{\#}^1(\mathbb{Y})$ of dimension j . Set for all $(\rho, X) \in [0, 1] \times S^{d-1}$ and $k \in \{1, \dots, d\}$,

$$\bar{\mathfrak{E}}^j(\rho, X, k) = \sqrt{\frac{|\mathbf{e}_k \cdot (\bar{D}^*(\rho, X) - \bar{D}^{*,j}(\rho, X))\mathbf{e}_k|}{\mathbf{e}_k \cdot \bar{D}^*(\rho, X)\mathbf{e}_k}}, \quad (3.12)$$

where, denoting by $\bar{\Phi}_k^j(\rho, X)$ the approximation of $\bar{\Phi}_k(\rho, X)$ in \mathcal{V}_j , we have

$$\begin{aligned} \mathbf{e}_k \cdot \bar{D}^*(\rho, X)\mathbf{e}_k &= \int_{\mathbb{Y}} (\mathbf{e}_k + \nabla \bar{\Phi}_k(\rho, X)) \cdot \mathcal{D}(\rho, X)(\mathbf{e}_k + \nabla \bar{\Phi}_k(\rho, X)), \\ \mathbf{e}_k \cdot \bar{D}^{*,j}(\rho, X)\mathbf{e}_k &= \int_{\mathbb{Y}} (\mathbf{e}_k + \nabla \bar{\Phi}_k^j(\rho, X)) \cdot \mathcal{D}(\rho, X)(\mathbf{e}_k + \nabla \bar{\Phi}_k^j(\rho, X)). \end{aligned} \quad (3.13)$$

Note that this estimator is consistent with the estimator associated with $\widehat{\mathcal{D}}$ since we have for all $\xi \in \mathbb{R}^d$,

$$\widehat{\mathfrak{E}}^j(\xi, k) = \bar{\mathfrak{E}}^j(\rho, X, k)$$

for $\rho = \frac{|\xi|}{1+|\xi|}$ and $X = \frac{\xi}{|\xi|}$, the estimator $\widehat{\mathfrak{E}}^j(\xi, k)$ (and the matrices $\widehat{D}^*(\xi)$, $\widehat{D}^{*,j}(\xi)$) being defined with the matrix $\widehat{\mathcal{D}}(\xi)$. Since we also have for all $\xi \in \mathbb{R}^d$

$$\begin{aligned} \bar{D}^*(\rho, X) &= \frac{1}{1+|\xi|} \widehat{D}^*(\xi), \\ \bar{D}^{*,j}(\rho, X) &= \frac{1}{1+|\xi|} \widehat{D}^{*,j}(\xi), \end{aligned}$$

for $\rho = \frac{|\xi|}{1+|\xi|}$ and $X = \frac{\xi}{|\xi|}$, it is equivalent to approximate \bar{D}^* and $\bar{D}^{*,j}$. We will focus on the former in what follows.

Before we turn to fast-assembly, let us make a comment of the RB method used here. The estimator (3.12) satisfies the second inequality of (3.6), namely there exists $C_2 > 0$ such that for all $j \in \mathbb{N}$, $(\rho, X) \in [0, 1] \times S^{d-1}$, and $k \in \{1, \dots, d\}$,

$$\|\nabla \bar{\Phi}_k(\rho, X) - \nabla \bar{\Phi}_k^j(\rho, X)\|_{L^2(\mathbb{Y})} \leq C_2 \bar{\mathfrak{E}}^j(\rho, X, k).$$

Yet the converse inequality only holds in a weaker sense. In particular, using that $\bar{D}^*(\rho, X)$ and $\bar{D}^{*,j}(\rho, X)$ can be defined as

$$\begin{aligned} \mathbf{e}_k \cdot \bar{D}^*(\rho, X)\mathbf{e}_k &= \int_{\mathbb{Y}} \mathbf{e}_k \cdot \mathcal{D}(\rho, X)(\mathbf{e}_k + \nabla \bar{\Phi}_k(\rho, X)), \\ \mathbf{e}_k \cdot \bar{D}^{*,j}(\rho, X)\mathbf{e}_k &= \int_{\mathbb{Y}} \mathbf{e}_k \cdot \mathcal{D}(\rho, X)(\mathbf{e}_k + \nabla \bar{\Phi}_k^j(\rho, X)), \end{aligned}$$

if $M \in L^2(\mathbb{Y}, \mathcal{M}_d(\mathbb{R}))$ is square-integrable but not *essentially bounded*, we end up with

$$C_1 \bar{\mathcal{E}}^j(\rho, X, k) \leq \|\nabla \bar{\Phi}_k(\rho, X) - \nabla \bar{\Phi}_k^j(\rho, X)\|_{L^2(\mathbb{Y})}^{1/2},$$

for some $C_1 > 0$, a weaker estimate than the first inequality of (3.6). As a consequence, the convergence of the RB method and of the greedy algorithm in this case does not follow from ^{7,12,13,3}. Filling the gap in the analysis for such unbounded coefficients is beyond the scope of the present work. Nevertheless, the numerical experiments show the efficiency of the algorithm to treat this case.

3.3.2. Fast-assembly procedure

In this section, we restrict our discussion to $d = 2$ for notational convenience. The case $d > 2$ can be treated similarly. In dimension 2, the unit sphere S^1 is parametrized by $[0, 2\pi]$, so that from now on, we write the element of S^1 as

$$X = \mathbf{e}(\theta) = \cos(\theta)\mathbf{e}_1 + \sin(\theta)\mathbf{e}_2, \quad (3.14)$$

and consider \mathcal{D} as a function of ρ and θ (instead of ρ and X). The diffusion matrix $\mathcal{D} : [0, 1] \times [0, 2\pi] \rightarrow L^2(\mathbb{Y}, \mathcal{M}_d(\mathbb{R}))$ given by (3.10), that is

$$\mathcal{D}(\rho, \theta) : y \mapsto (1 - \rho)D_0 + \rho \left(\alpha |M(y)\mathbf{e}(\theta)| \mathbb{I} + \beta \frac{M(y)\mathbf{e}(\theta) \otimes M(y)\mathbf{e}(\theta)}{|M(y)\mathbf{e}(\theta)|} \right),$$

is affine with respect to ρ , but not with respect to $\theta \in [0, 2\pi]$. The empirical interpolation technique has been successfully developed to deal with such problems, see for instance ²⁰. It amounts to constructing iteratively and adaptively a basis and interpolation points (called magic points) using a greedy algorithm. Yet the efficiency of this method heavily rests on the regularity of the coefficients with respect to both the space variable and the parameter. In the case under investigation here, the coefficients are not smooth in space, not even continuous (the coefficients are piecewise constant). As a direct consequence, the number of magic points to be considered grows at least linearly with the number of elements where the coefficients are constant. This is not a desired scaling property since its cost would increase with mesh refinement. This is observed in practice, even on an elementary one-dimensional example.

To circumvent this difficulty we use a partial Fourier series expansion in the θ -variable, and write:

$$\mathcal{D}(\rho, \theta)(y) = (1 - \rho)D_0 + \rho \left(\frac{a_0(y)}{2} + \sum_{n=1}^{\infty} (a_n(y) \cos(n\theta) + b_n(y) \sin(n\theta)) \right),$$

where the functions $y \mapsto a_n(y)$ and $y \mapsto b_n(y)$ are matrix fields which depend only on $y \mapsto M(y)$.

Given a finite-dimensional space $\mathcal{V}_N = \text{span}\{\Psi_1, \dots, \Psi_N\}$ of dimension $N \geq 1$, and some parameters $(\rho, \theta) \in [0, 1] \times [0, 2\pi]$ and $k \in \{1, \dots, d\}$, in order to approximate the corrector $\bar{\Phi}_k$ in \mathcal{V}_N , it is enough to solve the linear system

$$\mathbb{M}(\rho, \theta)U = B(\rho, \theta, k),$$

where U is the vector of coordinates of $\bar{\Phi}_k$ in V_N , $\mathbb{M}(\rho, \theta)$ is the $N \times N$ -matrix given for all $1 \leq j_1, j_2 \leq N$ by

$$\begin{aligned} \mathbb{M}(\rho, \theta)_{j_1 j_2} &= (1 - \rho) \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot D_0 \nabla \Psi_{j_2} + \rho \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot \frac{a_0}{2} \nabla \Psi_{j_2} \\ &\quad + \sum_{n=1}^{\infty} \rho \cos(n\theta) \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot a_n \nabla \Psi_{j_2} + \sum_{n=1}^{\infty} \rho \sin(n\theta) \int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot b_n \nabla \Psi_{j_2}, \end{aligned}$$

and the r. h. s. is the N -vector given for all $1 \leq j \leq N$ by

$$\begin{aligned} B(\rho, \theta, k)_j &= -(1 - \rho) \int_{\mathbb{Y}} \nabla \Psi_j \cdot D_0 \mathbf{e}_k - \rho \int_{\mathbb{Y}} \nabla \Psi_j \cdot \frac{a_0}{2} \mathbf{e}_k \\ &\quad - \sum_{n=1}^{\infty} \rho \cos(n\theta) \int_{\mathbb{Y}} \nabla \Psi_j \cdot a_n \mathbf{e}_k - \sum_{n=1}^{\infty} \rho \sin(n\theta) \int_{\mathbb{Y}} \nabla \Psi_j \cdot b_n \mathbf{e}_k. \end{aligned}$$

In particular, provided we truncate the Fourier series expansion up to some order $L \in \mathbb{N}$, a fast assembly procedure can be devised if the $2(L+1)$ following matrices of order N and $2Lk(L+1)$ following vectors of order N are stored:

$$\begin{aligned} &\left(\int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot D_0 \nabla \Psi_{j_2} \right)_{j_1, j_2}, \quad \left(\int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot \frac{a_0}{2} \nabla \Psi_{j_2} \right)_{j_1, j_2}, \\ &\left(\int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot a_n \nabla \Psi_{j_2} \right)_{j_1, j_2}, \quad \left(\int_{\mathbb{Y}} \nabla \Psi_{j_1} \cdot b_n \nabla \Psi_{j_2} \right)_{j_1, j_2} \quad \text{for } n \in \{1, \dots, L\}, \end{aligned} \quad (3.15)$$

and for $k \in \{1, \dots, d\}$,

$$\begin{aligned} &\left(\int_{\mathbb{Y}} \nabla \Psi_j \cdot D_0 \mathbf{e}_k \right)_j, \quad \left(\int_{\mathbb{Y}} \nabla \Psi_j \cdot \frac{a_0}{2} \mathbf{e}_k \right)_j, \\ &\left(\int_{\mathbb{Y}} \nabla \Psi_j \cdot a_n \mathbf{e}_k \right)_j, \quad \left(\int_{\mathbb{Y}} \nabla \Psi_j \cdot b_n \mathbf{e}_k \right)_j \quad \text{for } n \in \{1, \dots, L\}. \end{aligned} \quad (3.16)$$

Note that the number of real numbers to be stored for the fast-assembly only depends on L and N . In particular, if the reduced basis vectors Ψ_j are approximated in a finite-dimensional subspace of $H_{\#}^1(\mathbb{Y})$, this number is independent of the size of that subspace, as desired.

In practice, once we are given the reduced basis $\{\Psi_1, \dots, \Psi_N\}$, the matrices (3.15) and vectors (3.16) can be obtained by performing a fast Fourier transform of

$$\theta \mapsto \alpha |M(y)\mathbf{e}(\theta)| \mathbb{I} + \beta \frac{M(y)\mathbf{e}(\theta) \otimes M(y)\mathbf{e}(\theta)}{|M(y)\mathbf{e}(\theta)|}$$

at each Gauss point $y \in \mathbb{Y}$ to evaluate the values of $a_n(y)$ and $b_n(y)$.

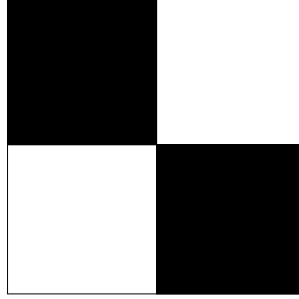


Fig. 2. Checkerboard structure

3.3.3. Numerical results

Let $d = 2$, $T_{\mathbb{Y}, h_1}, T_{\mathbb{Y}, \bar{h}_1}$ be regular tessellations of \mathbb{Y} of meshsize $h_1, \bar{h}_1 > 0$, and $\mathcal{V}_{\mathbb{Y}, h_1}^1, \mathcal{V}_{\mathbb{Y}, \bar{h}_1}^1$ be the subspaces of $H_{\#}^1(\mathbb{Y})$ made of \mathbb{P}_1 -periodic finite elements associated with $T_{\mathbb{Y}, h_1}$ and $T_{\mathbb{Y}, \bar{h}_1}$, respectively. The diffusion matrix $M \in L^2(\mathbb{Y}, \mathcal{M}_d(\mathbb{R}))$ is defined by

$$M(y) = K(y)(\mathbb{I} + \nabla \varphi^{\bar{h}_1}(y)),$$

where K is a standard checkerboard: for all $y = (y_1, y_2) \in \mathbb{Y}$,

$$K(y_1, y_2) = \begin{cases} 4.94064, & \text{if } \{y_1 \geq 0.5, y_2 \geq 0.5\} \text{ or } \{y_1 < 0.5, y_2 < 0.5\}, \\ 0.57816, & \text{elsewhere,} \end{cases}$$

see Figure 2, and $\varphi^{\bar{h}_1} = (\varphi_1^{\bar{h}_1}, \dots, \varphi_d^{\bar{h}_1})$ is defined as in (3.1). In this case, the correctors do not belong to finite element spaces, and shall take $h_1 \leq \bar{h}_1$. In the computations, we take $\bar{h}_1 \in \{1/10, 1/20, 1/40\}$ so that $\nu_{\bar{h}_1} := \dim \mathcal{V}_{\mathbb{Y}, \bar{h}_1}^1 \sim 100, 400, 1600$. The other parameters are the same as in Table 1. In the rest of this paragraph, we assume that the corrector equations (3.2) are solved in $\mathcal{V}_{\mathbb{Y}, h_1}^1$, so that the reduced basis will be a subspace of $\mathcal{V}_{\mathbb{Y}, h_1}^1$ as well.

For the reduced basis method we replace the compact space $\mathcal{P} = [0, 1] \times [0, 2\pi]$ by the finite set $\mathcal{P}_p := \{(\rho_i, \theta_j), (i, j) \in \{1, \dots, p\} \times \{1, \dots, p-1\}\}$, with $p \geq 2$, $\theta_j = (j-1)\frac{2\pi}{p-1}$, and $\rho_i = (i-1)\frac{1}{p-1}$, whose cardinal is denoted by \mathcal{N} . Let us denote by \mathcal{D}_L the diffusion matrix obtained by a truncation of the Fourier series expansion of \mathcal{D} at order L , and let \bar{D}^* denote the homogenized coefficients defined in (3.13) (where the correctors $\bar{\Phi}_k(\rho, X)$ is in fact approximated in $\mathcal{V}_{\mathbb{Y}, h_1}^1$, and with X related to θ through (3.14)), and let \bar{D}_L^* be defined by

$$\mathbf{e}_k \cdot \bar{D}_L^*(\rho, \theta) \mathbf{e}_k = \int_{\mathbb{Y}} (\mathbf{e}_k + \nabla \bar{\Phi}_k(\rho, \theta)) \cdot \mathcal{D}_L(\rho, \theta) (\mathbf{e}_k + \nabla \bar{\Phi}_k(\rho, \theta)).$$

We choose L such that

$$\sup_{i, j \in \{1, \dots, p\}} \frac{|\bar{D}^*(\rho_i, \theta_j) - \bar{D}_L^*(\rho_i, \theta_j)|}{|\bar{D}^*(\rho_i, \theta_j)|} \leq 10^{-6}.$$

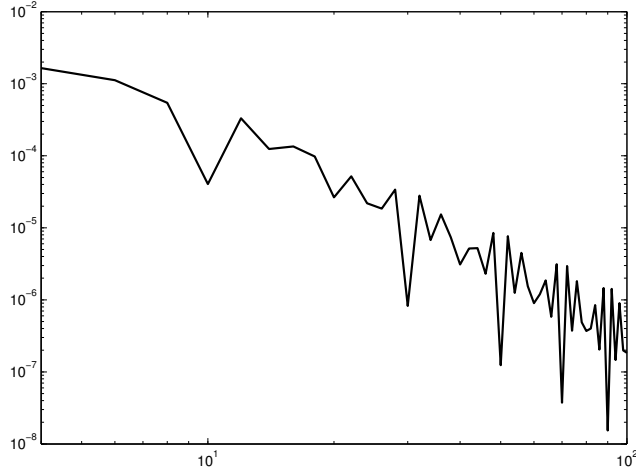


Fig. plot of the error due to the Fourier series expansion: $L \mapsto \sup_{i,j \in \{1, \dots, p\}} \frac{|\bar{D}^*(\rho_i, \theta_j) - \bar{D}_L^*(\rho_i, \theta_j)|}{|\bar{D}^*(\rho_i, \theta_j)|}$ 3. (slope of linear fitting: -3).

Note that in order to reduce the effect of the aliasing phenomena in the fast Fourier transform, we compute in practice twice as many coefficients as needed (that is, up to $2L$ for an effective truncation of order L). Numerical tests show that the convergence rate is 3, as can be seen on Figure 3, and that L depends both on the dimension $\nu_{\bar{h}_1}$ of $\mathcal{V}_{\bar{Y}, \bar{h}_1}^1$ and on the number \mathcal{N} of samples, but not on the dimension of $\mathcal{V}_{\bar{Y}, h_1}^1$ (associated with the discretization parameter h_1). As can be expected, the smaller \bar{h}_1 , the finer the approximation $\varphi^{\bar{h}_1}$ of the correctors φ of the Darcy equation, the more complex \mathcal{D} (it should however stabilize as $\bar{h}_1 \rightarrow 0$). We display the results of the numerical tests on L in Table 3.

For all $N \leq \mathcal{N} = p(p-1)$, we denote by \mathcal{V}_N the RB space of dimension N . We then choose N such that

$$\sup_{\mathcal{D}_p} (\mathfrak{E}_L^N(\rho, \theta))^2 \leq 10^{-6},$$

where \mathfrak{E}_L^N is the estimator associated with \mathcal{D}_L and the space \mathcal{V}_N , when the equations are solved in $\mathcal{V}_{\bar{Y}, h_1}^1$. As expected, N depends both on the dimension $\nu_{\bar{h}_1}$ of $\mathcal{V}_{\bar{Y}, \bar{h}_1}^1$ and on the number \mathcal{N} of samples, but not on the dimension of $\mathcal{V}_{\bar{Y}, h_1}^1$ (associated with the discretization parameter $h_1 \in \{1/10, 1/20, 1/40, 1/80, 1/160, 1/320\}$), which is the desired scaling property. The dimension N of the reduced basis in function of $\nu_{\bar{h}_1}$ and \mathcal{N} is displayed in Table 4. A more standard plot represents the RB error in function of the RB size. For completeness we have plotted such a graph on Figure 4, for $p = 10$, $L = 40$, $\bar{h}_1 = 1/20$, and $h_1 = 1/40$. As in simpler cases, the convergence

\mathcal{N} \backslash $\nu_{\bar{h}_1}$	100	400	1600
110	41	61	61
420	41	61	61
1640	49	95	175

Table 3. Dependence of the order L of the Fourier series expansion for an error less than 10^{-6} in function of the dimension $\nu_{\bar{h}_1}$ of $\mathcal{V}_{\bar{Y}, \bar{h}_1}^1$ and of the cardinal \mathcal{N} of \mathcal{P} .

\mathcal{N} \backslash $\nu_{\bar{h}_1}$	100	400	1600
110	21	25	25
420	23	38	44
1640	24	47	60

Table 4. Dependence of the size N of the reduced basis for an error less than 10^{-6} in function of the dimension $\nu_{\bar{h}_1}$ of $\mathcal{V}_{\bar{Y}, \bar{h}_1}^1$ and of the cardinal \mathcal{N} of \mathcal{P} .

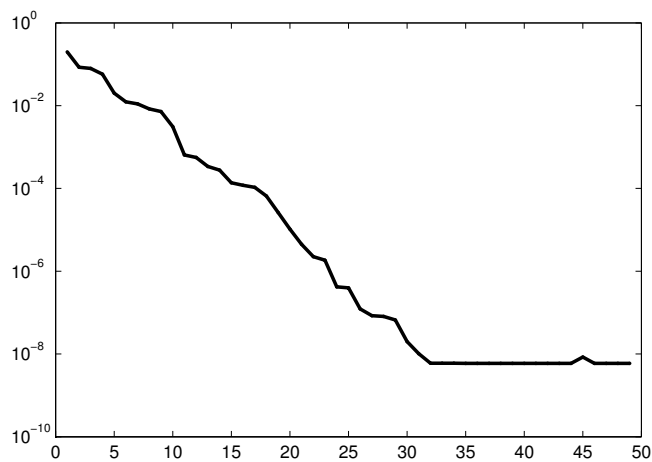


Fig. 4. RB basis error $\sup_{\mathcal{P}_p} (\mathfrak{e}_L^N(\rho, \theta))^2$ for $\mathcal{N} = 110$ in function of the size N of the reduced basis, exponential convergence.

is exponential (10^{-9} is the machine precision).

In order to check a posteriori the efficiency of the method (both in terms of L and

\mathcal{N} \backslash $\nu_{\bar{h}_1}$	100	400	1600
110	1.2e-04	9.0e-04	1.6e-03
420	3.4e-05	1.8e-04	2.6e-04
1640	7.4e-06	3.5e-05	2.0e-04

Table 5. Dependence of the RB error u in function of the dimension $\nu_{\bar{h}_1}$ of $\mathcal{V}_{\bar{y}, \bar{h}_1}^1$ and of the cardinal \mathcal{N} of \mathcal{P} , on a random sampling of 100 points

N), we have picked at random a set $\tilde{\mathcal{P}}$ of 100 pairs of parameters $(\rho, \theta) \in [0, 1] \times [0, 2\pi]$, computed the corresponding approximations $\bar{D}^*(\rho, \theta)$ of the homogenized coefficients in $\mathcal{V}_{\bar{y}, h_1}^1$, and compared them to the approximations $\bar{D}_L^{*,N}(\rho, \theta)$ using the reduced basis method of order N and a Fourier series expansion of \mathcal{D} truncated at order L . The numerical tests show that this error

$$\sup_{\tilde{\mathcal{P}}} \frac{|\bar{D}^*(\rho, \theta) - \bar{D}_L^{*,N}(\rho, \theta)|}{|\bar{D}^*(\rho, \theta)|}$$

does depend on the dimension $\nu_{\bar{h}_1}$ of $\mathcal{V}_{\bar{y}, \bar{h}_1}^1$ and on the number \mathcal{N} of samples, but not on the dimension of $\mathcal{V}_{\bar{y}, h_1}^1$ (associated with the discretization parameter h_1). We have chosen $p \in \{11, 21, 41\}$ so that the sample sets are included in one another, which ensures that the error due to the RB method decreases as p (and $\mathcal{N} = p(p-1)$) increases, as can be checked on Table 5. Note also that the error increases as $\bar{h}_1 \rightarrow 0$ (that is $\nu_{\bar{h}_1} \rightarrow \infty$).

A last comment is in order. For $\mathcal{N} = 1640$ and $\nu_{\bar{h}_1} = 1600$, the error is not reduced much with respect to $\mathcal{N} = 420$ in Table 5. On Figure 5 the points chosen by the greedy algorithm are plotted for $\mathcal{N} = 1640$ and $\nu_{\bar{h}_1} = 1600$ (circles denote points for the corrector in the direction \mathbf{e}_1 and crosses denote points for the corrector in the direction \mathbf{e}_2). This figure shows that most of the information for the RB lies in the region ρ close to 1 and θ in $[0, \pi]$ (this latter fact is indeed a consequence of the identity $\mathcal{D}(\rho, \theta) = \mathcal{D}(\rho, \pi - \theta)$). This motivates us to put more points in this region rather than in the rest of \mathcal{P} , and allows us to focus on the right region of the parameters. Taking for instance 5×168 points in the region $[0.9, 1] \times [0, \pi]$ and 10×30 in $[0, 1] \times [0, \pi]$, that is a total of 1140 points (to be compared to the 1640 uniformly chosen points in \mathcal{P}), the reduced basis has dimension $N = 68$ for $\nu_{\bar{h}_1} = 1600$, $L = 177$, and the error on the 100 random points of $\tilde{\mathcal{P}}$ is reduced to $4.1e - 05$ (instead of $2.0e - 04$).

For completeness we have also tested the influence of the choice of the scalar product (3.7) to orthogonalize the reduced basis vectors, compared to the canonical $L^2(\Omega)$ -scalar product. In practice this choice can only influence the conditioning of the linear system to be solved to approximate solutions in the reduced basis.

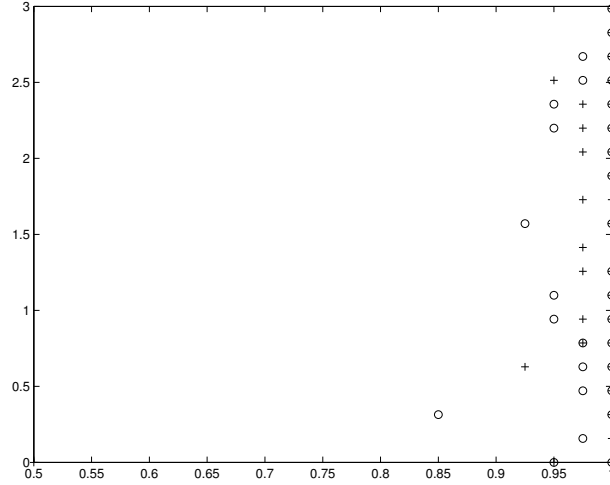


Fig. 5. Points chosen by the greedy algorithm for $\mathcal{N} = 1640$ and $\nu_{\bar{n}_1} = 1600$ (all the points chosen in $[0, 1] \times [0, 2\pi]$ lie in $[0.5, 1] \times [0, \pi]$).

N	$L^2(\Omega)$ scalar product (min. / max.)	scalar product (3.7) (min. / max.)
21	2.13 / 14.06	4.38 / 9.72
25	2.48 / 35.18	6.08 / 22.69
38	6.96 / 39.50	14.69 / 31.52
44	7.11 / 73.38	16.99 / 51.61

Table 6. Condition numbers (min and max) in function of the reduced basis dimension N , depending on the scalar product used.

We have compared the condition number of those matrices constructed with the scalar product (3.7) to the condition number of those matrices constructed with the canonical scalar product of $L^2(\Omega)$ on the 100 random pairs of parameters $(\rho, \theta) \in [0, 1] \times [0, 2\pi]$, as a function of the dimension of the reduced basis. The minimum and the maximum of the condition numbers are reported on Table 6. As can be seen the condition numbers are of the same orders, although they seem to depend less on the parameters for the scalar product (3.7).

In conclusion, these tests widely confirm the efficiency of the method. We do not observe any difficulty that could be due to the specific form of the coefficients. In particular the convergence properties that we observe seem not to be altered by the fact that the coefficients are unbounded. Note also that the numerical difficulty of the computation of the effective coefficients in itself would not be simplified by

considering bounded coefficients.

Remark 2. As in Remark 1, we could also consider locally periodic coefficients depending on both the slow and the fast variables, provided the dependence with respect to the slow variable is smooth enough. Of course the price to be paid is to increase the size of the set of parameters \mathcal{P} accordingly.

4. Conclusion

We have considered a simple model of radionuclide transport in porous media: the radionuclide concentration satisfies a convection–diffusion equation where the coefficients are determined through the Darcy law by solving a non-homogeneous elliptic equation. A remarkable feature of the model relies on the fact that the diffusion coefficients depend non linearly on the velocity and do not satisfy a uniform L^∞ -estimate.

Nevertheless, existence-uniqueness results can be established for this model. Although we can also perform the homogenization analysis, the effective coefficients remain non homogeneous due to the nonlinear coupling, even in the simple case of periodic oscillations. It impacts strongly the computational cost when using direct evaluations of the coefficients. We propose an algorithm based on the Reduced Basis method in order to speed up these computations. The method relies on a suitable parametrization of the problem, which in particular allows us to make use of the Fast Fourier Transform to construct efficiently stiffness matrices. Working with unbounded coefficients is clearly identified as a difficulty for analyzing the convergence properties of the method, but simulations demonstrate the efficiency of the scheme which is a valuable tool for the computation of such complex flows.

Appendix A. Proof of Theorem 2

We decompose the proof into two steps and homogenize the Darcy equation and the advection-diffusion equation separately.

Step 1. Homogenization of the Darcy equation, and two-scale convergence of U_ε and $D(U_\varepsilon)$.

By standard two-scale convergence arguments (see ¹, and see also ²²), the function Θ_ε two-scale converges to Θ_0 and $\nabla\Theta_\varepsilon$ two-scale converges to the function $(x, y) \in \Omega \times \mathbb{Y} \mapsto (\mathbb{I} + \nabla\varphi(y))\nabla\Theta_0(x)$. Likewise the flux $K_\varepsilon\nabla\Theta_\varepsilon$ two-scale converges to $(x, y) \mapsto K(y)(\mathbb{I} + \nabla\varphi(y))\nabla\Theta_0(x) = -\tilde{U}(x, y)$. In order to homogenize the advection-diffusion equation, we need the function $K_\varepsilon\nabla\Theta_\varepsilon$ to be an admissible test-function for two-scale convergence, see ¹. It is enough to prove that $K_\varepsilon\nabla\Theta_\varepsilon$ strongly two-scale converges to $-\tilde{U}$, that is, in addition of two-scale convergence, to prove that we have

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} |K_\varepsilon(x)\nabla\Theta_\varepsilon(x)|^2 dx = \int_{\Omega} \int_{\mathbb{Y}} |\tilde{U}(x, y)|^2 dy dx. \quad (\text{A.1})$$

This is essentially a consequence of the following convergence of the energy:

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \int_{\Omega} \nabla \Theta_{\varepsilon}(x) \cdot K_{\varepsilon}(x) \nabla \Theta_{\varepsilon} dx \\ = \int_{\Omega} \int_{\mathbb{Y}} (\nabla \Theta_0(x) + \nabla_y \Theta_1(x, y)) \cdot K(y) (\nabla \Theta_0(x) + \nabla_y \Theta_1(x, y)) dy dx, \end{aligned}$$

where $\Theta_1(x, y) = \sum_{i=1}^d \nabla_i \Theta_0(x) \varphi_i(y)$. In particular, since K is positive-definite, we may rewrite this identity as

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} |K_{\varepsilon}(x)^{1/2} \nabla \Theta_{\varepsilon}(x)|^2 dx = \int_{\Omega} \int_{\mathbb{Y}} |K(y)^{1/2} (\nabla \Theta_0(x) + \nabla_y \Theta_1(x, y))|^2 dy dx,$$

which upgrades the two-scale convergence of $K_{\varepsilon}^{1/2} \nabla \Theta_{\varepsilon}$ to $(x, y) \mapsto K(y)^{1/2} (\nabla \Theta_0(x) + \nabla_y \Theta_1(x, y))$ into strong two-scale convergence. We now consider a sequence $v_{\varepsilon} : \Omega \rightarrow \mathbb{R}^d$ which two-scale converges to $(x, y) \mapsto v_0(x, y)$. The sequence $V_{\varepsilon} := K_{\varepsilon}^{1/2} v_{\varepsilon}$ then two-scale converges to $(x, y) \mapsto K(y)^{1/2} v_0(x, y)$. Since we have proved that $K_{\varepsilon}^{1/2} \nabla \Theta_{\varepsilon}$ is an admissible test-function for the two-scale convergence, we have

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \int_{\Omega} K_{\varepsilon}(x)^{1/2} \nabla \Theta_{\varepsilon}(x) \cdot V_{\varepsilon}(x) dx \\ = \int_{\Omega} \int_{\mathbb{Y}} K(y)^{1/2} (\nabla \Theta_0(x) + \nabla_y \Theta_1(x, y)) \cdot K(y)^{1/2} v_0(x, y) dy dx. \end{aligned}$$

Taking $v_{\varepsilon} = K_{\varepsilon} \nabla \Theta_{\varepsilon}$ then proves (A.1).

We conclude this step by the proof of the strong two-scale convergence of $D(U_{\varepsilon})$ to $(x, y) \mapsto \tilde{D}(x, y)$. This is a direct consequence of ¹⁸ since $(x, U) \mapsto D(U)(x)$ is a Lipschitz function with respect to U uniformly in x , and U_{ε} strongly two-scale converges to \tilde{U} .

Step 2. Homogenization of the advection-diffusion equation.

In view of the results of Step 1, this is now standard matter to prove the two-scale convergence of C_{ε} to C_0 . The proof of Theorem 1 provides uniform bounds on C_{ε} which gives weak compactness. Hence, up to extraction, C_{ε} two-scale converges to some C_0 , and ∇C_{ε} to some $(x, y) \mapsto \nabla C_0(x) + \nabla_y C_1(x, y)$ (time is treated as a parameter). Since $D(U_{\varepsilon})$ and U_{ε} strongly two-scale converge to \tilde{D} and \tilde{U} , respectively, there is no difficulty to pass at the two-scale limit in the equation for C_{ε} tested with functions $(t, x) \mapsto \psi(t) \phi(x, x/\varepsilon)$ and $\phi \in C^{\infty}(\Omega, C_{\text{per}}^{\infty}(\mathbb{Y}))$ and $\psi \in C^{\infty}(0, T)$.

It remains to note that following the arguments of Step 5 in the proof of Theorem 1, we obtain

$$\begin{aligned} \int_0^T \int_{\Omega} \nabla C_0 \cdot D^* \nabla C_0 = \int_{\Omega} \int_{\mathbb{Y}} (\nabla C_0(x) + \nabla_y C_1(x, y)) \cdot \tilde{D}(x, y) (\nabla C_0(x) + \nabla_y C_1(x, y)) dy dx \\ \leq \|C_{\text{init}}\|_{L^2(\Omega)}^2 + \|S\|_{L^2(0, T; H^{-1}(\Omega))}^2, \end{aligned}$$

so that (Θ_0, C_0) is the unique weak solution of the homogenized system, and the whole sequence C_ε converges.

Acknowledgements

This work is supported by ANDRA (Direction Recherche et Développement/Service Evaluation et Analyse de Performances) through a specific Inria–ANDRA research partnership. We thank Guillaume Pépin and Marc Leconte for having introduced us to the problem and for their kind and constant encouragements.

References

1. G. Allaire. Homogenization and two-scale convergence. *SIAM J. Math. Anal.*, 23:1482–1518, 1992.
2. M. Bellieud and G. Bouchitté. Homogenization of elliptic problems in a fiber reinforced structure. Nonlocal effects. *Ann. Scuola Norm. Sup. Pisa Cl. Sci. (4)*, 26(3):407–436, 1998.
3. P. Binev, A. Cohen, W. Dahmen, R. DeVore, G. Petrova, and P. Wojtaszczyk. Convergence rates for greedy algorithms in reduced basis methods. 2010. Preprint AICES–2010/05-2.
4. A. Bourgeat, M. Kern, S. Schumacher, and J. Talandier. The complex test cases: Nuclear waste disposal simulation. *Computational Geosciences*, 8:83–98, 2004.
5. S. Boyaval. Reduced-basis approach for homogenization beyond the periodic setting. *Multiscale Model. Simul.*, 7(1):466–494, 2008.
6. M. Briane. Homogenization of non-uniformly bounded operators: critical barrier for nonlocal effects. *Arch. Ration. Mech. Anal.*, 164(1):73–101, 2002.
7. A. Buffa, Y. Maday, A. T. Patera, C. Prud'homme, and G. Turinici. A priori convergence of the greedy algorithm for the parametrized reduced basis. Preprint.
8. C. Chainais-Hillairet and J. Droniou. Convergence analysis of a mixed finite volume scheme for an elliptic-parabolic system modeling miscible fluid flows in porous media. *SIAM J. Numer. Anal.*, 45(5):2228–2258 (electronic), 2007.
9. C. Chainais-Hillairet, S. Krell, and A. Mouton. Convergence analysis and numerical results of a finite volume discretization for the peaceman model. In preparation.
10. Z. Chen and R. Ewing. Mathematical analysis for reservoir models. *SIAM J. Math. Anal.*, 30(2):431–453, 1999.
11. C. Choquet and A. Sili. Homogenization of a model of displacement with unbounded viscosity. *Netw. Heterog. Media*, 4(4):649–666, 2009.
12. A. Cohen, R. DeVore, and C. Schwab. Convergence rates of best N -term Galerkin approximations for a class of elliptic PDEs. *Found. Comput. Math.*, 10(6):615–646, 2010.
13. A. Cohen, R. DeVore, and C. Schwab. Analytic regularity and polynomial approximation of parametric and stochastic elliptic PDE's. *Anal. Appl. (Singap.)*, 9(1):11–47, 2011.
14. J. Douglas, R. E. Ewing, and M. F. Wheeler. A time-discretization procedure for a mixed finite element approximation of miscible displacement in porous media. *RAIRO Anal. Numér.*, 17(3):249–265, 1983.
15. R. E. Ewing, T. F. Russell, and M. F. Wheeler. Simulation of miscible displacement using mixed methods and a modified method of characteristics. *Proceedings of the 7th SPE Symposium on Reservoir Simulation, Dallas, TX, Paper SPE 12241, Society of Petroleum Engineers*, pages 71–81, 1983.

16. P. Fabrie and M. Langlais. Mathematical analysis of miscible displacement in porous medium. *SIAM J. Math. Anal.*, 23(6):1375–1392, 1992.
17. X. Feng. On existence and uniqueness results for a coupled system modeling miscible displacement in porous media. *J. Math. Anal. Appl.*, 194(3):883–910, 1995.
18. T. Goudon and F. Poupaud. Approximation by homogenization and diffusion of kinetic equations. *Comm. Partial Differential Equations*, 26(3-4):537–569, 2001.
19. F. Hecht, A. Le Hyaric, K. Ohtsuka, and O. Pironneau. Freefem++, finite elements software. <http://www.freefem.org/ff++/>.
20. Y. Maday, N. C. Nguyen, Anthony T. Patera, and G. S. H. Pau. A general multipurpose interpolation procedure: the magic points. *Commun. Pure Appl. Anal.*, 8(1):383–404, 2009.
21. Y. Maday, A. T. Patera, and G. Turinici. Global a priori convergence theory for reduced-basis approximations of single-parameter symmetric coercive elliptic partial differential equations. *C. R. Math. Acad. Sci. Paris*, 335(3):289–294, 2002.
22. G. Nguetseng. A general convergence result for a functional related to the theory of homogenization. *SIAM J. Math. Anal.*, 20:608–623, 1989.
23. A. Quarteroni and A. Valli. *Numerical Approximation of Partial Differential equations*, volume 23 of *Springer Series in Computational Mathematics*. Springer-Verlag, 1997.
24. T. F. Russell. Finite elements with characteristics for two-component incompressible miscible displacement. *Proceedings of the 6th SPE Symposium on Reservoir Simulation, New Orleans, TX, Paper SPE 10500, Society of Petroleum Engineers*, pages 123–135, 1982.
25. J. Simon. Compact sets in the space $L^p(0, T; B)$. *Ann. Mat. Pura Appl. (4)*, 146:65–96, 1987.
26. H. Wang, D. Liang, R. E. Ewing, S. L. Lyons, and G. Qin. An approximation to miscible fluid flows in porous media with point sources and sinks by an Eulerian-Lagrangian localized adjoint method and mixed finite element methods. *SIAM J. Sci. Comput.*, 22(2):561–581 (electronic), 2000.
27. H. Wang, D. Liang, R. E. Ewing, S. L. Lyons, and G. Qin. An improved numerical simulator for different types of flows in porous media. *Numer. Methods Partial Differential Equations*, 19(3):343–362, 2003.