



**HAL**  
open science

# Linearized numerical homogenization method for nonlinear monotone parabolic multiscale problems

Assyr Abdulle, Martin Huber, Gilles Vilmart

► **To cite this version:**

Assyr Abdulle, Martin Huber, Gilles Vilmart. Linearized numerical homogenization method for nonlinear monotone parabolic multiscale problems. 2014. hal-01017106v1

**HAL Id: hal-01017106**

**<https://hal.science/hal-01017106v1>**

Preprint submitted on 1 Jul 2014 (v1), last revised 4 Jun 2015 (v3)

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

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

# Linearized numerical homogenization method for nonlinear monotone parabolic multiscale problems

A. Abdulle<sup>1</sup>, M.E. Huber<sup>1</sup>, and G. Vilmart<sup>2</sup>

July 1, 2014

## Abstract

We introduce and analyze an efficient numerical homogenization method for a class of nonlinear parabolic problems of monotone type in highly oscillatory media. The new scheme avoids costly Newton iterations and is linear at both the macroscopic and the microscopic scales. It can be interpreted as a linearized version of a standard nonlinear homogenization method. We prove the stability of the method and derive optimal a priori error estimates which are fully discrete in time and space. Numerical experiments confirm the error bounds and illustrate the efficiency of the method for various nonlinear problems.

*Keywords:* monotone parabolic multiscale problem, linearized scheme, numerical homogenization method, fully discrete a priori error estimates.

*AMS subject classification (2010):* 65M60, 74Q10, 74D10.

## 1 Introduction

In this paper, we propose a linearized numerical homogenization method for the efficient approximation of multiscale parabolic problems of the form

$$\partial_t u^\varepsilon - \operatorname{div}(a^\varepsilon(x, \nabla u^\varepsilon) \nabla u^\varepsilon) = f, \quad \text{in } \Omega \times (0, T), \quad (1)$$

where  $\Omega \subset \mathbb{R}^d$  (with  $d \leq 3$ ) is a polygonal domain,  $T > 0$  is the final time and the  $d \times d$  tensor  $a^\varepsilon(x, \xi)$  rapidly fluctuates in space at a small scale  $\varepsilon$  and is both elliptic and bounded uniformly with respect to  $\varepsilon$ . For simplicity, we impose for (1) homogeneous Dirichlet boundary conditions and a non oscillatory initial condition at  $t = 0$ . We assume throughout this article that the equation (1) is of monotone type to guarantee the well-posedness of the problem, i.e., the maps  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi)\xi$  are strongly monotone and Lipschitz continuous with respect to  $\xi \in \mathbb{R}^d$ , with constants independent of  $\varepsilon$ . This class of problems arises in many applications, e.g., material laws in elasticity or constitutive relations in magnetodynamics [34, 35]. The nonlinearity of the tensor  $a^\varepsilon$  in (1) with respect to the solution gradient  $\nabla u^\varepsilon$  makes the problem challenging both computationally and for the analysis because it combines the difficulties of the finescale structure of the data (with small variations at the microscopic scale) and the nonlinearity of the problem.

The aim of this paper is to analyze the convergence of a linearized version of a nonlinear homogenization method proposed in [6] to approximate the effective solution to the multiscale problem  $\partial_t u^\varepsilon - \operatorname{div}(\mathcal{A}^\varepsilon(x, \nabla u^\varepsilon)) = f$  which includes the class of problems (1) for  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi)\xi$ . The method of [6] combines a nonlinear FE-HMM method (coupling macro and micro finite element methods) with the implicit Euler method in time. Although the computational cost of the nonlinear method in [6] is independent of the small scale  $\varepsilon$ , its upscaling procedure however relies on nonlinear elliptic cell problems which is computationally costly for practical simulations. The linearized version presented in this paper permits to avoid Newton iterations, which considerably improves the computational efficiency of the method.

The existence of a macroscopic effective model for (1) in the asymptotic regime  $\varepsilon \rightarrow 0$  is ensured by the homogenization theory in [36, 37]. Since the effective material properties are not explicitly available in

---

<sup>1</sup>ANMC, Mathematics Section, École Polytechnique Fédérale de Lausanne, Station 8, CH-1015 Lausanne, Switzerland; Assyr.Abdulle@epfl.ch, Martin.Huber@epfl.ch.

<sup>2</sup>Université de Genève, Section de mathématiques, 2-4 rue du Lièvre, CP 64, 1211 Genève 4, Switzerland. On leave from École Normale Supérieure de Rennes, INRIA Rennes, IRMAR, CNRS, UEB, av. Robert Schuman, F-35170 Bruz, France; Gilles.Vilmart@unige.fr

general, numerical methods are needed to estimate them. Following the methodology of the heterogeneous multiscale method (HMM) introduced in [18], we apply a finite element method (macro solver) on a spatial macro partition to approximate (in space) the solution of an effective equation whose material properties are recovered “on the fly“ by an appropriate numerical upscaling procedure. This is achieved by performing local micro simulations (using a finite element method as micro solver) within sampling domains which are of size comparable to  $\varepsilon$  (the size of spatial micro oscillations). Thus, the computational cost of the method is *independent* of the small scale  $\varepsilon$ . Combining the spatial macro solver with a new linearized implicit Euler scheme for the time integration leads to *linear* micro simulations. In turn, the resulting multiscale scheme does not involve *any* nonlinear equation to be solved, neither at the macro scale nor on the micro scale.

In the main results of this paper, we establish the stability of the method and derive optimal fully discrete a priori error estimates which hold without structural assumptions (like periodicity) on the tensor  $a^\varepsilon$ . The a priori estimates consist of explicit convergence rates for the time discretization error as well as the spatial finite element error on both macro and micro scale. Further, we derive error bounds that account for the modeling error depending on the parameters of the upscaling strategy, i.e., boundary conditions for micro simulations and size of the micro sampling domains.

**Literature overview.** Numerical homogenization methods are well developed for wide classes of multiscale problems, see e.g. the review [5] and references therein. However, the numerical literature for monotone parabolic problems (1) is less abundant. We refer to [38] for a multiscale method, which requires periodicity of the tensor  $a^\varepsilon$ , and to [21] for an extension of the multiscale finite element method (MsFEM). Further, we mention the discussions about linearization techniques for nonlinear monotone multiscale problems given in [20, 26]. As mentioned above, a nonlinear FE-HMM combined with the Euler implicit method has been proposed in [6] from where various results will be used in the present analysis.

For nonlinear parabolic singlescale problems  $\partial_t u - \operatorname{div} \mathcal{A}(\nabla u) = f$ , semi-discrete a priori error estimates (in space) have been derived in [13, 15]. Further, many strategies to linearize numerical methods for stiff ordinary differential equations (ODEs) and time-dependent singlescale PDEs are available in the literature. In particular, we mention Rosenbrock methods (only one Newton iteration per timestep) and  $W$ -methods (only one Newton iteration with inexact Jacobian per timestep), see [25, Section IV.7] for an overview. We emphasize that simply applying such linearized time integrators to the effective model associated to (1) does *not* yield the linearized multiscale scheme proposed in this article. Indeed, due to the nonlinearities arising at the microscopic level, the resulting scheme would remain nonlinear. For parabolic singlescale PDEs, already in the work [16] by Douglas and Dupont an extrapolated Crank-Nicholson time stepping scheme has been considered to avoid large nonlinear algebraic systems. Then, Nie and Thomée proposed in [33] linearized numerical schemes for singlescale problems of the form  $\partial_t u - \operatorname{div}(a(x, u)\nabla u) = f$  (in space dimension two) where  $a(x, s)$ , for  $s \in \mathbb{R}$ , is a strictly positive scalar function. Their results consist of optimal space-time a priori error estimates for numerical methods constructed by coupling a finite element method (with numerical quadrature) and linearized time integrators. Further, Makridakis studied in [32] a class of linearized space-time discrete methods (one linear system to solve per timestep) for a system of nonlinear hyperbolic PDEs from elastodynamics and derived optimal a priori error estimates (in both time and space). Finally, in [31], Lubich and Ostermann presented a semi-discrete analysis (in time) of linearly implicit integrators used for the time discretization of nonlinear parabolic PDEs, seen as evolution problems posed in Hilbert spaces.

**Outline.** The article is organized as follows. In Section 2, we recall the homogenization results for the model problem (1) and discuss conditions for the tensor  $a^\varepsilon$  that are sufficient to ensure the monotonicity of the model problem. Next, we introduce the linearized multiscale scheme in Section 3, prove the well-posedness of the numerical method and present the fully discrete a priori error estimates. The proofs of the error bounds are provided in Section 5. Further, in Section 6, we present several numerical experiments to illustrate the convergence results, the efficiency as well as the robustness of the method. The article ends with a conclusion in Section 7.

**Notations.** Let  $W^{k,p}(\Omega)$  denote the usual Sobolev spaces which we write as  $H^k(\Omega)$  for  $p = 2$ . Further, we use  $H_0^1(\Omega)$  for the space of  $H^1(\Omega)$ -functions with zero trace on the boundary  $\partial\Omega$ ,  $H^{-1}(\Omega)$  for its dual space and  $W_{per}^1(Y) = \{v \in H_{per}^1(Y) \mid \int_Y v \, dy = 0\}$  for periodic  $H^1(Y)$ -functions with zero mean on the unit cube  $Y = (0, 1)^d$  (with  $H_{per}^1(Y)$  being the closure of  $C_{per}^\infty(Y)$  for the  $H^1(Y)$  norm). Let  $g: [0, T] \rightarrow X$  be a function with values in a Banach space  $X$  (with norm  $\|\cdot\|_X$ ). The space of  $L^p$  functions

and continuous functions  $g$  with values in  $X$  is denoted by  $L^p(0, T; X)$  and  $C^0([0, T], X)$ , respectively. Both spaces form a Banach space when endowed with the norm  $\|g\|_{L^p(0, T; X)} = (\int_0^T \|g(t)\|_X^p dt)^{1/p}$  and  $\|g\|_{C^0([0, T], X)} = \sup_{t \in [0, T]} \|g(t)\|_X$ , respectively. For vectors  $b \in \mathbb{R}^d$  the Euclidean norm is denoted by  $|b|$  and the canonical basis of  $\mathbb{R}^d$  is represented by  $e_1, \dots, e_d$ . Further, we denote  $\|a\|_{\mathcal{F}}$  the Frobenius norm of matrices  $a \in \mathbb{R}^{d \times d}$ . Finally, the constant  $C$  is a generic constant whose value may differ at each occurrence. All constants considered in this paper are assumed independent of  $\varepsilon$ .

## 2 Model problem and homogenization

Let  $\Omega \times (0, T)$  be a space-time domain where  $\Omega \subset \mathbb{R}^d$  (with  $d \leq 3$ ) is a convex polygonal domain and  $T > 0$  is the final time. We study the parabolic quasilinear homogenization problem

$$\begin{aligned} \partial_t u^\varepsilon(x, t) - \operatorname{div}(a^\varepsilon(x, \xi) \nabla u^\varepsilon(x, t)) \nabla u^\varepsilon(x, t) &= f(x), & \text{in } \Omega \times (0, T), \\ u^\varepsilon(x, t) &= 0, & \text{on } \partial\Omega \times (0, T), \\ u^\varepsilon(x, 0) &= g(x), & \text{in } \Omega, \end{aligned} \quad (2)$$

where  $f \in L^2(\Omega)$  models the source term and  $g \in L^2(\Omega)$  prescribes the initial conditions.

**Assumptions on the tensor.** We assume that the family of tensors  $a^\varepsilon(x, \xi) \in (L^\infty(\Omega \times \mathbb{R}^d))^{d \times d}$  (indexed by  $\varepsilon$ ) is uniformly elliptic and bounded, i.e., there exist  $0 < \lambda_a \leq \Lambda_a$  such that

$$\lambda_a |\eta|^2 \leq a^\varepsilon(x, \xi) \eta \cdot \eta, \quad |a^\varepsilon(x, \xi) \eta| \leq \Lambda_a |\eta|, \quad \forall \xi, \eta \in \mathbb{R}^d, \text{ a.e. } x \in \Omega, \varepsilon > 0. \quad (3)$$

The parameter  $\varepsilon > 0$  denotes the characteristic length of the smallest scale in the problem (2). In particular, the tensors  $a^\varepsilon$  vary rapidly in space at this microscopic scale  $\varepsilon$ . We use homogeneous Dirichlet boundary conditions in (2) for simplicity. However, our analysis could be extended to other type of boundary conditions (as Neumann or mixed boundary conditions).

While the uniform ellipticity and boundedness in (3) are sufficient for the well-posedness of our algorithm, they are not sufficient in general to ensure the well-posedness of the exact problem (1). We therefore make the following standard hypotheses of strong monotonicity and Lipschitz continuity on the maps  $\mathcal{A}^\varepsilon: \Omega \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  defined by  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi) \xi$  for  $(x, \xi) \in \Omega \times \mathbb{R}^d$  and  $\varepsilon > 0$ . We assume that there exist  $L \geq \lambda > 0$  (independent of  $\varepsilon$ ) such that

$$(\mathcal{A}_1) \quad \text{Lipschitz continuity: } |\mathcal{A}^\varepsilon(x, \xi_1) - \mathcal{A}^\varepsilon(x, \xi_2)| \leq L |\xi_1 - \xi_2|, \text{ for } \xi_1, \xi_2 \in \mathbb{R}^d, \text{ a.e. } x \in \Omega;$$

$$(\mathcal{A}_2) \quad \text{Strong monotonicity: } [\mathcal{A}^\varepsilon(x, \xi_1) - \mathcal{A}^\varepsilon(x, \xi_2)] \cdot (\xi_1 - \xi_2) \geq \lambda |\xi_1 - \xi_2|^2, \text{ for } \xi_1, \xi_2 \in \mathbb{R}^d, \text{ a.e. } x \in \Omega.$$

Under the above assumptions, the well-posedness of monotone parabolic problems of the type (2) is classical, see [40, Theorem 30.A]: there exists a unique solution  $u^\varepsilon \in E$  in the Banach space  $E$  with norm  $\|u\|_E$  bounded independently of  $\varepsilon$ , where

$$E = \{v \in L^2(0, T; H_0^1(\Omega)) \mid \partial_t v \in L^2(0, T; H^{-1}(\Omega))\}, \quad \|v\|_E = \|v\|_{L^2(0, T; H_0^1(\Omega))} + \|\partial_t v\|_{L^2(0, T; H^{-1}(\Omega))}.$$

We note that the assumptions  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$  can be deduced from properties of the tensor  $a^\varepsilon$  itself as shown in the following remark. A proof is provided for completeness in Appendix.

**Remark 2.1.** Assume that  $a^\varepsilon$  is uniformly elliptic and bounded (3). If  $a^\varepsilon(x, \cdot) \in (W^{1, \infty}(\mathbb{R}^d))^{d \times d}$  for a.e.  $x \in \Omega$  and the following estimate from Babuška [9, Assumption 3.3 and 3.4] holds

$$\left( \sum_{i, j, k=1}^d \left| \frac{\partial a_{ij}^\varepsilon(x, \xi)}{\partial \xi_k} \right|^2 \right)^{1/2} \leq \frac{L_a}{1 + |\xi|}, \quad \forall \xi \in \mathbb{R}^d, \text{ a.e. } x \in \Omega, \varepsilon > 0, \quad (4)$$

with  $L_a < \lambda_a$ , where  $\lambda_a$  is the ellipticity constant from (3), then the maps  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi) \xi$  defined on  $\Omega \times \mathbb{R}^d$  are Lipschitz continuous and strongly monotone uniformly in  $\varepsilon > 0$ , i.e., satisfying  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$ .

For example, in the context of one-scale monotone elliptic problems, see [23, 28], one considers tensors of the form

$$a^\varepsilon(x, \xi) = \mu^\varepsilon(x, |\xi|) Id, \quad \text{with } \mu^\varepsilon: \bar{\Omega} \times [0, \infty) \rightarrow \mathbb{R} \text{ and } (x, \xi) \in \Omega \times \mathbb{R}^d,$$

where  $Id \in \mathbb{R}^{d \times d}$  is the identity matrix,  $\mu^\varepsilon$  is a continuous function on  $\bar{\Omega} \times [0, \infty)$  and  $\mu^\varepsilon(x, \cdot)$  is continuously differentiable for a.e.  $x \in \bar{\Omega}$ . If there exist  $M_\mu \geq m_\mu > 0$  such that

$$m_\mu(t-s) \leq \mu^\varepsilon(x, t)t - \mu^\varepsilon(x, s)s \leq M_\mu(t-s), \quad \text{for } t \geq s \geq 0, x \in \bar{\Omega},$$

then it is shown in [30, Lemma 2.1] that  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$  hold. This is satisfied in particular if  $\mu^\varepsilon(x, s)$  is bounded and there exists  $\lambda_\mu, L_\mu > 0$  such that  $L_\mu < \lambda_\mu \leq \mu^\varepsilon(x, s)$  and

$$\left| \frac{d\mu^\varepsilon}{ds}(x, s) \right| \leq \frac{L_\mu}{1+s}, \quad \forall s \in [0, \infty), x \in \bar{\Omega}, \varepsilon > 0.$$

Examples similar to the ones numerically investigated in [28] and fulfilling these assumptions are in particular  $\mu(x, s) = 2 + (1+s)^{-1}$  and  $\mu(x, s) = 2 + \exp(-s^2)$ .

**Homogenization for parabolic monotone problems.** The process of homogenization aims at characterizing the weak limit function  $u^0 \in E$  of the family of solutions  $\{u^\varepsilon\}$  of (2) as the solution of an effective partial differential equation, the homogenized equation. Although the weak convergence in  $E$  of a subsequence of  $\{u^\varepsilon\}$  follows directly from standard compactness arguments, the homogenization theory shows the existence of an effective model for (2) by means of parabolic  $G$ -convergence techniques. In [37, 36], the following convergence result has been derived: there exists a subsequence  $\{u^\varepsilon\}$  (again indexed by  $\varepsilon$ ), a map  $\mathcal{A}^0: \Omega \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  and  $u^0 \in E$  such that

$$\begin{aligned} u^\varepsilon &\rightharpoonup u^0 \text{ in } L^2(0, T; H_0^1(\Omega)), & \partial_t u^\varepsilon &\rightharpoonup \partial_t u^0 \text{ in } L^2(0, T; H^{-1}(\Omega)), \\ a^\varepsilon(x, \nabla u^\varepsilon) \nabla u^\varepsilon &\rightharpoonup \mathcal{A}^0(x, \nabla u^0) \text{ in } L^2(0, T; (L^2(\Omega))^d), \end{aligned}$$

where  $u^0 \in E$  can be characterized as the unique solution of the homogenized problem

$$\begin{aligned} \partial_t u^0(x, t) - \operatorname{div}(\mathcal{A}^0(x, \nabla u^0(x, t))) &= f(x), & \text{in } \Omega \times (0, T), \\ u^0(x, t) &= 0, & \text{on } \partial\Omega \times (0, T), \\ u^0(x, 0) &= g(x), & \text{in } \Omega. \end{aligned} \tag{5}$$

The existence and uniqueness of the solution  $u^0$  of (5) is deduced using  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$  which can be shown to hold also for the effective problem (5) (possibly with different constants).

**Remark 2.2.** An explicit representation of the map  $\mathcal{A}^0$  is only available for tensors  $a^\varepsilon$  with a particular structure (like periodicity or ergodicity), analogously to the linear case. In the case of locally periodic tensors  $a^\varepsilon(x, \xi) = a(x, \frac{x}{\varepsilon}, \xi)$  where  $a(x, y, \xi)$  is  $Y$ -periodic in  $y$ , it is shown in [27] for elliptic monotone problems, that the homogenized map  $\mathcal{A}^0(x, \xi)$  can be represented by

$$\mathcal{A}^0(x, \xi) = \int_Y a(x, y, \xi + \nabla \chi^\xi(x, y)) (\xi + \nabla \chi^\xi(x, y)) dy, \tag{6}$$

where  $\chi^\xi(x, \cdot) \in W_{per}^1(Y)$  solves

$$\int_Y a(x, y, \xi + \nabla \chi^\xi(x, y)) (\xi + \nabla \chi^\xi(x, y)) \cdot \nabla q(y) dy = 0, \quad \forall q \in W_{per}^1(Y). \tag{7}$$

This representation is true also in our context of parabolic problems because the homogenized map  $\mathcal{A}^0(x, \xi)$  is identical to the elliptic case, see [37]. Further, we note that for spatially periodic tensors  $a^\varepsilon$ , the homogenized map  $\mathcal{A}^0$  can be decomposed following  $\mathcal{A}^0(\xi) = a^0(\xi)\xi$  where  $a^0(\xi) \in \mathbb{R}^{d \times d}$  is the homogenized tensor, see [9].

### 3 Nonlinear and linearized multiscale methods

In this section, we recall the nonlinear multiscale method introduced in [6] and propose a new linearized method. Both methods rely on micro and macro finite element spaces.

### 3.1 Micro and macro finite element spaces

**Macroscopic spatial discretization.** Let  $\mathcal{T}_H$  be a shape-regular triangulation of the polygonal domain  $\Omega$  consisting of open simplices  $K \in \mathcal{T}_H$  with straight edges. The index  $H$  of the macroscopic triangulation  $\mathcal{T}_H$  denotes the macro mesh size  $H = \max_{K \in \mathcal{T}_H} \text{diam } K$  where  $\text{diam } K$  denotes the diameter of a simplex  $K \in \mathcal{T}_H$ . Further, for  $K \in \mathcal{T}_H$ , the measure and the barycenter of  $K$  are denoted by  $|K|$  and  $x_K$ , respectively.

Associated to the macro triangulation  $\mathcal{T}_H$  we introduce the finite element space  $S_0^1(\Omega, \mathcal{T}_H)$  consisting of piecewise affine functions

$$S_0^1(\Omega, \mathcal{T}_H) = \{v^H \in C^0(\overline{\Omega}) \cap H_0^1(\Omega) \mid v^H|_K \in \mathcal{P}^1(K) \text{ for all } K \in \mathcal{T}_H\}, \quad (8)$$

where  $\mathcal{P}^1(K)$  denotes the space of affine polynomials on  $K \in \mathcal{T}_H$ .

**Microscopic spatial discretization.** Let  $K \in \mathcal{T}_H$  be a macroscopic element. To perform localized microscopic simulations we introduce sampling domains  $K_\delta$  of microscopic size centered at the barycenter  $x_K$  of the macro element given by  $K_\delta = x_K + \delta(-\frac{1}{2}, \frac{1}{2})^d$ , with  $\delta \geq \varepsilon$ . The sampling domain  $K_\delta$  is discretized by a microscopic triangulation  $\mathcal{T}_h$  of open simplices  $T \in \mathcal{T}_h$  with straight edges. Here, the parameter  $h$  denotes the microscopic mesh size  $h = \max_{T \in \mathcal{T}_h} \text{diam } T$ . Further, let  $W(K_\delta) \subset H^1(K_\delta)$  be a Hilbert space. Then, the microscopic finite element space  $S^1(K_\delta, \mathcal{T}_h)$  is defined by

$$S^1(K_\delta, \mathcal{T}_h) = \{v^h \in C^0(\overline{K_\delta}) \cap W(K_\delta) \mid v^h|_T \in \mathcal{P}^1(T) \text{ for all } T \in \mathcal{T}_h\}, \quad (9)$$

where  $\mathcal{P}^1(T)$  is the set of affine polynomials on  $T \in \mathcal{T}_h$ .

**Time discretization.** The time domain  $(0, T)$  is discretized into  $N$  subintervals  $(t_{n-1}, t_n)$  of identical length  $\Delta t = T/N$ , where  $\Delta t$  is called the time step size and  $t_n = n\Delta t$ .

### 3.2 Nonlinear FE-HMM

We recall here the nonlinear FE-HMM proposed and analyzed in [6].

**Nonlinear macro method.** Let  $u_0^H \in S_0^1(\Omega, \mathcal{T}_H)$  be an approximation of the initial conditions  $g(x)$ . The sequence of numerical approximations  $\{u_n^H\} \subset S_0^1(\Omega, \mathcal{T}_H)$  generated by the nonlinear multiscale method proposed in [6], solves the nonlinear recursion

$$\int_{\Omega} \frac{1}{\Delta t} (u_{n+1}^H - u_n^H) w^H dx + \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \nabla \hat{u}_{K,n+1}^h) \nabla \hat{u}_{K,n+1}^h dx \cdot \nabla w^H(x_K) = \int_{\Omega} f w^H dx, \quad (10)$$

$$\forall w^H \in S_0^1(\Omega, \mathcal{T}_H),$$

where  $0 \leq n \leq N-1$  and  $\hat{u}_{K,n+1}^h$  is the solution to the nonlinear micro problem (11) constrained by the macro state  $v^H = u_{n+1}^H$ . In the analysis, we shall use the compact notation  $\bar{\partial}_t u_n^H = \frac{1}{\Delta t} (u_{n+1}^H - u_n^H)$  for the backward difference quotient in (10).

**Nonlinear micro problems.** For  $K \in \mathcal{T}_H$  and  $v^H \in S_0^1(\Omega, \mathcal{T}_H)$  fixed, consider the nonlinear micro problem: find  $\hat{v}_K^h$  such that  $\hat{v}_K^h - v^H \in S^1(K_\delta, \mathcal{T}_h)$  and

$$\int_{K_\delta} a^\varepsilon(x, \nabla \hat{v}_K^h) \nabla \hat{v}_K^h \cdot \nabla q^h dx = 0, \quad \forall q^h \in S^1(K_\delta, \mathcal{T}_h). \quad (11)$$

We recall that the nonlinear micro problem (11) is well-defined because  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi)\xi$  is assumed to satisfy  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$ .

### 3.3 Linearized FE-HMM

In contrast to the nonlinear FE-HMM described above, the idea of the linearized FE-HMM is to represent the solution at the macro and the micro scale using the following product of finite element spaces

$$\mathcal{S}^{H,h} = S_0^1(\Omega, \mathcal{T}_H) \times \prod_{K \in \mathcal{T}_H} S^1(K_\delta, \mathcal{T}_h), \quad (12)$$

where  $S_0^1(\Omega, \mathcal{T}_H)$  and  $S^1(K_\delta, \mathcal{T}_h)$  are defined in (8) and (9), respectively. An element  $\hat{z} = (z^H, \{z_K^h\}) \in \mathcal{S}^{H,h}$  thus consists of a macroscopic finite element function  $z^H \in S_0^1(\Omega, \mathcal{T}_H)$  and a family of microscopic functions  $\{z_K^h\}_{K \in \mathcal{T}_H}$  where  $z_K^h \in S^1(K_\delta, \mathcal{T}_h)$  for every sampling domain  $K_\delta$ . Further, for  $x \in K_\delta$ , we define  $\hat{v}_K^h(x) = v^H(x) + v_K^h(x)$ .

**Modified macro bilinear form.** For a given  $\hat{z} = (z^H, \{z_K^h\}) \in \mathcal{S}^{H,h}$  we introduce the bilinear form  $B^H(\hat{z}; \cdot, \cdot)$  for macroscopic functions  $v^H, w^H \in S_0^1(\Omega, \mathcal{T}_H)$  by

$$B^H(\hat{z}; v^H, w^H) = \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \nabla \hat{z}_K^h(x)) \nabla \hat{v}_K^{h,\hat{z}}(x) dx \cdot \nabla w^H(x_K), \quad (13)$$

where  $\hat{v}_K^{h,\hat{z}} = v^H(x) + v_K^{h,\hat{z}}(x)$ , with  $v_K^{h,\hat{z}}(x) \in S^1(K_\delta, \mathcal{T}_h)$ , is the solution of the micro problem (14) with parameter  $\hat{z}_K^h$  and macro constraint  $v^H$ .

**Micro problems.** The proposed multiscale strategy is driven by simulations at the microscopic scale. To upscale the microscopic behavior linked to a given macroscopic state we introduce constrained micro problems on the sampling domains. For  $K \in \mathcal{T}_H$ ,  $\hat{z} \in \mathcal{S}^{H,h}$  and  $v^H \in S_0^1(\Omega, \mathcal{T}_H)$  fixed, we introduce the micro problem: find  $\hat{v}_K^{h,\hat{z}}(x) = v^H(x) + v_K^{h,\hat{z}}(x)$  with  $v_K^{h,\hat{z}} \in S^1(K_\delta, \mathcal{T}_h)$  such that

$$\int_{K_\delta} a^\varepsilon(x, \nabla \hat{z}_K^h(x)) \nabla \hat{v}_K^{h,\hat{z}}(x) \cdot \nabla q^h(x) dx = 0, \quad \forall q^h \in S^1(K_\delta, \mathcal{T}_h). \quad (14)$$

Notice that problem (14) is linear, in contrast to problem (11). The coupling between the macroscopic state  $v^H$  and the solution  $\hat{v}_K^{h,\hat{z}}$  to the micro problem (14) is imposed by the choice of the subspace  $W(K_\delta) \subset H^1(K_\delta)$  implicitly encoded into the micro finite element space  $S^1(K_\delta, \mathcal{T}_h)$  defined in (9). In this article we consider two different coupling conditions

- periodic coupling:  $W(K_\delta) = W_{per}^1(K_\delta) = \{v \in H_{per}^1(K_\delta) \mid \int_{K_\delta} v(x) dx = 0\}$  and  $\delta/\varepsilon \in \mathbb{N}_{>0}$ ;
- Dirichlet coupling:  $W(K_\delta) = H_0^1(K_\delta)$  and  $\delta > \varepsilon$ .

We next explain here the construction of the linearized FE-HMM solution  $u_n^H$  approximating the homogenized solution  $u^0$  in (5) at time  $t = n\Delta t$ . We first describe the scheme starting from  $\hat{u}_1$  given at time  $t_1 = \Delta t$ . The procedure to construct  $\hat{u}_1$  is discussed afterwards.

**Linearized macro method.** Let  $\hat{u}_1 = (u_1^H, \{u_{1,K}^h\}) \in \mathcal{S}^{H,h}$  be given, then the sequence  $\{\hat{u}_n\}$  is defined by the following *linear* recursion. For  $1 \leq n \leq N-1$ , each time step of the multiscale method corresponds to the map  $\hat{u}_n \mapsto \hat{u}_{n+1} = (u_{n+1}^H, \{u_{n+1,K}^h\}) \in \mathcal{S}^{H,h}$  defined as

- (i) *evolution of the macroscopic state:* find  $u_{n+1}^H \in S_0^1(\Omega, \mathcal{T}_H)$ , the solution of the linear problem

$$\int_{\Omega} \frac{1}{\Delta t} (u_{n+1}^H - u_n^H) w^H dx + B^H(\hat{u}_n; u_{n+1}^H, w^H) = \int_{\Omega} f w^H dx, \quad \forall w^H \in S_0^1(\Omega, \mathcal{T}_H); \quad (15)$$

- (ii) *update the microscopic states:* for  $K \in \mathcal{T}_H$ , compute

$$u_{n+1,K}^h := v_K^{h,\hat{u}_n} \quad (16)$$

the solution to the micro problem (14) with parameter  $\hat{z} = \hat{u}_n$  and macro constraint  $v^H = u_{n+1}^H$ .

**Initialization procedure.** We next discuss how to define  $\hat{u}_1$  for the linearized scheme (15). Let  $u_0^H \in S_0^1(\Omega, \mathcal{T}_H)$  be an approximation of the initial state  $g(x)$ . For instance, a natural choice is  $u_0^H = \mathcal{I}_H g$  where  $\mathcal{I}_H$  is the nodal interpolant (39), but our analysis is valid for general initial conditions  $u_0^H$ . To be able to start the linearized multiscale method (15) an element  $\hat{u}_1 = (u_1^H, \{u_{1,K}^h\}) \in \mathcal{S}^{H,h}$  with micro functions  $u_{1,K}^h \in S^1(K_\delta, \mathcal{T}_h)$  is required. A trivial initialization would be to set  $\hat{u}_0 = (u_0^H, \{0\})$  and to calculate  $\hat{u}_1$  using the linearized multiscale method (15) with  $n = 0$ , but this would deteriorate the accuracy. We thus propose to use one single time step of the fully nonlinear multiscale method (10), which allows to prove optimal convergence of the temporal error. Let  $u_1^H$  be the numerical solution of (10) at time  $t_1 = \Delta t$  and  $\hat{u}_{1,K}^h(x)$  the associated solutions to the nonlinear micro problems (11). We then initialize the linearized multiscale method at time  $t_1 = \Delta t$  with

$$\hat{u}_1 = (u_1^H, \{\hat{u}_{1,K}^h - u_1^H\}) \in \mathcal{S}^{H,h}. \quad (17)$$

**Remark 3.1.** We emphasize once again that in the linearized FE-HMM defined above, both the macroscopic state equation (15) and the independent micro problems (14) are linear, in contrast to the nonlinear FE-HMM (10) which involves nonlinear and coupled problems at both the macro and micro scales. Indeed, observe in (15) that the form  $B^H$  is evaluated with  $B^H(\hat{u}_n; u_{n+1}^H, w^H)$  instead of  $B^H(\hat{u}_{n+1}; u_{n+1}^H, w^H)$ , where the nonlinear parameter  $\hat{u}_n$  is already known. Since  $B^H(\hat{u}_n; \cdot, \cdot)$  is a bilinear form, this means that the cost of solving (15) is analogous to that of the implicit Euler method applied to a linear parabolic finite element problem. In terms of memory storage, notice that the space  $\mathcal{S}^{H,h}$  used to represent the numerical solution of the linearized FE-HMM  $\hat{u}_n \mapsto \hat{u}_{n+1}$  is the macro state  $u_n^H$  and a vector  $\nabla \hat{u}_n$  for each sampling domain  $K_\delta$ , whereas only the macro state  $u_n^H$  is needed for the nonlinear FE-HMM (10).

## 4 Main results

In this section, we derive the well-posedness and a priori convergence estimates of the proposed linearized FE-HMM.

### 4.1 Well-posedness of the numerical method

The well-posedness of the linearized FE-HMM relies on the following lemma.

**Lemma 4.1.** *Assume that (3) holds. Then, the form defined in (13) satisfies for all  $\hat{z} \in \mathcal{S}^{H,h}$ ,  $v^H, w^H \in S_0^1(\Omega, \mathcal{T}_H)$ ,*

$$B^H(\hat{z}; v^H, v^H) \geq \lambda_a \|\nabla v^H\|_{L^2(\Omega)}^2, \quad |B^H(\hat{z}; v^H, w^H)| \leq \frac{\Lambda_a^2}{\lambda_a} \|\nabla v^H\|_{L^2(\Omega)} \|\nabla w^H\|_{L^2(\Omega)},$$

with the constants  $\lambda_a$  and  $\Lambda_a$  from (3). Thus,  $B^H(\hat{z}; \cdot, \cdot)$  is elliptic and bounded on  $S_0^1(\Omega, \mathcal{T}_H) \times S_0^1(\Omega, \mathcal{T}_H)$  (uniformly in  $\hat{z}$ ).

*Proof.* First, the existence and uniqueness of a solution to the constrained linear micro problems (14) is clear as the tensor  $a^\varepsilon$  is uniformly elliptic and bounded, see (3). Further, we note that Lemma 4.1 is a generalization of a result known for FE-HMM applied to linear elliptic problems, see [1, Proposition 3.2]. The proof relies on the fundamental energy equivalence

$$\|\nabla v^H\|_{L^2(K_\delta)} \leq \left\| \nabla \hat{v}_K^{h, \hat{z}} \right\|_{L^2(K_\delta)} \leq \frac{\Lambda_a}{\lambda_a} \|\nabla v^H\|_{L^2(K_\delta)},$$

where  $\hat{v}_K^{h, \hat{z}}$  is the solution to the linear micro problem (14) with parameter  $\hat{z}$  and constraint  $v^H$ .  $\square$

**Lemma 4.2.** *Let  $u_0^H \in S_0^1(\Omega, \mathcal{T}_H)$ ,  $f \in L^2(\Omega)$  and assume that (3) and  $(\mathcal{A}_{1-2})$  hold. Then, for all  $H, h$  and  $\Delta t$  the sequence  $\{u_n^H\}_{1 \leq n \leq N}$  defined by the linearized method (15) using the nonlinear initialization (17) exists, is unique and satisfies the a priori bound*

$$\max_{1 \leq n \leq N} \|u_n^H\|_{L^2(\Omega)} + \min\{\lambda, \lambda_a\} \left( \sum_{n=1}^N \Delta t \|\nabla u_n^H\|_{L^2(\Omega)}^2 \right)^{1/2} \leq C(\|f\|_{L^2(\Omega)} + \|u_0^H\|_{L^2(\Omega)}),$$

where  $C$  depends on the ellipticity constant  $\lambda_a$  of  $a^\varepsilon$ , the monotonicity constant  $\lambda$  of  $\mathcal{A}^\varepsilon$ , the final time  $T$  and the Poincaré constant  $C_p$  of the domain  $\Omega$ .

*Proof.* First, we note that the existence, uniqueness and boundedness of the nonlinear initialization (17) has been proved in [6, Theorem 3.5] using the hypotheses  $(\mathcal{A}_{1-2})$ . In particular, we have the bound

$$\|u_1^H\|_{L^2(\Omega)}^2 - \|u_0^H\|_{L^2(\Omega)}^2 + \lambda \Delta t \|\nabla u_1^H\|_{L^2(\Omega)}^2 \leq \frac{1}{\lambda} \Delta t C_p^2 \|f\|_{L^2(\Omega)}^2. \quad (18)$$

Next, for  $\hat{z} \in \mathcal{S}^{H,h}$  and  $v^H, w^H \in S_0^1(\Omega, \mathcal{T}_H)$ , we introduce the bilinear form  $A^{H, \Delta t}(\hat{z}; \cdot, \cdot)$  and the linear form  $F^{H, \Delta t}(\hat{z}; \cdot)$  by

$$A^{H, \Delta t}(\hat{z}; v^H, w^H) = \frac{1}{\Delta t} \int_{\Omega} v^H w^H dx + B^H(\hat{z}; v^H, w^H), \quad F^{H, \Delta t}(\hat{z}; w^H) = \frac{1}{\Delta t} \int_{\Omega} z^H w^H dx + \int_{\Omega} f w^H dx.$$



Thus, for  $1 \leq n \leq N - 1$ , the evolution of the macroscopic state (15) is equivalent to

$$A^{H,\Delta t}(\hat{u}_n; u_{n+1}^H, w^H) = F^{H,\Delta t}(\hat{u}_n; w^H), \quad \forall w^H \in S_0^1(\Omega, \mathcal{T}_H).$$

The ellipticity and boundedness of  $A^{H,\Delta t}(\hat{z}; \cdot, \cdot)$  and the continuity of  $F^{H,\Delta t}(\hat{z}; \cdot)$  follow from Lemma 4.1. Thus, the variational problem (15) has a unique solution for every  $1 \leq n \leq N - 1$ .

Next, we prove the boundedness of the numerical solution  $\{u_n^H\}$ . First, we observe that for  $0 \leq n \leq N - 1$  it holds

$$\int_{\Omega} \bar{\partial}_t u_n^H u_{n+1}^H dx \geq \frac{1}{2} \bar{\partial}_t \|u_n^H\|_{L^2(\Omega)}^2. \quad (19)$$

Thus, for  $1 \leq n \leq N - 1$ , the inequality (19) and the uniform ellipticity of  $B^H$ , see Lemma 4.1, lead to

$$\begin{aligned} \frac{1}{2\Delta t} \left( \|u_{n+1}^H\|_{L^2(\Omega)}^2 - \|u_n^H\|_{L^2(\Omega)}^2 \right) + \lambda_a \|\nabla u_{n+1}^H\|_{L^2(\Omega)}^2 &\leq \int_{\Omega} \bar{\partial}_t u_n^H u_{n+1}^H dx + B^H(\hat{u}_n; u_{n+1}^H, u_{n+1}^H) \\ &= \int_{\Omega} f u_{n+1}^H dx \leq \frac{C_p^2}{2\lambda_a} \|f\|_{L^2(\Omega)}^2 + \frac{\lambda_a}{2} \|\nabla u_{n+1}^H\|_{L^2(\Omega)}^2, \end{aligned} \quad (20)$$

where we used the definition of the method (15), the Poincaré inequality (with constant  $C_p$ ) and Young's inequality. We conclude by combining the inequalities (18) and (20) summed from  $n = 1$  to  $n = N - 1$ .  $\square$

## 4.2 A priori error estimates

In this section, we derive rigorous a priori error estimates for the proposed linearized FE-HMM with two different sets of assumptions. In the first case, we assume directly the monotonicity and Lipschitz continuity of the map  $\mathcal{A}^\varepsilon$  and we make a smallness assumption on the size of the nonlinearity of the problem. We note that such type of smallness assumption is commonly used in the numerical analysis of nonlinear PDEs, e.g., see [7, Theorem 4] or [11, Section 8.7]. In the second case, under the conditions on the tensor  $a^\varepsilon(x, \xi)$  derived in Remark 2.1 error estimates are shown without this smallness assumption. However, the result in the second case is obtained at the expense of assuming that a certain linearization error denoted by  $e_{n,K}$ , see (31), is small enough. In Section 6.1 we illustrate with numerical tests that this hypothesis is indeed satisfied for sufficiently fine discretizations of the space-time domain.

To estimate the error introduced by the numerical upscaling procedure built into the multiscale strategy (15), we define

$$r_{HMM}(\nabla v^H) = \left( \sum_{K \in \mathcal{T}_H} |K| \left| \mathcal{A}^0(x_K, \nabla v^H(x_K)) - \mathcal{A}_K^{0,h}(\nabla v^H) \right|^2 \right)^{1/2}, \quad \text{for } v^H \in S_0^1(\Omega, \mathcal{T}_H), \quad (21)$$

where  $\mathcal{A}^0$  is the exact homogenized map from the homogenized equation (5) and  $\mathcal{A}_K^{0,h}$  is the numerically homogenized map defined in (35). In particular, in the a priori error estimates of Theorems 4.3 and 4.4 the upscaling error is quantified by  $e_{HMM}$  given by

$$e_{HMM} = \max_{1 \leq n \leq N} r_{HMM}(\nabla \mathcal{U}_n^H), \quad (22)$$

where  $\mathcal{U}_n^H \in S_0^1(\Omega, \mathcal{T}_H)$  is a finite element approximation of the homogenized solution  $u^0$  at time  $t_n$ , for  $1 \leq n \leq N$ . In the analysis, we consider either  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$  the nodal interpolant (39) or  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$  the elliptic projection (41), as detailed in Section 5.1.

**Error estimates using a smallness assumption on the nonlinearity.** To derive our first a priori error estimate, we assume additionally that the tensor  $a^\varepsilon(x, \xi)$  is uniformly Lipschitz continuous in the second variable  $\xi$ , i.e., there exists a constant  $\tilde{L}_a > 0$  such that

$$\|a^\varepsilon(x, \xi_1) - a^\varepsilon(x, \xi_2)\|_{\mathcal{F}} \leq \tilde{L}_a |\xi_1 - \xi_2|, \quad \forall \xi_1, \xi_2 \in \mathbb{R}^d, \text{ a.e. } x \in \Omega. \quad (23)$$

For  $\xi \in \mathbb{R}^d$  and  $K \in \mathcal{T}_H$ , we introduce the exact micro function  $\bar{\chi}_K^\xi$  solving the variational problem: find  $\bar{\chi}_K^\xi \in W(K_\delta)$  such that

$$\int_{K_\delta} a^\varepsilon(x, \xi + \nabla \bar{\chi}_K^\xi)(\xi + \nabla \bar{\chi}_K^\xi) \cdot \nabla q dx = 0, \quad \forall q \in W(K_\delta), \quad (24)$$

and its finite element approximation  $\chi_K^{\xi,h} \in S^1(K_\delta, \mathcal{T}_h)$  satisfying

$$\int_{K_\delta} a^\varepsilon(x, \xi + \nabla \chi_K^{\xi,h})(\xi + \nabla \chi_K^{\xi,h}) \cdot \nabla q^h dx = 0, \quad \forall q^h \in S^1(K_\delta, \mathcal{T}_h). \quad (25)$$

Note that for  $v^H \in S_0^1(\Omega, \mathcal{T}_H)$  and  $\xi = \nabla v^H(x_K)$  we recover  $\xi + \nabla \chi_K^{\xi,h} = \hat{v}_K^h$  where  $\hat{v}_K^h$  solves (11).

Analogously to linear elliptic problems, the exact solution  $\bar{\chi}_K^\xi$  satisfies the bound  $\|\nabla \bar{\chi}_K^\xi\|_{L^2(K_\delta)} \leq C\sqrt{|K_\delta|}|\xi|$ , where  $C$  is independent of  $\varepsilon$  and  $\xi$ . Under additional regularity of the data of the nonlinear micro problem (24) the Lipschitz continuity of its solution  $\bar{\chi}_K^\xi$  can be shown, e.g., see [29, Theorem 4.1]. For our analysis it is necessary to know the explicit dependence of the Lipschitz constant of  $\bar{\chi}_K^\xi$  with respect to  $\varepsilon$  and  $\xi$ . We assume that

$$\mathbf{(R1)} \quad \left\| \nabla \bar{\chi}_K^\xi \right\|_{L^\infty(K_\delta)} \leq C^* |\xi| \text{ for } K \in \mathcal{T}_H, \xi \in \mathbb{R}^d.$$

Further, we use the affine bijection  $G_{K_\delta} : Y \rightarrow K_\delta$  between the micro cell domain  $K_\delta$  and the unit cell  $Y = (0, 1)^d$ , and define  $\chi_{K,Y}^{\xi,\hat{h}} = \chi_K^{\xi,h} \circ G_{K_\delta}$  and  $\bar{\chi}_{K,Y}^\xi = \bar{\chi}_K^\xi \circ G_{K_\delta}$  where  $\hat{h}$  is the mesh size of the rescaled partition  $\mathcal{T}_{\hat{h}}$  on  $Y$  obtained from  $\mathcal{T}_h$  via the bijection  $G_{K_\delta}$ . If the partition  $\mathcal{T}_{\hat{h}}$  is quasi-uniform<sup>1</sup> and the rescaled micro problems (with solution  $\bar{\chi}_{K,Y}^\xi$ ) are regular enough, then the maximum norm estimates for nonlinear monotone elliptic problems derived in [24] combined with an inverse inequality, see [12, Theorem 3.2.6], yield

$$\left\| \chi_{K,Y}^{\xi,\hat{h}} - \bar{\chi}_{K,Y}^\xi \right\|_{W^{1,\infty}(Y)} \leq C\hat{h}^{-1} \left\| \chi_{K,Y}^{\xi,\hat{h}} - \bar{\chi}_{K,Y}^\xi \right\|_{L^\infty(Y)} \leq C\hat{h} \left| \log \hat{h} \right|^{\frac{d}{4}+1}, \quad (26)$$

where  $C$  is independent of  $\hat{h}$  and  $\varepsilon$ , with unknown explicit dependence on  $|\xi|$ . Analogously to **(R1)** we postulate that  $C$  scales linearly with  $|\xi|$ . By transferring the bound (26) back to the sampling domain  $K_\delta$  and observing that  $\delta = \mathcal{O}(\varepsilon)$  we obtain that

$$\mathbf{(R2)} \quad \left\| \nabla \chi_K^{\xi,h} - \nabla \bar{\chi}_K^\xi \right\|_{L^\infty(K_\delta)} \leq C^* \frac{h}{\varepsilon} |\xi| \text{ for } K \in \mathcal{T}_H, \xi \in \mathbb{R}^d,$$

where  $C^* = C|\log(h/\varepsilon)|^{\frac{d}{4}+1}$  is weakly depending on  $h/\varepsilon$ .

We may now state our first a priori error estimate on the linearized FE-HMM based on the smallness assumption (28) on the nonlinearity.

**Theorem 4.3.** *Let  $u^0$  be the solution to the homogenized problem (5) and  $u_n^H$  the approximations defined by the linearized multiscale method (15) using the nonlinear initialization (17). Assume that the tensor  $a^\varepsilon$  satisfies the assumptions (3), (23) and that the map  $\mathcal{A}^\varepsilon$  given by  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi)\xi$  satisfies assumptions  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$ . Let the following conditions be valid for  $\mu = 1$  and some constant  $L_0 > 0$ ,*

$$\begin{aligned} u^0, \partial_t u^0 &\in \mathcal{C}^0([0, T], H^2(\Omega)), \quad \partial_t^2 u^0 \in \mathcal{C}^0([0, T], L^2(\Omega)), \\ \mathcal{A}^0(\cdot, \xi) &\in W^{\mu, \infty}(\Omega; \mathbb{R}^d) \quad \text{with} \quad \|\mathcal{A}^0(\cdot, \xi)\|_{W^{\mu, \infty}(\Omega; \mathbb{R}^d)} \leq C(L_0 + |\xi|), \quad \forall \xi \in \mathbb{R}^d. \end{aligned} \quad (27)$$

Assume further that  $\bar{\chi}_K^\xi$  and  $\chi_K^{\xi,h}$  given by (24) and (25), respectively, satisfy assumptions **(R1)** and **(R2)**. If the exact solution  $u^0$  verifies

$$u^0 \in \mathcal{C}^0([0, T], W^{2, \infty}(\Omega)), \quad \sqrt{2}(1 + C^*)\tilde{L}_a \max_{t \in [0, T]} |u^0(x, t)|_{W^{1, \infty}(\Omega)} < \lambda_a, \quad (28)$$

(where  $C^*$  is the constant from **(R1)**) then there exist  $H_0, h_0 > 0$  such that for any  $H < H_0, h < h_0$ , we have the a priori error estimates

$$\max_{1 \leq n \leq N} \|u_n^H - u^0(x, t_n)\|_{L^2(\Omega)} \leq C(\Delta t + H^\mu + e_{HMM} + \|u_0^H - g\|_{L^2(\Omega)}), \quad (29a)$$

$$\left( \Delta t \sum_{n=1}^N \|\nabla u_n^H - \nabla u^0(x, t_n)\|_{L^2(\Omega)}^2 \right)^{1/2} \leq C(\Delta t + H + e_{HMM} + \|u_0^H - g\|_{L^2(\Omega)}), \quad (29b)$$

where  $\mu = 1$ , the upscaling error  $e_{HMM}$  (evaluated at the nodal interpolant  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$ , see (39)) is given by (22), and  $C$  is independent of  $\Delta t, H$  and  $e_{HMM}$ .

Further, the estimate (29a) holds with  $\mu = 2$  and  $e_{HMM}$  evaluated at the elliptic projection  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$ , see (41), if additionally (27) holds for  $\mu = 2$  and conditions (42) are satisfied.

<sup>1</sup>Precisely, there exists a constant  $C > 0$  such that  $\hat{h}_k \geq C\hat{h}$  for all sizes  $\hat{h}_k$  of the finite elements  $k \in \mathcal{T}_{\hat{h}}$ .

**Error estimates without the smallness assumption on the nonlinearity.** We note that conditions (3) and (23) together do not imply  $(\mathcal{A}_1)$  or  $(\mathcal{A}_2)$  in general for the map  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi)\xi$ . For our second error estimate, we make the assumptions

$$a^\varepsilon(x, \cdot) \in (W^{1,\infty}(\mathbb{R}^d))^{d \times d} \quad \text{with} \quad \left( \sum_{i,j,k=1}^d \left| \frac{\partial a_{ij}^\varepsilon(\xi)}{\partial \xi_k} \right|^2 \right)^{1/2} \leq \frac{L_a}{1 + |\xi|} \quad \text{and} \quad L_a < \frac{\lambda_a}{2\sqrt{2}}, \quad (30)$$

$$\forall \xi \in \mathbb{R}^d, \text{ a.e. } x \in \Omega.$$

We recall that condition (30) combined with (3) is sufficient to imply  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$  for  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi)\xi$ , as shown in Remark 2.1.

For  $0 \leq n \leq N$  and  $K \in \mathcal{T}_H$ , we consider the error term  $e_{n,K} \in (L^\infty(K_\delta))^{d \times d}$  given by

$$e_{n,K}(x) = a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) - \int_0^1 a^\varepsilon(x, \nabla \hat{u}_{n,K}^h - \tau \nabla \hat{\theta}_{n,K}^h) d\tau, \quad \text{a.e. } x \in K_\delta, \quad (31)$$

where  $\hat{u}_n = (u_n^H, \{u_{n,K}^h\}) \in \mathcal{S}^{H,h}$  is obtained from the numerical method (15) and  $\hat{\theta}_n = (\theta_n^H, \{\theta_{n,K}^h\})$  denotes the difference  $\hat{\theta}_n = \hat{u}_n - \hat{U}_n$  between the numerical solution  $\hat{u}_n$  and an approximation  $\hat{U}_n$  of the homogenized solution  $u^0$  (and its associated first order correctors), see (47) for details. Thus, if  $\hat{u}_n$  is a good approximation to  $u^0$  one can expect that  $e_{n,K}$  is small.

**Theorem 4.4.** *Let  $u^0$  be the solution to the homogenized problem (5) and  $u_n^H$  the approximations defined by the linearized multiscale method (15) using the nonlinear initialization (17). Assume that the tensor  $a^\varepsilon$  satisfies (3) and (30). Further, let the homogenized solution  $u^0$  and the homogenized map  $\mathcal{A}^0$  satisfy the regularity assumptions (27) for  $\mu = 1$ . If*

$$\max_{\substack{K \in \mathcal{T}_H \\ 1 \leq n \leq N-1}} \|e_{n,K}\|_{(L^\infty(K_\delta))^{d \times d}} < \frac{\lambda_a}{2\sqrt{2}}, \quad (32)$$

then we have the a priori error estimates

$$\max_{1 \leq n \leq N} \|u_n^H - u^0(x, t_n)\|_{L^2(\Omega)} \leq C (\Delta t + H^\mu + e_{HMM} + \|u_0^H - g\|_{L^2(\Omega)}), \quad (33a)$$

$$\left( \Delta t \sum_{n=1}^N \|\nabla u_n^H - \nabla u^0(x, t_n)\|_{L^2(\Omega)}^2 \right)^{1/2} \leq C (\Delta t + H + e_{HMM} + \|u_0^H - g\|_{L^2(\Omega)}), \quad (33b)$$

where  $\mu = 1$ , the upscaling error  $e_{HMM}$  (evaluated at the nodal interpolant  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$ , see (39)) is given by (22), and  $C$  is independent of  $\Delta t, H$  and  $e_{HMM}$ .

If additionally  $u^0 \in C^0([0, T], W^{2,\infty}(\Omega))$ , conditions (27) hold for  $\mu = 2$  and hypotheses (42) are satisfied then there exists  $H_0 > 0$  such that for any  $H < H_0$  the estimate (33a) holds with  $\mu = 2$  and  $e_{HMM}$  evaluated at the elliptic projection  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$ , see (41).

**Remark 4.5.** We observe that the monotonicity of the model problem allows us to avoid the smallness assumption (28) in the above theorem. The alternative condition (32) assumes that the quantity  $e_{n,K}$  defined in (31) is small enough. We note that this assumption automatically holds for linear problems. For nonlinear problems however  $\|e_{n,K}\|_{(L^\infty(K_\delta))^{d \times d}}$  can be bounded by  $C \|\nabla \hat{u}_{n,K}^h - \nabla \hat{U}_{n,K}^h\|_{L^\infty(K_\delta)}$  for which convergence results are difficult to derive. However, in Section 6.1, we provide numerical evidence that the condition (32) is verified for spatial and temporal discretizations that are fine enough.

**Fully discrete a priori error estimates.** To derive fully discrete estimates taking into account the finescale discretization and upscaling errors, it remains to bound the error  $e_{HMM}$  defined in (22) and involved in Theorems 4.3 and 4.4. Explicit estimates for  $e_{HMM}$  rely on a decomposition of the total upscaling error into modeling and microscopic error denoted by  $e_{mod}$  and  $e_{mic}$ , respectively,

$$e_{HMM} \leq \max_{1 \leq n \leq N} \left( \sum_{K \in \mathcal{T}_H} |K| |\mathcal{A}^0(x_K, \nabla \mathcal{U}_n^H(x_K)) - \bar{\mathcal{A}}_K^0(\nabla \mathcal{U}_n^H)|^2 \right)^{1/2} \quad (= e_{mod})$$

$$+ \max_{1 \leq n \leq N} \left( \sum_{K \in \mathcal{T}_H} |K| |\bar{\mathcal{A}}_K^0(\nabla \mathcal{U}_n^H) - \mathcal{A}_K^{0,h}(\nabla \mathcal{U}_n^H)|^2 \right)^{1/2}, \quad (= e_{mic}) \quad (34)$$

where  $\mathcal{A}^0$  is the exact homogenized map and the approximated maps  $\bar{\mathcal{A}}_K^0$  and  $\mathcal{A}_K^{0,h}$  are given by

$$\bar{\mathcal{A}}_K^0(\xi) = \frac{1}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \xi + \nabla \bar{\chi}_K^\xi)(\xi + \nabla \bar{\chi}_K^\xi) dx, \quad \mathcal{A}_K^{0,h}(\xi) = \frac{1}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \xi + \nabla \chi_K^{\xi,h})(\xi + \nabla \chi_K^{\xi,h}) dx, \quad (35)$$

where  $\bar{\chi}_K^\xi$  and  $\chi_K^{\xi,h}$  solve the micro problem (24) and (25), respectively, for  $K \in \mathcal{T}_H$ ,  $\xi \in \mathbb{R}^d$ .

For explicit convergence rates of the micro error  $e_{mic}$  with respect to the micro mesh size  $h$  appropriate regularity of the exact solution  $\bar{\chi}_K^\xi$  of the nonlinear micro problem (24) is required, e.g., see [3]. Further, to obtain *optimal* convergence rates for  $e_{mic}$  the adjoint micro problems (36) are crucial. Such adjoint problems have already been introduced for linear multiscale problems with non-symmetric tensors, see [17].

In view of those results, we assume that  $\bar{\chi}_K^\xi$  solving the nonlinear micro problems (24) satisfy

$$\mathbf{(H1)} \quad \bar{\chi}_K^\xi \in H^2(K_\delta) \text{ and } \left| \bar{\chi}_K^\xi \right|_{H^2(K_\delta)} \leq C\varepsilon^{-1} |\xi| \sqrt{|K_\delta|} \text{ for } K \in \mathcal{T}_H, \xi \in \mathbb{R}^d,$$

and introduce the linear variational problem: find  $\bar{X}_K^{\xi,j} \in W(K_\delta)$  such that

$$\int_{K_\delta} \left( D_\xi \mathcal{A}^\varepsilon(x, \xi + \nabla \bar{\chi}_K^\xi) \right)^T (e_j + \nabla \bar{X}_K^{\xi,j}) \cdot \nabla z \, dx = 0, \quad \forall z \in W(K_\delta), \quad (36)$$

where  $\bar{\chi}_K^\xi$  solves the nonlinear micro problem (24),  $1 \leq j \leq d$ ,  $K \in \mathcal{T}_H$  and  $\xi \in \mathbb{R}^d$ . Note that the micro problem (36) is well-defined as the Jacobian  $D_\xi \mathcal{A}^\varepsilon(x, \xi)$  is uniformly elliptic and bounded if  $\mathcal{A}^\varepsilon$  satisfies  $(\mathcal{A}_1)$ ,  $(\mathcal{A}_2)$  and is smooth enough. We assume further that the solution  $\bar{X}_K^{\xi,j}$  of the adjoint micro problem (36) fulfills

$$\mathbf{(H1^*)} \quad \begin{cases} \bar{X}_K^{\xi,j} \in H^2(K_\delta) & \text{and} & \left| \bar{X}_K^{\xi,j} \right|_{H^2(K_\delta)} \leq C\varepsilon^{-1} \sqrt{|K_\delta|} & \text{for } K \in \mathcal{T}_H, \xi \in \mathbb{R}^d, j = 1, \dots, d; \\ \bar{X}_K^{\xi,j} \in W^{1,\infty}(K_\delta) & \text{and} & \left| \bar{X}_K^{\xi,j} \right|_{W^{1,\infty}(K_\delta)} \leq C & \text{for } K \in \mathcal{T}_H, \xi \in \mathbb{R}^d, j = 1, \dots, d. \end{cases}$$

We note that the micro error  $e_{mic}$  can be estimated independently of the structure of the spatial variations of the tensor  $a^\varepsilon(x, \xi)$ . In contrast, to derive explicit estimates for the modeling error  $e_{mod}$  structural assumptions on the spatial heterogeneities of the tensor  $a^\varepsilon(x, \xi)$  are necessary. For linear multiscale problems, i.e., with tensors  $a^\varepsilon(x, \xi)$  independent of  $\xi \in \mathbb{R}^d$ , such results have been derived assuming local periodicity or random stationarity of the tensor, see [19]. In this article, we present explicit estimates for

$\mathbf{(H2)}$  locally periodic tensor  $a^\varepsilon(x, \xi) = a^\varepsilon(x, \frac{x}{\varepsilon}, \xi) = a(x, y, \xi)$  which is  $Y$ -periodic in  $y$  and satisfies

$$\|a(x_1, y, \xi) - a(x_2, y, \xi)\|_{\mathcal{F}} \leq C|x_1 - x_2|, \quad \forall x_1, x_2 \in \Omega, \xi \in \mathbb{R}^d, \text{ a.e. } y \in Y.$$

**Theorem 4.6.** *Let  $u^0$  and  $u_n^H$  be the exact homogenized solution and the numerical solution defined by the linearized multiscale method (15) with nonlinear initialization (17), respectively. Assume the hypotheses of Theorem 4.3 and Theorem 4.4 such that optimal error estimates in the  $L^2(H^1)$  and  $C^0(L^2)$  hold for the time discretization and macroscopic spatial error.*

*If additionally  $a_{ij}^\varepsilon(x, \cdot) \in W^{2,\infty}(\mathbb{R}^d)$  for  $1 \leq i, j \leq d$  and assumptions  $\mathbf{(H1)}$  and  $\mathbf{(H1^*)}$  are satisfied then we obtain the optimal fully discrete error estimates*

$$\begin{aligned} \max_{1 \leq n \leq N} \|u^0(\cdot, t_n) - u_n^H\|_{L^2(\Omega)} &\leq C \left[ \Delta t + H^2 + \left(\frac{h}{\varepsilon}\right)^2 + e_{mod} + \|u_0^H - g\|_{L^2(\Omega)} \right], \\ \left( \Delta t \sum_{n=1}^N \|\nabla u^0(\cdot, t_n) - \nabla u_n^H\|_{L^2(\Omega)}^2 \right)^{1/2} &\leq C \left[ \Delta t + H + \left(\frac{h}{\varepsilon}\right)^2 + e_{mod} + \|u_0^H - g\|_{L^2(\Omega)} \right], \end{aligned}$$

where  $e_{mod}$  is the modeling error and  $C$  is independent of  $\Delta t$ ,  $H$ ,  $h$ ,  $\varepsilon$  and  $\delta$ . Further, if  $a^\varepsilon$  satisfies  $\mathbf{(H2)}$  with  $a^\varepsilon(x, \xi)$  replaced by  $a(x_K, x/\varepsilon, \xi)$  in (13) and (14), and if  $\chi^\xi(x_K, \cdot) \in W^{1,\infty}(Y)$  in (7), then

$$e_{mod} \leq \begin{cases} 0, & W(K_\delta) = W_{per}^1(K_\delta), \delta/\varepsilon \in \mathbb{N}_{>0}, \\ C\left(\frac{\varepsilon}{\delta}\right)^{1/2}, & W(K_\delta) = H_0^1(K_\delta), \delta > \varepsilon, \end{cases}$$

where  $C$  is independent of  $\varepsilon$  and  $\delta$ .

## 5 Analysis

This section is devoted to the proofs of our main results stated in the previous section.

### 5.1 Preliminaries

We introduce two semi-norms on the product space  $\mathcal{S}^{H,h}$  defined in (12). For  $\hat{v} = (v^H, \{v_K^h\}) \in \mathcal{S}^{H,h}$ , we define

$$\|\nabla \hat{v}\|_{\mathcal{S}^{H,h}} = \left( \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \|\nabla \hat{v}_K^h\|_{L^2(K_\delta)}^2 \right)^{1/2}, \quad \|\nabla \hat{v}\|_{\mathcal{S}_\infty^{H,h}} = \max_{K \in \mathcal{T}_H} \|\nabla \hat{v}_K^h\|_{L^\infty(K_\delta)}.$$

We observe that for all  $K \in \mathcal{T}_H$  it holds

$$\|\nabla \hat{v}_K^h\|_{L^2(K_\delta)}^2 = \|\nabla v^H(x_K)\|_{L^2(K_\delta)}^2 + \|\nabla v_K^h\|_{L^2(K_\delta)}^2, \quad (37)$$

because  $\int_{K_\delta} \nabla v_K^h dx \cdot \nabla v^H(x_K) = 0$  due to cancelling and vanishing values of  $v_K^h$  on the boundary  $\partial K_\delta$  for micro spaces  $S^1(K_\delta, \mathcal{T}_h)$  with periodic and Dirichlet boundary conditions, respectively. Combining that with the Poincaré inequality (for  $H_0^1(\Omega)$  and  $W(K_\delta) = H_0^1(K_\delta)$ ) and the Poincaré-Wirtinger inequality (if  $W(K_\delta) = W_{per}^1(K_\delta)$  instead) proves that  $\|\cdot\|_{\mathcal{S}^{H,h}}$  is a norm. Further, we note that the identity (37) yields

$$\|\nabla v^H\|_{L^2(\Omega)} \leq \|\nabla \hat{v}\|_{\mathcal{S}^{H,h}}, \quad \text{for all } \hat{v} = (v^H, \{v_K^h\}) \in \mathcal{S}^{H,h}. \quad (38)$$

We shall use in our analysis the nodal interpolant

$$\mathcal{I}_H: \mathcal{C}^0(\bar{\Omega}) \rightarrow S^1(\Omega, \mathcal{T}_H), \quad \text{see [12, Section 2.4]}, \quad (39)$$

where  $S^1(\Omega, \mathcal{T}_H)$  is the space of continuous, piecewise affine functions on the macro mesh  $\mathcal{T}_H$  (compared to  $S_0^1(\Omega, \mathcal{T}_H)$  we do not impose Dirichlet boundary conditions on the boundary  $\partial\Omega$ ). We recall that the choice  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$  (for  $0 \leq n \leq N$ ) combined with the regularity assumptions (27) for  $\mu = 1$  is sufficient to derive optimal error estimates in the  $L^2(0, T; H_0^1(\Omega))$  norm.

However, to obtain sharp error estimates in the  $\mathcal{C}^0([0, T], L^2(\Omega))$  norm, more involved techniques have been used in [6]. In particular, an appropriate elliptic projection  $\tilde{u}^{H,0}$  has been introduced. Herein, we just recall its definition and refer to [6] for the details. For  $t \in [0, T]$  and  $v, w \in H_0^1(\Omega)$ , let the bilinear form  $B_\pi$  be given by

$$B_\pi(t; v, w) = \int_{\Omega} \mathcal{A}^0(x, t) \nabla v \cdot \nabla w \, dx, \quad \text{with } \mathcal{A}^0(x, t) = D_\xi \mathcal{A}^0(x, \nabla u^0(x, t)), \quad (40)$$

where the homogenized map  $\mathcal{A}^0$  and the homogenized solution  $u^0$  are assumed to be smooth enough. Let  $t \in [0, T]$ , the elliptic projection  $\tilde{u}^{H,0}(\cdot, t)$  of  $u^0(\cdot, t)$  solves the variational problem: find  $\tilde{u}^{H,0}(\cdot, t) \in S_0^1(\Omega, \mathcal{T}_H)$  such that

$$B_\pi(t; \tilde{u}^{H,0}(\cdot, t), w^H) = B_\pi(t; u^0(\cdot, t), w^H), \quad \forall w^H \in S_0^1(\Omega, \mathcal{T}_H). \quad (41)$$

Choosing  $\mathcal{U}_n^H = \tilde{u}^{H,0} = \tilde{u}^{H,0}(\cdot, t_n)$  in (47), optimal a priori error estimates in the  $\mathcal{C}^0([0, T], L^2(\Omega))$  norm have been derived in [6] under appropriate regularity assumptions. Those estimates are based on the following lemma summarizing Lemmas 5.4, 5.5 from [6].

**Lemma 5.1.** *Let  $u^0$  be the exact homogenized solution and  $\tilde{u}^{H,0}$  its elliptic projection (41). Assume that  $\mathcal{A}^0$  satisfies  $(\mathcal{A}_{1-2})$  and consider  $B_\pi$  and  $\mathcal{A}^0$  defined in (40). If  $u^0 \in \mathcal{C}^0([0, T], W^{2,\infty}(\Omega))$ ,  $\partial_t u^0 \in \mathcal{C}^0([0, T], H^2(\Omega))$  and the following holds,*

$$\begin{aligned} \mathcal{A}^0(x, \cdot) \in W^{2,\infty}(\mathbb{R}^d; \mathbb{R}^d), \quad \text{a.e. } x \in \Omega, \quad \mathcal{A}_{ij}^0, \partial_t \mathcal{A}_{ij}^0 \in \mathcal{C}^0([0, T], W^{1,\infty}(\Omega)), \quad 1 \leq i, j, \leq d, \\ \text{quasi-uniformity of meshes } \mathcal{T}_H, \text{ e.g., see [12, Eq. (3.2.28)], and elliptic regularity (44),} \end{aligned} \quad (42)$$

then for all  $t \in [0, T]$  and  $k \in \{0, 1\}$ ,  $s \in \{1, 2\}$ , we have the optimal error estimates

$$\|\partial_t^k u^0(\cdot, t) - \partial_t^k \tilde{u}^{H,0}(\cdot, t)\|_{H^{2-s}(\Omega)} \leq CH^s, \quad \|u^0(\cdot, t) - \tilde{u}^{H,0}(\cdot, t)\|_{W^{1,\infty}(\Omega)} \leq CH \|u^0(\cdot, t)\|_{W^{2,\infty}(\Omega)}, \quad (43)$$

where  $C$  is independent of  $H$  and the  $W^{1,\infty}$  estimate is valid for  $H < H_0$  for some  $H_0 > 0$ .

Note that the elliptic regularity assumed in (42) reads as: for  $1 < p < \sigma$  with some  $\sigma > d$  we have

$$\|u^0(\cdot, t)\|_{W^{2,p}(\Omega)} + \|u^{0,*}(\cdot, t)\|_{W^{2,p}(\Omega)} \leq C \|\operatorname{div}(\mathcal{A}^0(\cdot, t) \nabla u^0(\cdot, t))\|_{L^p(\Omega)}, \quad \forall t \in [0, T], \quad (44)$$

where  $u^{0,*}(\cdot, t)$  solves the dual problem  $B_\pi(t; w, u^{0,*}(\cdot, t)) = B_\pi(t; u^0(\cdot, t), w)$  for all  $w \in H_0^1(\Omega)$ .

We now recall several estimates on the nonlinear FE-HMM (10) from [6] that will be useful for the analysis of the proposed linearized version of the method.

**Lemma 5.2.** *Consider  $v^H, w^H \in S_0^1(K_\delta, \mathcal{T}_h)$  and  $K \in \mathcal{T}_H$ . Assume that  $\mathcal{A}^\varepsilon$  satisfies  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$ . Let  $\hat{v}_K^h$  and  $\hat{w}_K^h$  be the solutions to the nonlinear micro problem (11) associated to the macro functions  $v^H$  and  $w^H$ , respectively. Further, let the map  $\mathcal{A}_K^{0,h}: S_0^1(\Omega, \mathcal{T}_H) \rightarrow \mathbb{R}^d$  be given by (35), then*

$$\begin{aligned} \|\nabla \hat{v}_K^h - \nabla \hat{w}_K^h\|_{L^2(K_\delta)} &\leq \frac{L}{\lambda} \sqrt{|K_\delta|} |\nabla v^H(x_K) - \nabla w^H(x_K)|, \\ |\mathcal{A}_K^{0,h}(\nabla v^H) - \mathcal{A}_K^{0,h}(\nabla w^H)| &\leq \frac{L^2}{\lambda} |\nabla v^H(x_K) - \nabla w^H(x_K)|, \end{aligned}$$

where  $L$  and  $\lambda$  are the Lipschitz and monotonicity constants of the map  $\mathcal{A}^\varepsilon(x, \xi)$ , see  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$ .

*Proof.* The monotonicity  $(\mathcal{A}_2)$ , the definition of the nonlinear micro problems (11) and the Lipschitz continuity  $(\mathcal{A}_1)$  of the map  $\mathcal{A}^\varepsilon$  lead to

$$\begin{aligned} \lambda \|\nabla \hat{v}_K^h - \nabla \hat{w}_K^h\|_{L^2(K_\delta)}^2 &\leq \int_{K_\delta} [\mathcal{A}^\varepsilon(x, \nabla \hat{v}_K^h) - \mathcal{A}^\varepsilon(x, \nabla \hat{w}_K^h)] \cdot (\nabla \hat{v}_K^h - \nabla \hat{w}_K^h) dx \\ &= \int_{K_\delta} [\mathcal{A}^\varepsilon(x, \nabla \hat{v}_K^h) - \mathcal{A}^\varepsilon(x, \nabla \hat{w}_K^h)] \cdot (\nabla v^H(x_K) - \nabla w^H(x_K)) dx \\ &\leq L \sqrt{|K_\delta|} \|\nabla \hat{v}_K^h - \nabla \hat{w}_K^h\|_{L^2(K_\delta)} |\nabla v^H(x_K) - \nabla w^H(x_K)|, \end{aligned}$$

from where the first inequality follows. Further, combining the definition (35), the Lipschitz continuity of  $\mathcal{A}^\varepsilon$  with the first part of Lemma 5.2 concludes the proof.  $\square$

The following lemma summarizes Lemmas 5.6, 5.7, 5.8 and Corollary 5.10 derived in [6], its proof is omitted.

**Lemma 5.3.** *Let  $u^0$  be the homogenized solution and  $\mathcal{U}_n^H$  be given either by the nodal interpolant  $\mathcal{I}_H u^0(\cdot, t_n)$  or the elliptic projection  $\tilde{u}_n^{H,0}$ , see (41). Further, assume that (3) holds for the tensor  $a^\varepsilon$  and that the map  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi)\xi$  satisfies  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$ . Then, under the regularity assumptions (27) for  $\mu = 1$ , we have that for  $0 \leq n \leq N$*

$$\|\nabla \mathcal{U}_{n+1}^H - \nabla \mathcal{U}_n^H\|_{L^2(\Omega)} \leq C \Delta t. \quad (45)$$

Further, for  $0 \leq n \leq N - 1$  and  $w^H \in S_0^1(\Omega, \mathcal{T}_H)$ , the following estimates hold with  $\mu = 1$

$$\left| \int_{\Omega} [\partial_t u^0(x, t_{n+1}) - \bar{\partial}_t \mathcal{U}_n^H] w^H dx \right| \leq C(\Delta t + H^2) \|w^H\|_{L^2(\Omega)}, \quad (46a)$$

$$\begin{aligned} \left| \int_{\Omega} \mathcal{A}^0(x, \nabla u^0(x, t_{n+1})) \cdot \nabla w^H dx - \sum_{K \in \mathcal{T}_H} |K| \mathcal{A}_K^{0,h}(\nabla \mathcal{U}_{n+1}^H) \cdot \nabla w^H(x_K) \right| \\ \leq C(H^\mu + r_{HMM}(\nabla \mathcal{U}_{n+1}^H)) \|\nabla w^H\|_{L^2(\Omega)}, \end{aligned} \quad (46b)$$

where  $r_{HMM}$  is defined in (21),  $\mathcal{U}_n^H$  is given by the nodal interpolant  $\mathcal{I}_H u^0$  and  $C$  is independent of  $\Delta t$ ,  $H$  and  $r_{HMM}$ .

If alternatively  $\mathcal{U}_n^H$  is given by the elliptic projection  $\tilde{u}_n^{H,0}$ , assume that hypotheses (27) hold for  $\mu = 2$  and additionally  $u^0 \in \mathcal{C}^0([0, T], W^{2,\infty}(\Omega))$  as well as hypotheses (42) are satisfied. Then, the first estimate (45) still holds and there exist some  $H_0 > 0$  such that for  $H < H_0$  the estimates (46) hold for  $\mu = 2$ .

Optimal convergence rates for the error between upscaled maps  $\mathcal{A}_K^{0,h}(\xi)$  and  $\bar{\mathcal{A}}_K^0(\xi)$  were first presented in [1] for linear elliptic problems, generalized to high order in [2, Lemma 10], [4, Corollary 10]. It is extended to the case of non-symmetric tensors in [17] introducing appropriate adjoint cell problems (see also [8, Lemma 4.6] in the context of nonlinear nonmonotone problems). For the class of nonlinear problems (1), the following lemma estimates the modeling error  $e_{mod}$  and the micro error  $e_{mic}$  defined in (34). The estimates on  $e_{mod}$  and  $e_{mic}$  are shown in Section 5.4.2 and Section 5.4.1 in [6]. The proof of Lemma 5.4 is thus omitted.

**Lemma 5.4.** Let  $e_{mod}$  and  $e_{mic}$  be the modeling and micro error introduced in (34), respectively. Assume that the tensor  $a^\varepsilon$  satisfies (3) and the map  $\mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi)\xi$  satisfies the conditions  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$ . Further, if  $e_{HMM}$  is evaluated at the nodal interpolant  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$  of the homogenized solution  $u^0$  we assume the regularity  $u^0 \in C^0([0, T], H^2(\Omega))$ . However, if  $e_{HMM}$  is evaluated at the elliptic projection  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$ , see (41), then we assume that  $\tilde{u}_n^{H,0}$  is uniformly bounded in the  $W^{1,\infty}$  norm.

If the multiscale tensor  $a^\varepsilon$  satisfies  $a_{ij}^\varepsilon(x, \cdot) \in W^{2,\infty}(\mathbb{R}^d)$  for  $1 \leq i, j \leq d$  (for a.e.  $x \in \Omega$ ) and hypotheses **(H1)** and **(H1\*)** hold then the micro error  $e_{mic}$  can be explicitly estimated by

$$e_{mic} \leq C \left( \frac{h}{\varepsilon} \right)^2$$

where  $C$  is independent of  $h$  and  $\varepsilon$ . Further, assume that  $a^\varepsilon(x, \xi)$  satisfies **(H2)** and is replaced by  $a(x_K, x/\varepsilon, \xi)$  in (13) and (14). If  $\chi^\varepsilon(x_K, \cdot) \in W^{1,\infty}(Y)$  in (6), then the modeling error is bounded by

$$e_{mod} \leq \begin{cases} 0, & W(K_\delta) = W_{per}^1(K_\delta), \delta/\varepsilon \in \mathbb{N}_{>0}, \\ C(\frac{\varepsilon}{\delta})^{1/2}, & W(K_\delta) = H_0^1(K_\delta), \delta > \varepsilon, \end{cases}$$

where  $C$  is independent of  $\varepsilon$  and  $\delta$ .

**Proof of Theorem 4.6.** Combining Lemma 5.4 with Theorems 4.3 and 4.4, we immediately deduce Theorem 4.6.  $\square$

It remains to prove Theorems 4.3 and 4.4, this is the purpose of the next section.

## 5.2 Proof of the a priori error estimates

For  $0 \leq n \leq N$ , let  $\hat{u}_n = (u_n^H, \{u_{n,K}^h\}) \in \mathcal{S}^{H,h}$  be the numerical solution obtained by the *linearized* multi-scale method (15),(16). In our analysis, we shall first consider the case where the nonlinear initialization (17) is not used, i.e. for a given  $\hat{u}_0 = (u_0^H, \{u_{0,K}^h\}) \in \mathcal{S}^{H,h}$  at time  $t_0 = 0$ , the sequence  $\{\hat{u}_n\}$  is defined using (15) for all  $n \geq 0$ , including the first step  $\hat{u}_1$  with  $n = 0$  in (15),(16). We shall derive our a priori error estimates in terms of an initialization error  $e_{init}$  defined below. Then, we will show how to take advantage of the nonlinear initialization defined in (17) to derive the claimed error estimates.

Consider the elements  $\hat{\mathcal{U}}_n = (\mathcal{U}_n^H, \{\mathcal{U}_{n,K}^h\}) \in \mathcal{S}^{H,h}$  such that  $\hat{\mathcal{U}}_{n,K}^h$  is the solution to the *nonlinear* micro problem (11) constrained by the macro function  $\mathcal{U}_n^H$ . We define  $\hat{\theta}_n \in \mathcal{S}^{H,h}$  by

$$\hat{\theta}_n = \hat{u}_n - \hat{\mathcal{U}}_n, \quad \text{i.e.,} \quad \theta_n^H = u_n^H - \mathcal{U}_n^H, \quad \hat{\theta}_{n,K}^h = \hat{u}_{n,K}^h - \hat{\mathcal{U}}_{n,K}^h, \quad 0 \leq n \leq N, K \in \mathcal{T}_H. \quad (47)$$

Using notation (47), we define the initialization error  $e_{init}$ ,

$$e_{init} = \|\theta_0^H\|_{L^2(\Omega)} + \sqrt{\Delta t} \|\nabla \hat{\theta}_0\|_{\mathcal{S}^{H,h}}. \quad (48)$$

Our analysis will show in particular that the  $\sqrt{\Delta t}$  term in (48) can be removed from the a priori error estimates. In what follows, we take  $\mathcal{U}_n^H$  as either the nodal interpolant  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$  defined in (39) or the elliptic projection  $\mathcal{U}_n^H = \tilde{u}^{H,0}(\cdot, t_n)$  in (41).

**Lemma 5.5.** Assume that  $a^\varepsilon$  satisfies (3) and that conditions (27), which depend on  $\mu \in \{1, 2\}$ , hold. Let either  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$  be the nodal interpolant (39) and  $\mu = 1$  or  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$  be the elliptic projection (41) and  $\mu = 2$  (when additionally assuming  $u^0 \in C^0([0, T], W^{2,\infty}(\Omega))$  and (43)). Then, for any  $0 \leq n \leq N - 1$  and  $\hat{w} \in \mathcal{S}^{H,h}$

$$\begin{aligned} \int_{\Omega} \bar{\partial}_t \theta_n^H w^H dx + \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \nabla \hat{\theta}_{n+1,K}^h \cdot \nabla \hat{w}_K^h dx \\ \leq C(\Delta t + H^\mu + r_{HMM}(\nabla \mathcal{U}_{n+1}^H)) \|\nabla \hat{w}\|_{\mathcal{S}^{H,h}} + L_n(\nabla \hat{w}), \end{aligned} \quad (49)$$

where the constant  $C$  is independent of  $H, \Delta t$ , the upscaling error  $r_{HMM}$  is defined in (21) and the linearization error functional  $L_n: \mathcal{S}^{H,h} \rightarrow \mathbb{R}$  is given by

$$L_n(\nabla \hat{w}) = \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} \left[ a^\varepsilon(x, \nabla \hat{\mathcal{U}}_{n,K}^h) - a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \right] \nabla \hat{\mathcal{U}}_{n,K}^h \cdot \nabla \hat{w}_K^h dx. \quad (50)$$

*Proof.* As first step, we derive an error propagation formula for the sequence  $\{\hat{\theta}_n\}$ . Let  $0 \leq n \leq N-1$  and  $\hat{w} \in \mathcal{S}^{H,h}$ . Using the definition of the multiscale method (15) and the weak formulation of the effective equation (5) at time  $t_{n+1}$  we obtain

$$\begin{aligned} & \int_{\Omega} \bar{\partial}_t \theta_n^H w^H dx + \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \nabla \hat{\theta}_{n+1,K}^h \cdot \nabla \hat{w}_K^h dx \\ &= \int_{\Omega} \bar{\partial}_t u_n^H w^H dx + B^H(\hat{u}_n^H; u_{n+1}^H, w^H) \\ & \quad - \int_{\Omega} \bar{\partial}_t \mathcal{U}_n^H w^H dx - \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \nabla \hat{\mathcal{U}}_{n+1,K}^h \cdot \nabla \hat{w}_K^h dx \\ &= \int_{\Omega} [\partial_t u^0(x, t_{n+1}) - \bar{\partial}_t \mathcal{U}_n^H] w^H dx \end{aligned} \quad (51a)$$

$$+ \int_{\Omega} \mathcal{A}^0(x, \nabla u^0(x, t_{n+1})) \cdot \nabla w^H dx - \sum_{K \in \mathcal{T}_H} |K| \mathcal{A}_K^{0,h}(\nabla \mathcal{U}_{n+1}^H) \cdot \nabla w^H(x_K) \quad (51b)$$

$$+ \sum_{K \in \mathcal{T}_H} |K| \left[ \mathcal{A}_K^{0,h}(\nabla \mathcal{U}_{n+1}^H) \cdot \nabla w^H(x_K) - \frac{1}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \nabla \hat{\mathcal{U}}_{n+1,K}^h \cdot \nabla \hat{w}_K^h dx \right], \quad (51c)$$

where the numerically homogenized nonlinear map  $\mathcal{A}_k^{0,h}$  is given by (35). We note that the terms (51a) and (51b) can be bounded using Lemma 5.3.

Next, we estimate the term (51c), which is due to the linearization applied in the proposed multiscale method (15). Decomposing the error term (51c) and combining that with the results of Lemma 5.2 as well as the boundedness of the tensor  $a^\varepsilon$  postulated in (3) yields

$$\begin{aligned} (51c) &= \sum_{K \in \mathcal{T}_H} |K| \left[ \mathcal{A}_K^{0,h}(\nabla \mathcal{U}_{n+1}^H) - \mathcal{A}_K^{0,h}(\nabla \mathcal{U}_n^H) \right] \cdot \nabla w^H(x_K) \\ & \quad + \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} \left[ a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) - a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \right] \nabla \hat{\mathcal{U}}_{n,K}^h \cdot \nabla \hat{w}_K^h dx \\ & \quad + \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \left[ \nabla \hat{\mathcal{U}}_{n,K}^h - \nabla \hat{\mathcal{U}}_{n+1,K}^h \right] \cdot \nabla \hat{w}_K^h dx \\ & \leq \frac{L}{\lambda} (L + \Lambda_a) \|\nabla \mathcal{U}_{n+1}^H - \nabla \mathcal{U}_n^H\|_{L^2(\Omega)} \|\nabla \hat{w}\|_{\mathcal{S}^{H,h}} + L_n(\nabla \hat{w}), \end{aligned}$$

where  $L_n(\nabla \hat{w})$  is defined in (50). Further, in both cases  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$  and  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$ , see (39) and (41), respectively, the estimate  $\|\nabla \mathcal{U}_{n+1}^H - \nabla \mathcal{U}_n^H\|_{L^2(\Omega)} \leq C \Delta t$  holds due to the estimate (45) from Lemma 5.3 and the regularity assumptions (27) and (43).  $\square$

The following lemma states a priori error estimates analogous to those of Theorem 4.3 in the case where the nonlinear initialization (17) is not used.

**Lemma 5.6.** *Under the hypotheses of Theorem 4.3, consider the linearized FE-HMM (15) where in contrast to the nonlinear initialization (17), for a given  $\hat{u}_0 = (u_0^H, \{u_{0,K}^h\}) \in \mathcal{S}^{H,h}$  at time  $t_0 = 0$ , the value  $\hat{u}_1 = (u_1^H, \{u_{1,K}^h\}) \in \mathcal{S}^{H,h}$  at time  $t_1 = h$  is defined using (15), (16) with  $n = 0$ . Then, the conclusion of Theorem 4.3 holds, with  $\|u_0^H - g\|_{L^2(\Omega)}$  replaced by  $e_{init}$  defined in (48).*

*Proof.* Let  $0 \leq n \leq N-1$  and  $\hat{w} \in \mathcal{S}^{H,h}$ . We first estimate the linearization error  $L_n$ , see (50),

$$|L_n(\nabla \hat{w})| \leq \tilde{L}_a \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} \left| \nabla \hat{\theta}_{n,K}^h \right| \left| \nabla \hat{\mathcal{U}}_{n,K}^h \right| \left| \nabla \hat{w}_K^h \right| dx \leq \tilde{L}_a \left\| \nabla \hat{\mathcal{U}}_n \right\|_{\mathcal{S}_\infty^{H,h}} \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}^{H,h}} \|\nabla \hat{w}\|_{\mathcal{S}^{H,h}},$$

where we used the Lipschitz continuity (23) of  $a^\varepsilon(x, \cdot)$ . Thus,

$$|L_n(\nabla \hat{w})| \leq \mathcal{L}_n \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}^{H,h}} \|\nabla \hat{w}\|_{\mathcal{S}^{H,h}}, \quad \text{with } \mathcal{L}_n = \tilde{L}_a \left\| \nabla \hat{\mathcal{U}}_n \right\|_{\mathcal{S}_\infty^{H,h}}. \quad (52)$$



Then, we follow the lines of the proof of [6, Theorem 4.1]. Let us choose  $\hat{w} = \hat{\theta}_{n+1}$  in inequality (49). Due to (19), (38) and the uniform ellipticity of the tensor  $a^\varepsilon$  we obtain

$$\begin{aligned}
& \frac{1}{2\Delta t} \left( \|\theta_{n+1}^H\|_{L^2(\Omega)}^2 - \|\theta_n^H\|_{L^2(\Omega)}^2 \right) + \lambda_a \left\| \nabla \hat{\theta}_{n+1} \right\|_{\mathcal{S}^{H,h}}^2 \\
& \leq \int_{\Omega} \bar{\partial}_t \theta_n^H \theta_{n+1}^H dx + \sum_{K \in \mathcal{T}_H} \frac{|K|}{|K_\delta|} \int_{K_\delta} a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \nabla \hat{\theta}_{n+1,K}^h \cdot \nabla \hat{\theta}_{n+1,K}^h dx \\
& \leq C(\Delta t + H^\mu + r_{HMM}(\nabla \mathcal{U}_{n+1}^H)) \left\| \nabla \hat{\theta}_{n+1} \right\|_{\mathcal{S}^{H,h}} + \mathcal{L}_n \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}^{H,h}} \left\| \nabla \hat{\theta}_{n+1} \right\|_{\mathcal{S}^{H,h}} \\
& \leq C(\Delta t^2 + H^{2\mu} + r_{HMM}(\nabla \mathcal{U}_{n+1}^H)^2) + \frac{\mathcal{L}_n^2}{\lambda_a} \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}^{H,h}}^2 + \frac{\lambda_a}{2} \left\| \nabla \hat{\theta}_{n+1} \right\|_{\mathcal{S}^{H,h}}^2, \quad (53)
\end{aligned}$$

where we used Young's inequality for the last estimate. Let  $1 \leq K \leq N$ , then summing inequality (53) from  $n = 0$  to  $n = K - 1$  yields

$$\begin{aligned}
\left\| \theta_K^H \right\|_{L^2(\Omega)}^2 + \lambda_a \Delta t \sum_{n=1}^K \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}^{H,h}}^2 & \leq C(\Delta t^2 + H^{2\mu} + \max_{1 \leq n \leq K} r_{HMM}(\nabla \mathcal{U}_n^H)^2) + \left\| \theta_0^H \right\|_{L^2(\Omega)}^2 \\
& \quad + \frac{2}{\lambda_a} \Delta t \sum_{n=0}^{K-1} \mathcal{L}_n^2 \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}^{H,h}}^2. \quad (54)
\end{aligned}$$

As inequality (54) holds for any  $1 \leq K \leq N$  we derive

$$\begin{aligned}
\max_{1 \leq n \leq N} \left\| \theta_n^H \right\|_{L^2(\Omega)}^2 + \left( \lambda_a - \frac{2}{\lambda_a} \max_{1 \leq n \leq N-1} \mathcal{L}_n^2 \right) \Delta t \sum_{n=1}^N \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}^{H,h}}^2 \\
\leq C \left( \Delta t^2 + H^{2\mu} + \max_{1 \leq n \leq N} r_{HMM}(\nabla \mathcal{U}_n^H)^2 + \left\| \theta_0^H \right\|_{L^2(\Omega)}^2 + \Delta t \mathcal{L}_0^2 \left\| \nabla \hat{\theta}_0 \right\|_{\mathcal{S}^{H,h}}^2 \right), \quad (55)
\end{aligned}$$

which proves the convergence under the condition that

$$\lambda_a - \frac{2}{\lambda_a} \max_{1 \leq n \leq N-1} \mathcal{L}_n^2 > 0 \quad \Leftrightarrow \quad \sqrt{2} \tilde{L}_a \max_{1 \leq n \leq N-1} \left\| \nabla \hat{\mathcal{U}}_n \right\|_{\mathcal{S}_\infty^{H,h}} < \lambda_a, \quad (56)$$

where the explicit expression for  $\mathcal{L}_n$  is given in (52).

Next, we have to relate the condition (56) (a smallness assumption on  $\tilde{L}_a$  and the micro solutions to the nonlinear cell problem (11) constrained by  $\mathcal{U}_n^H$ ) to the condition (28) (a smallness assumption on  $\tilde{L}_a$  and the exact effective solution  $u^0$ ). As **(R1)**, **(R2)** and (43) hold we apply the result of Corollary 5.8. Thus, for every  $\eta > 0$  there exist  $H_0 > 0$  and  $h_0 > 0$  such that for  $H < H_0$  and  $h < h_0$  it holds

$$\sqrt{2} \tilde{L}_a \max_{1 \leq n \leq N-1} \left\| \nabla \hat{\mathcal{U}}_n \right\|_{\mathcal{S}_\infty^{H,h}} \leq \sqrt{2} \tilde{L}_a (1 + C^* + \eta) \sup_{t \in [0, T]} |u^0(\cdot, t)|_{W^{1,\infty}(\Omega)}, \quad (57)$$

where  $C^*$  is the constant from **(R1)**. Thus, for  $\eta > 0$  small enough the condition (56) follows from the smallness assumption (28).

Further, for the same parameters  $\eta, H_0$  and  $h_0$  as above one can show analogously that  $\mathcal{L}_0$  is bounded by the right-hand side of (57). Thus, using the boundedness of  $\mathcal{L}_0$  the terms of the right-hand side of (55) depending on  $\hat{\theta}_0$  can be estimated by

$$\left\| \theta_0^H \right\|_{L^2(\Omega)}^2 + \Delta t \mathcal{L}_0^2 \left\| \nabla \hat{\theta}_0 \right\|_{\mathcal{S}^{H,h}}^2 \leq C e_{init}^2, \quad (58)$$

where the initialization error  $e_{init}$  is defined in (48).

Combining the estimates  $\|u^0(\cdot, t_n) - \mathcal{U}_n^H\|_{H^{2-s}(\Omega)} \leq CH^s$  for  $s = 1, 2$ , which hold due to the regularity (27) (for  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$ ) and the additional assumption (43) (for  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$ ), and estimate (55) concludes the proof of Theorem 4.3.  $\square$

**Lemma 5.7.** For  $K \in \mathcal{T}_H$ ,  $\xi \in \mathbb{R}^d$ , let  $\chi_K^{\xi,h}$  and  $\bar{\chi}_K^\xi$  be the solution to the nonlinear micro problems (25) and (24), respectively. If **(R1)** and **(R2)** hold then for every  $\eta > 0$  there exists some  $h_0 > 0$  such that for all  $h < h_0$  we have

$$\left\| \xi + \nabla \chi_K^{\xi,h} \right\|_{L^\infty(K_\delta)} \leq (1 + C^* + \eta) |\xi|,$$

where  $C^*$  is the constant from **(R1)**.

*Proof.* The result follows by applying assumptions **(R1)** and **(R2)** to the decomposition

$$\left\| \xi + \nabla \chi_{K_K}^{\xi, h} \right\|_{L^\infty(K_\delta)} \leq |\xi| + \left\| \nabla \tilde{\chi}_K^\xi \right\|_{L^\infty(K_\delta)} + \left\| \nabla \chi_{K_K}^{\xi, h} - \nabla \tilde{\chi}_K^\xi \right\|_{L^\infty(K_\delta)}. \quad \square$$

**Corollary 5.8.** *Let  $0 \leq n \leq N$  and  $\mathcal{U}_n^H$  either be given by the nodal interpolant (39) of  $u^0(\cdot, t_n)$  or the elliptic projection (41). Assume that (43) additionally holds if  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$ . If **(R1)**, **(R2)** hold and  $u^0 \in \mathcal{C}^0([0, T], W^{2,\infty}(\Omega))$ , then for every  $\eta > 0$  there exist some  $H_0, h_0 > 0$  such that for  $H < H_0$  and  $h < h_0$  it holds that*

$$\max_{1 \leq n \leq N-1} \left\| \nabla \hat{\mathcal{U}}_n \right\|_{\mathcal{S}_{\infty}^{H,h}} \leq (1 + C^* + \eta) \sup_{t \in [0, T]} |u^0(\cdot, t)|_{W^{1,\infty}(\Omega)},$$

where  $C^*$  is the constant from **(R1)**.

*Proof.* If  $u^0 \in \mathcal{C}^0([0, T], W^{2,\infty}(\Omega))$  we have that  $\|\mathcal{U}_n^H - u^0(\cdot, t_n)\|_{W^{1,\infty}(\Omega)} \leq CH$  for  $\mathcal{U}_n^H = \mathcal{I}_H u^0(\cdot, t_n)$  (see [12, Theorem 3.1.6]) as well as for  $\mathcal{U}_n^H = \tilde{u}_n^{H,0}$  (if additionally (43) is satisfied). Combining that with Lemma 5.7 (for  $\xi = \nabla \mathcal{U}_n^H(x_K)$  and  $h$  small enough) yields

$$\begin{aligned} \max_{1 \leq n \leq N-1} \left\| \nabla \hat{\mathcal{U}}_n \right\|_{\mathcal{S}_{\infty}^{H,h}} &\leq (1 + C^* + \frac{\eta}{2}) \max_{1 \leq n \leq N-1} \left\| \nabla \mathcal{U}_n^H \right\|_{L^\infty(\Omega)} \\ &\leq (1 + C^* + \frac{\eta}{2}) \max_{1 \leq n \leq N-1} \left[ |\mathcal{U}_n^H - u^0(\cdot, t_n)|_{W^{1,\infty}(\Omega)} + |u^0(\cdot, t_n)|_{W^{1,\infty}(\Omega)} \right] \\ &\leq (1 + C^* + \eta) \sup_{t \in [0, T]} |u^0(\cdot, t)|_{W^{1,\infty}(\Omega)}, \end{aligned}$$

where the last step holds if  $H$  is small enough.  $\square$

Analogously to Lemma 5.6, we have the following lemma in the case where, in contrast to Theorem 4.4, the nonlinear initialization procedure (17) is not applied.

**Lemma 5.9.** *Under the hypotheses of Theorem 4.4, consider the linearized FE-HMM (15) where in contrast to the nonlinear initialization (17), for a given  $\hat{u}_0 = (u_0^H, \{u_{0,K}^h\}) \in \mathcal{S}^{H,h}$  at time  $t_0 = 0$ , the value  $\hat{u}_1 = (u_1^H, \{u_{1,K}^h\}) \in \mathcal{S}^{H,h}$  at time  $t_1 = h$  is defined using (15), (16) with  $n = 0$ . Then, the conclusion of Theorem 4.4 holds, with  $\|u_0^H - g\|_{L^2(\Omega)}$  replaced by  $e_{init}$  defined in (48).*

*Proof.* Compared to the proof of Theorem 4.3, this proof relies on a different estimate of the linearization functional  $L_n$  defined in (50). Let  $0 \leq n \leq N-1$  and  $\hat{w} \in \mathcal{S}_{H,h}$ . Instead of the estimate (52) we use the result of the technical Lemma 5.10, i.e.,

$$|L_n(\nabla \hat{w})| \leq \mathcal{L}_n \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}_{H,h}} \left\| \nabla \hat{w} \right\|_{\mathcal{S}_{H,h}}, \quad \text{with } \mathcal{L}_n = L_a + \max_{K \in \mathcal{T}_H} \|e_{n,K}\|_{(L^\infty(K_\delta))^{d \times d}}, \quad (59)$$

where  $L_a$  is the constant from (30) and  $e_{n,K}$  is given by (31).

Following the lines of the proof of Theorem 4.3 the convergence result can be shown again (cf. condition (56)) under the condition that

$$\lambda_a - \frac{2}{\lambda_a} \max_{1 \leq n \leq N-1} \mathcal{L}_n^2 > 0 \quad \Leftrightarrow \quad L_a + \max_{\substack{K \in \mathcal{T}_H \\ 1 \leq n \leq N-1}} \|e_{n,K}\|_{(L^\infty(K_\delta))^{d \times d}} < \frac{\lambda_a}{\sqrt{2}}, \quad (60)$$

where the right-hand side of the equivalence (60) is due to the definition (59) of  $\mathcal{L}_n$ . Then, it is easily seen that condition (60) holds due to the hypothesis  $L_a < \lambda_a/(2\sqrt{2})$  (ensuring monotonicity) from (30) and assumption (32). Finally, we observe that  $\mathcal{L}_0$  is bounded due to the boundedness (3) of  $a^\varepsilon$ , i.e.,  $\|e_{0,K}(x)\|_{\mathcal{F}} \leq C\Lambda_a$  for  $K \in \mathcal{T}_H$  and a.e.  $x \in K_\delta$ . Thus, the error terms depending on  $\hat{\theta}_0$  can again be bounded by  $e_{init}$ , cf. (58).  $\square$

**Lemma 5.10.** *Let  $\hat{w} \in \mathcal{S}^{H,h}$  and  $0 \leq n \leq N-1$ . If the tensor  $a^\varepsilon$  satisfies (30) then*

$$|L_n(\nabla \hat{w})| \leq \left( L_a + \max_{K \in \mathcal{T}_H} \|e_{n,K}\|_{(L^\infty(K_\delta))^{d \times d}} \right) \left\| \nabla \hat{\theta}_n \right\|_{\mathcal{S}_{H,h}} \left\| \nabla \hat{w} \right\|_{\mathcal{S}_{H,h}},$$

where  $L_a$  is the constant from (30) and  $e_{n,K}$  is defined in (31).

*Proof.* Let  $\hat{w} \in \mathcal{S}^{H,h}$  and  $0 \leq n \leq N-1$ , we recall that the explicit representation of  $L_n(\nabla \hat{w})$  is given in (50). Then, for  $1 \leq i, j \leq d$ ,  $K \in \mathcal{T}_H$  and a.e.  $x \in K_\delta$ , we have

$$a_{ij}^\varepsilon(x, \nabla \hat{U}_{n,K}^h) - a_{ij}^\varepsilon(x, \nabla \hat{u}_{n,K}^h) = - \int_0^1 (\nabla_\xi a^\varepsilon(x, \hat{d}_{n,K}(\tau)) [\nabla \hat{\theta}_{n,K}^h])_{ij} d\tau,$$

where  $\hat{d}_{n,K}(\tau) = \nabla \hat{U}_{n,K}^h + \tau \nabla \hat{\theta}_{n,K}^h$  (for  $0 \leq \tau \leq 1$ ) and  $\nabla_\xi a^\varepsilon(x, \xi) [\eta] \in \mathbb{R}^{d \times d}$  is defined by

$$(\nabla_\xi a^\varepsilon(x, \xi) [\eta])_{ij} = \nabla_\xi a_{ij}^\varepsilon(x, \xi) \cdot \eta, \quad \xi, \eta \in \mathbb{R}^d, \text{ a.e. } x \in K_\delta. \quad (61)$$

Thus, the integral in (50) can be expressed as

$$\begin{aligned} \int_{K_\delta} \left[ a^\varepsilon(x, \nabla \hat{U}_{n,K}^h) - a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) \right] \nabla \hat{U}_{n,K}^h \cdot \nabla \hat{w}_K^h dx \\ = - \int_{K_\delta} \int_0^1 \nabla_\xi a^\varepsilon(x, \hat{d}_{n,K}(\tau)) [\nabla \hat{\theta}_{n,K}^h] \nabla \hat{U}_{n,K}^h d\tau \cdot \nabla \hat{w}_K^h dx \end{aligned} \quad (62)$$

$$= - \int_{K_\delta} \int_0^1 I_{n,K}(x, \tau) d\tau \cdot \nabla \hat{w}_K^h dx + \int_{K_\delta} \int_0^1 \tilde{I}_{n,K}(x, \tau) d\tau \cdot \nabla \hat{w}_K^h dx, \quad (63)$$

where in the last line we decompose (62) into two parts with  $I_{n,K}, \tilde{I}_{n,K}: K_\delta \times (0, 1) \rightarrow \mathbb{R}^d$  given by

$$I_{n,K}(x, \tau) = \nabla_\xi a^\varepsilon(x, \hat{d}_{n,K}(\tau)) [\nabla \hat{\theta}_{n,K}^h] \hat{d}_{n,K}(\tau), \quad \tilde{I}_{n,K}(x, \tau) = \nabla_\xi a^\varepsilon(x, \hat{d}_{n,K}(\tau)) [\nabla \hat{\theta}_{n,K}^h] \tau \nabla \hat{\theta}_{n,K}^h, \quad (64)$$

which is well-defined for a.e.  $\tau \in (0, 1)$  and a.e.  $x \in K_\delta$ .

Recalling the definition (61) of  $\nabla_\xi a^\varepsilon$  and applying repeatedly the Cauchy-Schwarz inequality leads to

$$\begin{aligned} |I_{n,K}(\tau)| &\leq \sum_{i,j,k=1}^d \left| \frac{\partial a_{ij}^\varepsilon}{\partial \xi_k}(x, \hat{d}_{n,K}(\tau)) \right|^2 |\hat{d}_{n,K}(\tau)|^2 |\nabla \hat{\theta}_{n,K}^h|^2 \leq L_a^2 \frac{|\hat{d}_{n,K}(\tau)|^2}{(1 + |\hat{d}_{n,K}(\tau)|)^2} |\nabla \hat{\theta}_{n,K}^h|^2 \\ &\leq L_a^2 |\nabla \hat{\theta}_{n,K}^h|^2, \end{aligned} \quad (65)$$

for a.e.  $\tau \in [0, 1]$ ,  $x \in K_\delta$ , where the assumption (30) yields the second last inequality.

Finally, we study the term  $\tilde{I}_{n,K}$  defined in (64). Let  $1 \leq i, j \leq d$  we observe that for a.e.  $\tau \in (0, 1)$

$$\frac{\partial}{\partial \tau} [a_{ij}^\varepsilon(x, \hat{d}_{n,K}(\tau))] = \sum_{k=1}^d \frac{\partial a_{ij}^\varepsilon}{\partial \xi_k}(x, \hat{d}_{n,K}(\tau)) \frac{\partial}{\partial \tau} [\hat{d}_{n,K}(\tau) \cdot e_k] = (\nabla_\xi a^\varepsilon(x, \hat{d}_{n,K}(\tau)) [\nabla \hat{\theta}_{n,K}^h])_{ij}, \quad (66)$$

a.e. in  $K_\delta$ . Thus, the definition of  $\tilde{I}_{n,K}$  and  $\nabla_\xi a^\varepsilon$ , see (64) and (61), respectively, and identity (66) yield

$$\tilde{I}_{n,K}(\tau) \cdot e_i = \tau \sum_{j=1}^d (\nabla_\xi a^\varepsilon(x, \hat{d}_{n,K}(\tau)) [\nabla \hat{\theta}_{n,K}^h])_{ij} (\nabla \hat{\theta}_{n,K}^h \cdot e_j) = \sum_{j=1}^d \tau \frac{\partial}{\partial \tau} [a_{ij}^\varepsilon(x, \hat{d}_{n,K}(\tau))] (\nabla \hat{\theta}_{n,K}^h \cdot e_j). \quad (67)$$

By combining (67), integrating by parts and using the variable  $s = 1 - \tau$ , we obtain the representation

$$\begin{aligned} \int_0^1 \tilde{I}_{n,K}(\tau) \cdot e_i d\tau &= \sum_{j=1}^d \left[ a_{ij}^\varepsilon(x, \nabla \hat{u}_{n,K}^h) - \int_0^1 a_{ij}^\varepsilon(x, \nabla \hat{u}_{n,K}^h - s \nabla \hat{\theta}_{n,K}^h) ds \right] (\nabla \hat{\theta}_{n,K}^h \cdot e_j) \\ &= e_{n,K} \nabla \hat{\theta}_{n,K}^h \cdot e_i, \end{aligned} \quad (68)$$

for a.e.  $x \in K_\delta$ , where the definition (31) of  $e_{n,K}$  is used in the last line.

Thus, we conclude the proof by combining the definition (50), the decomposition (63), the estimate (65) for  $I_{n,K}$  and the exact representation (68) of  $\tilde{I}_{n,K}$ .  $\square$

Based on Lemma 5.6 and Lemma 5.9 involving the initialization error  $e_{init}$ , we may now prove Theorem 4.3 and Theorem 4.4 by taking advantage of the nonlinear initialization (17).

*Proof of Theorem 4.3 and Theorem 4.4.* Consider the sequence  $\{\hat{u}_n\}$  generated by the linearized multi-scale scheme (15) for all  $n \geq 2$  and where  $\hat{u}_1^H$  is defined using the nonlinear initialization (17). For the error analysis, we define in view of (47)

$$\theta_0^H = u_0^H - \mathcal{U}_0^H, \quad \hat{\theta}_n = \hat{u}_n - \hat{\mathcal{U}}_n, \quad \text{for } 1 \leq n \leq N. \quad (69)$$

Since the first step  $\hat{u}_1^H$  is defined using the nonlinear FE-HMM, the error propagation formula for the first step of the nonlinear scheme, see [6, Eq. (40)], yields

$$\|\theta_1^H\|_{L^2(\Omega)}^2 - \|\theta_0^H\|_{L^2(\Omega)}^2 + \lambda \Delta t \|\nabla \theta_1^H\|_{L^2(\Omega)}^2 \leq C \Delta t (\Delta t^2 + H^{2\mu} + r_{HMM}(\nabla \mathcal{U}_1^H)^2), \quad (70)$$

where  $\lambda$  is the monotonicity constant of  $\mathcal{A}^\varepsilon$ . In particular we observe that

$$\Delta t \|\nabla \theta_1^H\|_{L^2(\Omega)}^2 \leq C \Delta t (\Delta t^2 + H^{2\mu} + r_{HMM}(\nabla \mathcal{U}_1^H)^2) + C \|\theta_0^H\|_{L^2(\Omega)}^2. \quad (71)$$

From (54), we have that for  $2 \leq K \leq N$  (where  $K \geq 2$  instead of  $K \geq 1$  due to the initialization step)

$$\begin{aligned} \|\theta_{n+1}^H\|_{L^2(\Omega)}^2 - \|\theta_n^H\|_{L^2(\Omega)}^2 + \lambda_a \Delta t \|\nabla \hat{\theta}_{n+1}\|_{\mathcal{S}^{H,h}}^2 &\leq C \Delta t (\Delta t^2 + H^{2\mu} + r_{HMM}(\nabla \mathcal{U}_{n+1}^H)^2) \\ &\quad + \frac{2}{\lambda_a} \Delta t \mathcal{L}_n^2 \|\nabla \hat{\theta}_n\|_{\mathcal{S}^{H,h}}^2, \end{aligned} \quad (72)$$

where  $\lambda_a$  is the ellipticity constant of the tensor  $a^\varepsilon$  and  $\mathcal{L}_n$  is the linearization error defined in the proofs of Lemma 5.6 and Lemma 5.9. Analogously to the proof of Lemma 5.6, summing (72) from  $n = 1$  to  $n = N - 1$  and adding (70) yields

$$\begin{aligned} \max_{1 \leq n \leq N} \|\theta_n^H\|_{L^2(\Omega)}^2 + \lambda \Delta t \|\nabla \theta_1^H\|_{L^2(\Omega)}^2 + \left( \lambda_a - \frac{2}{\lambda_a} \max_{2 \leq n \leq N-1} \mathcal{L}_n^2 \right) \Delta t \sum_{n=2}^N \|\nabla \hat{\theta}_n\|_{\mathcal{S}^{H,h}}^2 &\quad (73) \\ \leq C \left( \Delta t^2 + H^{2\mu} + \max_{1 \leq n \leq N} r_{HMM}(\nabla \mathcal{U}_n^H)^2 \right) + \|\theta_0^H\|_{L^2(\Omega)}^2 + \frac{2}{\lambda_a} \Delta t \mathcal{L}_1^2 \|\nabla \hat{\theta}_1\|_{\mathcal{S}^{H,h}}^2. \end{aligned}$$

Further, analogously to the boundedness of  $\mathcal{L}_0$  in the proofs of Lemma 5.6 and Lemma 5.9, we deduce that  $\mathcal{L}_1$  is bounded. However, in contrast to case of a general initialization of the linearized multiscale method (see (55), where  $\|\nabla \hat{\theta}_0\|_{\mathcal{S}^{H,h}}$  cannot be estimated), we are now able to bound the term  $\|\nabla \hat{\theta}_1\|_{\mathcal{S}^{H,h}}$  explicitly. Let  $K \in \mathcal{T}_H$  and  $K_\delta$  be its associated sampling domain. Then, according to (69), we have that  $\theta_{K,1}^h = \hat{u}_{K,1}^h - \hat{\mathcal{U}}_{K,1}^h$  where  $\hat{u}_{K,1}^h$  and  $\hat{\mathcal{U}}_{K,1}^h$  is the solution to the nonlinear micro problem (11) constrained by  $u_1^H$  and  $\mathcal{U}_1^H$ , respectively. Thus, Lemma 5.2 yields

$$\|\nabla \hat{\theta}_{1,K}^h\|_{L^2(K_\delta)} \leq \frac{L}{\lambda} \sqrt{|K_\delta|} |\nabla u_1^H(x_K) - \nabla \mathcal{U}_1^H(x_K)|, \quad \text{i.e.,} \quad \|\nabla \hat{\theta}_1\|_{\mathcal{S}^{H,h}} \leq \frac{L}{\lambda} \|\nabla \theta_1^H\|_{L^2(\Omega)}. \quad (74)$$

Finally, by combining inequalities (71), (73) and (74) we obtain

$$\begin{aligned} \max_{1 \leq n \leq N} \|\theta_n^H\|_{L^2(\Omega)}^2 + \lambda \Delta t \|\nabla \theta_1^H\|_{L^2(\Omega)}^2 + \left( \lambda_a - \frac{2}{\lambda_a} \max_{2 \leq n \leq N-1} \mathcal{L}_n^2 \right) \Delta t \sum_{n=2}^N \|\nabla \hat{\theta}_n\|_{\mathcal{S}^{H,h}}^2 & \\ \leq C \left( \Delta t^2 + H^{2\mu} + \max_{1 \leq n \leq N} r_{HMM}(\nabla \mathcal{U}_n^H)^2 \right) + \|\theta_0^H\|_{L^2(\Omega)}^2 + C \Delta t \|\nabla \theta_1^H\|_{L^2(\Omega)}^2 & \\ \leq C \left( \Delta t^2 + H^{2\mu} + \max_{1 \leq n \leq N} r_{HMM}(\nabla \mathcal{U}_n^H)^2 \right) + C \|\theta_0^H\|_{L^2(\Omega)}^2. & \end{aligned}$$

Estimating the initialization error by combining  $\|\theta_0^H\|_{L^2(\Omega)} \leq \|u_0^H - g\|_{L^2(\Omega)} + \|g - \mathcal{U}_0^H\|_{L^2(\Omega)}$  with the error bound  $\|g - \mathcal{U}_0^H\|_{L^2(\Omega)} \leq CH^\mu$  concludes the proof.  $\square$

## 6 Numerical results

In this section, we compare the performances of the *nonlinear* multiscale method (10) whose upscaling procedure relies on nonlinear micro problems, and the *linearized* version (15) which is based on linear

micro problems. Using various sample problems in 2D, we show that this linearized version yields analogous numerical errors compared to the nonlinear version, but it is much faster because it avoids Newton iterations.

To measure the quality of the numerical solution  $\{u_n^H\}$ , we calculate the relative error measures  $e_{C^0(L^2)}$  and  $e_{L^2(H^1)}$  given by<sup>2</sup>

$$e_{C^0(L^2)} = \left( \max_{0 \leq n \leq N} \|u^{ref}(\cdot, t_n) - u_n^H\|_{L^2(\Omega)} \right) \left( \max_{0 \leq k \leq N_{ref}} \|u^{ref}(\cdot, t_k)\|_{L^2(\Omega)} \right)^{-1}, \quad (75a)$$

$$e_{L^2(H^1)} = \left( \Delta t \sum_{n=0}^N' \|\nabla u^{ref}(\cdot, t_n) - \nabla u_n^H\|_{L^2(\Omega)}^2 \right)^{1/2} \left( \Delta t \sum_{k=0}^{N_{ref}}' \|\nabla u^{ref}(\cdot, t_k)\|_{L^2(\Omega)}^2 \right)^{-1/2}, \quad (75b)$$

and  $u^{ref}$  denotes a reference solution for the homogenized equation (5). Since in general it is not available in analytical form, an algorithm to obtain an accurate reference solution  $u^{ref}$  is discussed below.

## 6.1 Convergence rates and performance comparisons

**Test problem.** We consider the square domain  $\Omega = (0, 1)^2$  and the final time  $T = 2$ . To investigate first the spatial (on macro and micro scale) and temporal discretization errors, we choose a test problem with a periodic tensor, which yields a modeling error  $e_{mod} = 0$  using a periodic coupling. We consider the multiscale problem (2) with the periodic tensor  $a^\varepsilon$  given by

$$a^\varepsilon(x, \xi) = a\left(\frac{x}{\varepsilon}, \xi\right) = a(y, \xi) = \left[ \frac{8}{5} + \frac{1}{3} \cdot \frac{1}{\left(\frac{1}{4} + |\xi|^2\right)^\gamma} \right] Id + \begin{pmatrix} \frac{\frac{9}{8} + \sin(2\pi y_1)}{\frac{9}{8} + \cos(2\pi y_2)} & 0 \\ 0 & \frac{\frac{9}{8} + \sin(2\pi y_2)}{\frac{9}{8} + \sin(2\pi y_1)} \end{pmatrix}, \quad (76)$$

for  $(x, \xi) \in \Omega \times \mathbb{R}^d$  and we choose  $\varepsilon = 10^{-4}$  as period of the micro oscillations. Throughout Section 6.1 we use  $\gamma = 1/2$  (with a single exception specified in the text). We note that (76) satisfies the assumption (30) of Theorem 4.4. Further, the right-hand side term in (2) is defined by

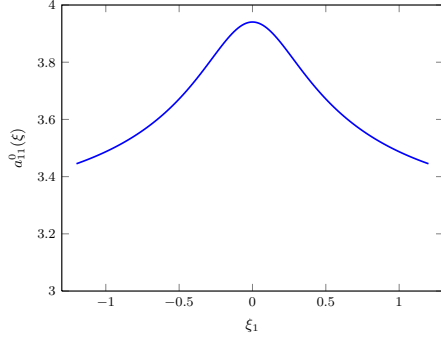
$$f(x, t) = \frac{1}{2}(1 + 2 \sin(2\pi x_1 t))(2 + 10x_2^2 + \cos(\pi t)), \quad (x, t) \in \Omega \times (0, 2). \quad (77)$$

Although our analysis is presented for a time independent  $f$ , it could be generalized straightforwardly to the case of a time dependent source term  $f$ .

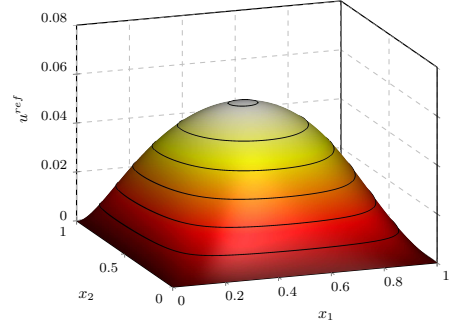
**Reference solution computation.** The reference solution  $u^{ref}$  is obtained by homogenizing the multiscale problem (2) following an iterative approach. First, we precompute the homogenized map  $\mathcal{A}^0(\xi)$  given by (6) for  $\xi$  within some bounded box  $Q \subset \mathbb{R}^d$ . We note that the homogenized map  $\mathcal{A}^0$  from (6) and the cell problems (7) are independent of the spatial variable  $x \in \Omega$  because the tensor (76) is periodic. The box  $Q$  has to be adjusted such that the gradient  $\nabla u^{ref}$  of the reference solution lies in  $Q$ . Within  $Q$  we choose uniformly distributed points  $\xi_i \in Q$  for which the nonlinear cell problems (7) are solved by a finite element method using piecewise affine basis functions. Using the numerical solutions to (7), an approximation of  $\mathcal{A}^0(\xi_i)$  is then calculated following the formula (6). For a general  $\xi \in Q$  we use bilinear interpolation of the values  $\mathcal{A}^0(\xi_i)$  at the uniformly distributed points  $\xi_i$ . In Figure 1.(a) the nonlinearity of homogenized map  $\mathcal{A}^0(\xi)$  is illustrated. For  $\xi = (\xi_1, 0)^T$  with  $-\frac{6}{5} \leq \xi_1 \leq \frac{6}{5}$ , we plot the entry  $a_{11}^0(\xi)$  of the homogenized tensor  $a^0(\xi)$  satisfying  $\mathcal{A}^0(\xi) = a^0(\xi)\xi$ .

Using this precomputed approximation of  $\mathcal{A}^0(\xi)$  we solve the effective equation (5) by combining the implicit Euler method in time and a finite element method (again with piecewise affine functions) in space. In Figure 1.(b)–(d) the reference solution  $u^{ref}$  at time  $t = 0, 1, 2$  is plotted. We note that the evolution of the local maxima of  $u^{ref}$  over time is mainly driven by the time-dependency of the right-hand side function  $f(x, t)$  while the nonlinearity of  $\mathcal{A}^0(\xi)$  leads to edge sharpening effects.

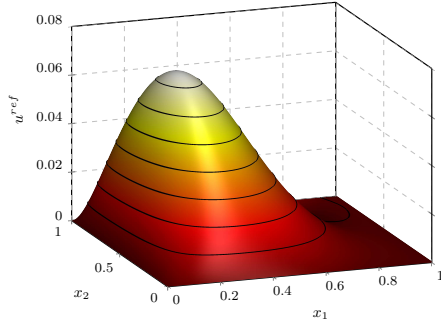
<sup>2</sup>The prime in  $\sum'$  indicates the use of the trapezoidal rule for the quadrature in time, i.e., the first and the last terms of the sums are multiplied by 1/2.



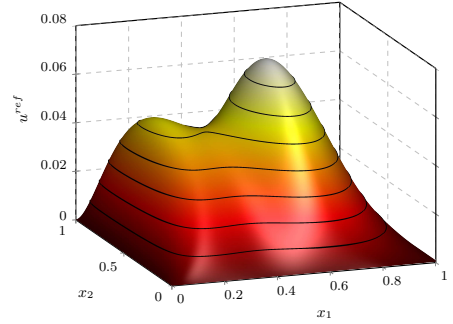
(a) Entry  $a_{11}^0(\xi_1, 0)$  of the homogenized tensor.



(b) Reference solution  $u^{ref}$  at  $t = 0$ .



(c) Reference solution  $u^{ref}$  at  $t = 1$ .



(d) Reference solution  $u^{ref}$  at  $t = 2$ .

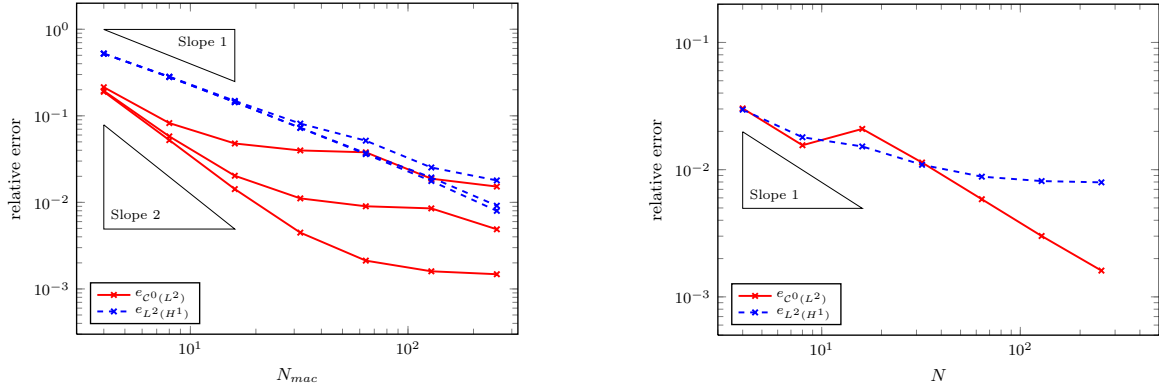
Figure 1: Reference solution  $u^{ref}$  for the homogenized solution  $u^0$  of the test problem of Section 6.1. Reference solution  $u^{ref}$  obtained as described in Section 6.1. Homogenized map approximated at  $601 \times 601$  uniformly distributed points  $\xi_i$  within the box  $Q = [-\frac{6}{5}, \frac{6}{5}]^2$ . Cell problems (7) solved on uniform triangular mesh with  $128^2$  degrees of freedom. Reference solution  $u^{ref}$  calculated at  $N_{ref} = 4096$  equidistant times on uniform triangulation of  $\Omega$  with  $512^2$  degrees of freedom.

**Initial conditions.** To avoid regularity issues for the initial condition, which are a classical issue already for linear parabolic singlescale problems, see [39, Chapter 3], we apply the following methodology. We calculate the reference solution  $u^{ref}$  on the extended time interval  $(-1/2, 2)$  with initial conditions at  $t = -1/2$  given by  $u^0(x, -1/2) = (x_1 - x_1^2)(x_2 - x_2^2)$ . Then, we use  $g(x) = u^{ref}(x, 0)$  as initial conditions for the test problem (2). Thus, the effects of incompatible or non-smooth initial data are negligible as the linearized multiscale scheme (15) is studied on  $(0, 2)$ , i.e., on a time interval safely bounded away from  $t = -1/2$ .

**Convergence rates.** We study the convergence of the linearized multiscale method (15) when solving the multiscale problem (2),(76),(77) with the initial condition at  $t = 0$  defined above. We perform the tests for the microscopic period  $\varepsilon = 10^{-4}$  and we choose the periodic coupling  $W(K_\delta) = W_{per}^1(K_\delta)$  and sampling domain size  $\delta = \varepsilon$  to obtain a vanishing modeling error  $e_{mod} = 0$ . For the discretization of the spatial macro domain  $\Omega$  and the sampling domains  $K_\delta$  we use uniform triangular meshes with  $N_{mac}$  and  $N_{mic}$  the number of elements in each spatial dimension, respectively. Further, we note that the mesh size of the macro and micro triangulations behave like  $H \sim N_{mac}^{-1}$  and  $h/\varepsilon \sim N_{mic}^{-1}$ , respectively.

First, we study the convergence with respect to the spatial discretizations. The influence of the time discretization is made negligible by choosing a fine time grid with  $N = 2000$  uniform time steps. The error measures (75) are plotted in Figure 2.(a) in dependence of  $N_{mac}$  while the micro discretizations are kept fixed with  $N_{mic} = 4, 8$  or  $16$ . We observe that the error measures (75) indicate a saturation of the error for fine macro discretizations (with some additional effects for  $N_{mic} = 4, 8$ ). However, the saturation levels clearly depend on the micro discretization  $N_{mic}$ . Thus, we conclude that for small macroscopic error, i.e., large  $N_{mac}$ , the microscopic error gets dominant as predicted by Theorem 4.6. We note that the micro error decreases superlinearly in  $h/\varepsilon$  (the rescaled micro mesh size). Further, the convergence rates with respect to the macro mesh size  $H \sim 1/N_{mac}$  are in coincidence with Theorem 4.6.

In Figure 2.(b), we take fine spatial macro and micro meshes with  $N_{mac} = 256$  and  $N_{mic} = 32$ , respectively, and analyze the dependence of the error measures (75) with respect to the time step size  $\Delta t \sim N^{-1}$ . While the error measure  $e_{C^0(L^2)}$  shows a linear convergence, the error measured by  $e_{L^2(H^1)}$  quickly approaches a constant value. Thus, despite the (relatively) fine spatial macro discretization the macroscopic error is still dominant (the micro error can be excluded as  $e_{C^0(L^2)}$  does not get saturated at a comparable level). In summary, the numerical tests presented in Figure 2 largely corroborate the fully discrete a priori bounds of Theorem 4.6.



(a) Space discretization error. The different lines correspond to a constant micro mesh  $N_{mic} = 4, 8, 16$ . Number of time steps  $N = 1024$ . Macro meshes with  $N_{mac} = 4, 8, 16, 32, 64, 128, 256$ .

(b) Time discretization error. Macro and micro space discretization with constant meshes  $N_{mac} = 256$ ,  $N_{mic} = 32$ . Number of time steps  $N = 4, 8, 16, 32, 64, 128, 256$ .

Figure 2: Convergence tests for linearized multiscale scheme applied to test problem of Section 6.1. Relative error measured by  $e_{C^0(L^2)}$  (solid line) and  $e_{L^2(H^1)}$  (dashed line) as a function of  $N_{mac}$  (in part (a)) and  $N$  (in part (b)), respectively. Comparison to the reference solution  $u^{ref}$  defined in Section 6.1.

**Refinement strategies for spatial discretization.** As proved in Theorem 4.6 and observed in Figure 2.(a) the spatial meshes have to be refined simultaneously to obtain an overall convergence of the spatial error. Therefore, optimal refinement strategies of the spatial meshes are essential to achieve an optimal computational cost, analogously to the linear case. Using  $H \sim N_{mac}^{-1}$  and  $h/\varepsilon \sim N_{mic}^{-1}$  (where  $N_{mac}$  and  $N_{mic}$  denote the number of elements in each spatial dimension for the macro and micro mesh, respectively) we recall the two  $L^2(H^1)$  and  $C^0(L^2)$  refinement strategies, which yield linear and quadratic error decays with respect to  $H$ , respectively ,

$$\begin{aligned} \text{error in } L^2(H^1) \text{ norm: } \quad H \sim \left(\frac{h}{\varepsilon}\right)^2 &\implies N_{mic} \sim \sqrt{N_{mac}} && \text{as } H^1 \text{ refinement strategy,} \\ \text{error in } C^0(L^2) \text{ norm: } \quad H^2 \sim \left(\frac{h}{\varepsilon}\right)^2 &\implies N_{mic} \sim N_{mac} && \text{as } L^2 \text{ refinement strategy.} \end{aligned} \quad (78)$$

**Study of the linearization error.** In view of Theorem 4.4 and Remark 4.5 it is important to study the error term  $e_{n,K}$  ( $0 \leq n \leq N$ ,  $K \in \mathcal{T}_H$ ) given in (31) by

$$e_{n,K}(x) = a^\varepsilon(x, \nabla \hat{u}_{n,K}^h) - \int_0^1 a^\varepsilon(x, \nabla \hat{u}_{n,K}^h - \tau \nabla \hat{\theta}_{n,K}^h) d\tau, \quad \text{a.e. } x \in K_\delta, \quad (79)$$

where  $\hat{u}_n \in \mathcal{S}^{H,h}$  is the approximation obtained by the linearized multiscale method (15) and  $\hat{\theta}_n \in \mathcal{S}^{H,h}$  denotes the difference  $\hat{\theta}_n = \hat{u}_n - \hat{U}_n \in \mathcal{S}^{H,h}$  between the numerical solution and an approximation of the exact solution  $u^0$  and its associated first order oscillations. In particular, we have that  $\hat{U}_n = (\mathcal{U}_n^H, \{\mathcal{U}_{n,K}^h\})$  where the macro function  $\mathcal{U}_n^H$  is an approximation of the homogenized solution  $u^0$  at time  $t_n$  (either the nodal interpolant  $\mathcal{I}_H u^0(\cdot, t_n)$  or the elliptic projection  $\tilde{u}_n^{H,0}$  defined in (41)) and  $\mathcal{U}_{n,K}^h$  is the solution of the nonlinear micro problem (11) constrained by  $\mathcal{U}_n^H$ . We choose  $\mathcal{U}_n^H = \mathcal{I}_H u^{ref}(\cdot, t_n)$  the nodal interpolant of the reference solution  $u^{ref}$ . The integral in (79) is evaluated using the Gauss quadrature formula with 10 nodes (to ensure a negligible quadrature error). In what follows, we study numerically the term

$\max \|e_{n,K}\|$  given by

$$\max \|e_{n,K}\| = \max_{\substack{K \in \mathcal{T}_H \\ 0 \leq n \leq N}} \|e_{n,K}\|_{(L^\infty(K_\delta))^{d \times d}}.$$

We apply the linearized multiscale method (15) to the test problem with tensor  $a^\varepsilon$  given in (76). For the spatial discretizations, we use the optimal simultaneous refinement of macro and micro grids derived in (78), denoted as  $H^1$  or  $L^2$  refinement. For the explicit choice of the parameters  $N_{mac}$  and  $N_{mic}$  we refer to Table 1.

	$N_{mac}$	4	8	16	32	64
$H^1$ refinement	$N_{mic}$	—	3	4	6	8
$L^2$ refinement	$N_{mic}$	4	8	16	32	64

Table 1: Discretization parameters for the refinement strategies of  $H^1$  and  $L^2$  refinement. The parameters  $N_{mac}$  and  $N_{mic}$  denote the number of elements in each spatial dimension when discretizing the macro domain  $\Omega$  and the sampling domains  $K_\delta$ , respectively, by uniform triangular meshes.

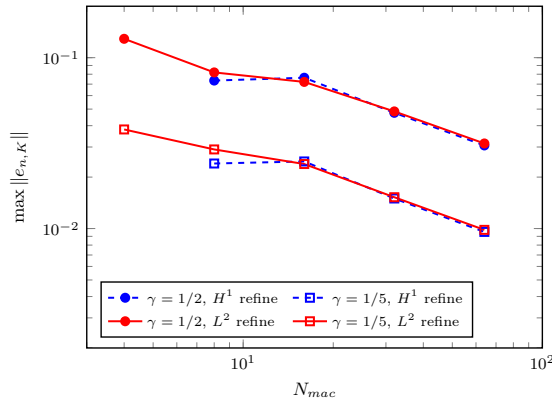


Figure 3: Tests for linearization error  $e_{n,K}$ , see (31), for test problem of Section 6.1 using tensor  $a^\varepsilon$  with  $\gamma = 1/2$  (circle marks) and  $\gamma = 1/5$  (square marks). Error term measured by  $\max \|e_{n,K}\|$  as a function of  $N_{mac}$ . Constant number of time steps  $N = 1024$ . Macro and micro meshes according to Table 1 using  $H^1$  refinement (dashed line) and  $L^2$  refinement (solid line), respectively.

In Figure 3 we plot the linearization error  $\max \|e_{n,K}\|$  under  $H^1$  and  $L^2$  refinement strategy in space for a given fine time grid with  $N = 1024$ . To gain insights on the linearization error  $\max \|e_{n,K}\|$  for different test problems we study the tensor  $a^\varepsilon$  given in (76) for  $\gamma = 1/2$  and  $\gamma = 1/5$ . First, we note in both cases that the linearization error converges with respect to  $N_{mac}$  (at a rate less than linear) and that the linearization error is comparable for both refinement strategies. Thus for this test setting, the linearization error is dictated by the spatial macro discretization. Further, we observe that the absolute value of the linearization error is (roughly) three times smaller for  $\gamma = 1/5$  than for  $\gamma = 1/2$ . This is reasonable as the tensor  $a^\varepsilon$  with  $\gamma = 1/5$  has a weaker nonlinearity than for  $\gamma = 1/2$  (smaller Lipschitz constant). In summary, the tests of Figure 3 suggest that the smallness assumption (32) for the linearization error is numerically satisfied for tensors  $a^\varepsilon$  given in (76) for  $0 \leq \gamma \leq 1/2$  if the spatial and temporal discretization parameters are fine enough.

**Performance comparisons.** The efficiency of the linearized multiscale method (15) compared to the nonlinear multiscale method (10) of [6] is the main feature of the proposed method, and is thus carried out carefully. The methods are implemented as similar as possible in MATLAB (version R2013b, 64-bit) and are run on one single thread of an Intel Xeon E5620 @2.4GHz CPU (with hyperthreading disabled in MATLAB). We apply both methods to the test problem of Section 6.1 for a given set of spatial and temporal discretizations. Further, for each set of parameters the reported CPU time  $\mathfrak{t}$  is obtained as the mean of the measured CPU time for 10 runs.

First, we consider the  $e_{L^2(H^1)}$  refinement strategy  $\Delta t \sim H \sim (h/\varepsilon)^2$  for a set of discretization parameters  $N$ ,  $N_{mac}$  and  $N_{mic}$ . We report the obtained errors and CPU time in Table 2 for both the nonlinear and the linearized method. As expected by our convergence analysis, we obtain analogous



errors for both methods. More interesting is the CPU time  $\mathfrak{t}$  which is smaller by about a factor 10 for the linearized method.

$N$	$N_{mac}$	$N_{mic}$	linearized method			nonlinear method		
			$\mathfrak{t}$	$e_{C^0(L^2)}$	$e_{L^2(H^1)}$	$\mathfrak{t}$	$e_{C^0(L^2)}$	$e_{L^2(H^1)}$
8	8	3	0.06	0.0560	0.2794	0.25	0.0563	0.2793
16	16	4	0.35	0.0618	0.1501	2.23	0.0559	0.1495
32	32	6	2.56	0.0189	0.0734	19.65	0.0138	0.0731
64	64	8	20.25	0.0131	0.0374	165.75	0.0111	0.0372

Table 2: Performance comparison between the linearized multiscale method and nonlinear multiscale method for test problem of Section 6.1. Simultaneous refinement of  $\Delta t$ ,  $H$  and  $h$  according to  $\Delta t \sim H \sim \sqrt{h/\varepsilon}$ . CPU time  $\mathfrak{t}$  measured in minutes. Error measures defined in (75).

Analogously, we consider an overall refinement of spatial and temporal discretizations such that  $e_{C^0(L^2)}$  converges quadratically with respect to  $H$ . Thus, we choose  $\Delta t \sim H^2$  and use the  $L^2$  refinement in space, see Table 1. In Table 3, we observe that the error measures  $e_{C^0(L^2)}$  and  $e_{L^2(H^1)}$  show quadratic and linear convergence, respectively. Further, both error measures indicate a comparable accuracy of the linearized and nonlinear method for the given set of parameters. However, the computational cost for the linearized scheme are again significantly smaller. We conclude, that for the test problem of this section, the linearized method needs 4-9 times less execution time.

$N$	$N_{mac}$	$N_{mic}$	linearized method			nonlinear method		
			$\mathfrak{t}$	$e_{C^0(L^2)}$	$e_{L^2(H^1)}$	$\mathfrak{t}$	$e_{C^0(L^2)}$	$e_{L^2(H^1)}$
4	4	4	0.01	0.2220	0.5042	0.04	0.2220	0.5044
16	8	8	0.11	0.0697	0.2807	0.74	0.0646	0.2805
64	16	16	2.07	0.0174	0.1440	19.95	0.0159	0.1439
256	32	32	76.59	0.0043	0.0724	789.70	0.0040	0.0724

Table 3: Performance comparison between the linearized multiscale method and nonlinear multiscale method for test problem of Section 6.1. Simultaneous refinement of  $\Delta t$ ,  $H$  and  $h$  according to  $\Delta t \sim H^2 \sim (h/\varepsilon)^2$ . CPU time  $\mathfrak{t}$  measured in minutes. Error measures defined in (75).

Finally, we perform a series of tests where we search parameters  $N$ ,  $N_{mac}$ ,  $N_{mic}$  for both linearized and nonlinear methods such that a given accuracy measured by  $e_{C^0(L^2)}$  is obtained at minimal computational cost. As set of possible parameters we take  $N \geq 2$ ,  $N_{mac} \in \{4, 8, 16, 32\}$  and  $N_{mic} = N_{mac}$  (according to  $L^2$  refinement in (78)). The results are given in Table 4. While the spatial parameters  $N_{mac}$  and  $N_{mic}$  are identical for both methods, the linearized scheme requires roughly twice as many timesteps to obtain a given precision. This is due to a slightly larger error constant  $C$  in the a priori estimate of Theorem 4.6 for the linearized method compared to the constant in the error estimates for the nonlinear method, see [6, Theorem 4.2]. We emphasize that this factor is independent of the spatial discretizations. However, this still leads to computational savings of a factor 3-6 for the linearized scheme. Thus, for the test problem studied in this section, the linearized multiscale method (15) indeed is drastically more efficient than the nonlinear multiscale scheme (10).

precision	linearized method					nonlinear method				
	$N$	$N_{mac}$	$N_{mic}$	$\mathfrak{t}$	$e_{C^0(L^2)}$	$N$	$N_{mac}$	$N_{mic}$	$\mathfrak{t}$	$e_{C^0(L^2)}$
0.1000	4	8	8	0.06	0.0792	4	8	8	0.19	0.0792
0.0750	6	8	8	0.07	0.0722	6	8	8	0.29	0.0722
0.0500	4	16	16	0.46	0.0410	4	16	16	1.48	0.0410
0.0250	26	16	16	1.06	0.0248	12	16	16	4.34	0.0240
0.0100	49	32	32	18.86	0.0099	21	32	32	91.21	0.0100
0.0075	73	32	32	25.55	0.0075	34	32	32	139.75	0.0074
0.0050	157	32	32	48.88	0.0050	87	32	32	293.01	0.0049

Table 4: Performance comparison between the linearized multiscale method and nonlinear multiscale method for test problem of Section 6.1. Given precision (measured in  $e_{C^0(L^2)}$ ) attained at optimal computational cost. CPU time  $\mathfrak{t}$  measured in minutes. Error measures defined in (75).

## 6.2 Case of a degenerated problem

Many physical applications, e.g., non-Newtonian fluids, problems in elasticity and magnetodynamics, are modeled as monotone nonlinear parabolic problems (2) with a tensor  $a^\varepsilon(x, \xi)$  degenerated in  $\xi \in \mathbb{R}^d$  (typically  $\|a^\varepsilon(x, \xi)\|_{\mathcal{F}} \rightarrow 0$  or  $\infty$  for either  $\xi \rightarrow 0$  or  $|\xi| \rightarrow \infty$ ). A widely studied example is the  $p$ -Laplacian, on which we now focus in a multiscale context. Such degenerated parabolic problems are particularly challenging numerically and for the analysis, due to the poor regularity of the exact solutions, see e.g. [10, 14]. However, the homogenization results of [36, 37] cited in Section 2 (for monotone operators on  $H^1(\Omega)$ ) hold as well for monotone operators on  $W^{1,p}(\Omega)$  for  $p \geq 2$ , e.g., for operators with nonlinearities similar to the  $p$ -Laplacian. Hence, to study the applicability of the linearized numerical homogenization method (15) for homogenization problems (2) with a degenerated multiscale tensor  $a^\varepsilon(x, \xi)$  we consider the problem of a multiscale  $p$ -Laplacian on the space time domain  $\Omega \times (0, T) = (0, 1)^2 \times (0, 1/2)$ . For  $p > 2$ , we introduce the periodic tensor  $a^\varepsilon(x, \xi)$  given by

$$a^\varepsilon(x, \xi) = a\left(\frac{x}{\varepsilon}, \xi\right) = a(y, \xi) = \left( \frac{11}{10} + \sin(2\pi(x_1 + x_2)) + \frac{\frac{9}{8} + \sin(2\pi y_1)}{\frac{9}{8} + \cos(2\pi y_2)} + \frac{\frac{9}{8} + \sin(2\pi y_2)}{\frac{9}{8} + \sin(2\pi y_1)} \right) |\xi|^{p-2} Id, \quad (80)$$

which is equal to the zero matrix for  $\xi = 0$  and unbounded for  $|\xi| \rightarrow \infty$ . In this section, we consider the tensor  $a^\varepsilon$  for  $p = 3$ . Further, we choose  $\varepsilon = 10^{-4}$ , the right-hand side function  $f \equiv 1$  and the initial condition  $u^\varepsilon(x, 0) = \frac{1}{2}x_2(1 - x_2)\cos(\pi x_1)$ . We employ mixed boundary conditions on the spatial boundary  $\partial\Omega$

$$u^\varepsilon(x, t) = 0, \quad \text{on } \Gamma_D \times (0, \frac{1}{2}), \quad a^\varepsilon(x, \nabla u^\varepsilon(x, t)) \nabla u^\varepsilon(x, t) \cdot n = 0 \quad \text{on } \Gamma_N \times (0, \frac{1}{2}),$$

where  $\Gamma_D = [0, 1] \times \{x_2 = 0, 1\}$ ,  $\Gamma_N = \partial\Omega \setminus \Gamma_D$  and  $n$  is the outer normal vector.

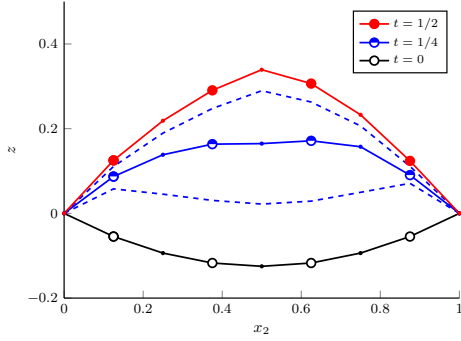
**Numerical studies.** We apply the linearized multiscale method (15) to the degenerated parabolic multiscale problem with the tensor  $a^\varepsilon(x, \xi)$  given in (80) for  $p = 3$ . The solutions obtained by the linearized scheme are compared to a numerical solution computed by using the nonlinear multiscale method (10). For both methods, we choose an optimal coupling of macro and micro solvers, i.e., periodic coupling  $W(K_\delta) = W_{per}^1(K_\delta)$  and sampling domain size  $\delta = \varepsilon$ . Further, to avoid singular linear systems due to degenerated (linear and nonlinear) micro problems we regularize the tensor (80) by replacing  $|\xi|$  by  $\sqrt{|\xi|^2 + \eta}$  with  $\eta = 10^{-10}$ .

The spatial points  $x \in \Omega$  where  $\nabla u^0(x, t) = 0$  for some  $t \in (0, 1/2)$  are of particular interest due to the degeneracy of the tensor  $a^\varepsilon(x, \xi)$  in  $\xi = 0$ . As  $\Gamma_D \cap \{u^\varepsilon(x, 0) < 0\}$  is a set of points where degeneracy occurs, we present the profiles of the numerical solutions at  $x_1 = 1$  at times  $t = 0, 1/8, 1/4, 3/8, 1/2$ . The plot of Figure 4.(d) shows the numerical solution given by the nonlinear multiscale scheme calculated with  $N = 256$  time steps and  $N_{mac} = N_{mic} = 64$  elements in each spatial dimension of the uniform macro and micro mesh. Then, in Figure 4.(a-c) the solutions obtained by the linearized scheme for parameters  $N = N_{mac} = N_{mic} = 8, 16, 64$  are presented. We recall that in view of Theorem 4.6, a simultaneous refinement of temporal and spatial discretization parameters is needed to obtain robust convergence at optimal computational cost.

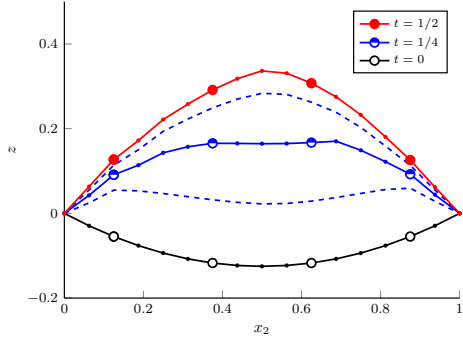
We observe that the numerical solutions obtained by the linearized scheme approximate well the reference solution calculated by the nonlinear scheme. Further, we note that in Figure 4 small oscillations can be noticed for values of  $x_2$  close to the boundary. Thus, at the points (in space) where the tensor  $a^\varepsilon$  degenerates, stability issues may appear. However, these become small when appropriately refining the temporal and spatial discretization.

**Dirichlet coupling.** As for practical problems the exact period  $\varepsilon$  of the spatial micro oscillations often cannot be determined exactly, one might use coupling conditions with non-periodic boundary conditions for micro sampling. A popular choice is to use Dirichlet boundary conditions  $W(K_\delta) = H_0^1(K_\delta)$  for the micro problems and a sampling domain size  $\delta$  larger than the actual period  $\varepsilon$ . Herein, we present numerical results illustrating how the modeling error due to those non-optimal coupling conditions behaves when the sampling domain size  $\delta$  is increased.

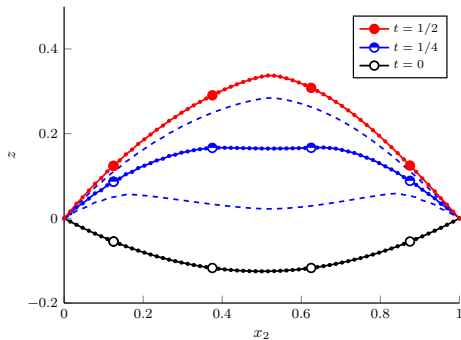
For all tests of this paragraph, we use the linearized multiscale method (15), take a fixed number of time steps  $N = 160$  and use a spatial macro mesh with  $N_{mac} = 32$  (the number elements in each spatial dimension). First, we compute a reference solution  $\{\hat{u}_n^{per}\}$  by using optimal periodic coupling, i.e.,  $W(K_\delta) = W_{per}^1(K_\delta)$  and sampling domain size  $\delta = \varepsilon$ . For the micro discretization, we use  $N_{mic} = 32$  elements in each spatial dimension. We emphasize that this solution  $\{\hat{u}_n^{per}\}$  is free of any modeling



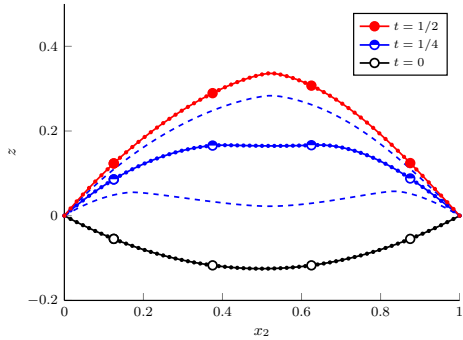
(a) Linearized scheme.  $N_{mac} = N_{mic} = N = 8$ .



(b) Linearized scheme.  $N_{mac} = N_{mic} = N = 16$ .



(c) Linearized scheme.  $N_{mac} = N_{mic} = N = 64$ .



(d) Nonlinear scheme as reference.  $N_{mac} = N_{mic} = 64$  and  $N = 256$ .

Figure 4: Degenerated parabolic multiscale problem of Section 6.2. Numerical solutions obtained by linearized multiscale method (15) and nonlinear multiscale method (10), respectively. Profiles of numerical solution as a function of  $x_2$  at  $x_1 = 1$  for times  $t = 0, 1/8, 1/4, 3/8, 1/2$ . Simultaneous refinement of temporal and spatial discretization for linearized scheme. To facilitate comparisons, the bullets indicate the solutions for  $x_2 = 1/8, 3/8, 5/8, 7/8$  and  $t = 0, 1/4, 1/2$ .

error, i.e., satisfies the estimates of Theorem 4.6 with  $e_{mod} = 0$ . Then, we apply the linearized multiscale scheme (15) with Dirichlet boundary conditions  $W(K_\delta) = H_0^1(K_\delta)$  and sampling domain size  $\delta = 2^k \log(3)\varepsilon$  for  $k = 0, \dots, 4$  (solutions denoted by  $\{\hat{u}_n^{\delta,k}\}$ ). The micro domain discretization is adapted to the sampling domain size  $\delta$  such that the micro mesh size  $h$  is constant, i.e., the micro error is constant. In particular we take  $N_{mic} = 35, 70, 141, 281, 562$ .

As the time step size  $\Delta t$  as well as the macro and micro mesh sizes  $H, h$  used for the solutions  $\{\hat{u}_n^{\delta,k}\}$  (with Dirichlet coupling) and the reference solution  $\{\hat{u}_n^{per}\}$  (with periodic coupling) are identical, the difference  $\hat{u}_n^{\delta,k} - \hat{u}_n^{per}$  is solely due to the modeling error. In Figure 5, we compare the numerical solutions  $\{\hat{u}_n^{\delta,k}\}$  to the reference solution  $\{\hat{u}_n^{per}\}$  using the relative error measures (75).

In Figure 5 we observe a convergence of  $e_{C^0(L^2)}$  of linear order  $\mathcal{O}(\varepsilon/\delta)$ . A similar trend can be identified for  $e_{L^2(H^1)}$ . This suggests that the modeling error for the studied nonlinear and degenerated test problem behaves like for linear homogenization problems, see [19]. Further, the estimate from Theorem 4.6 predicting a convergence of order  $\mathcal{O}(\sqrt{\varepsilon/\delta})$  (for non-degenerated tensors  $a^\varepsilon$ ) seems to be non-optimal for the studied test problem.

## 7 Conclusion

We presented a new linearized multiscale method to solve a class of nonlinear monotone parabolic homogenization problems and we derived fully discrete a priori error estimates (in time and space). The assumptions for the convergence results are twofold. Either we make a smallness assumption for the strength of the nonlinearity, or we suppose that the linearization error itself is sufficiently small. Numerical results show that the linearization error is indeed small for sufficiently fine discretizations in time

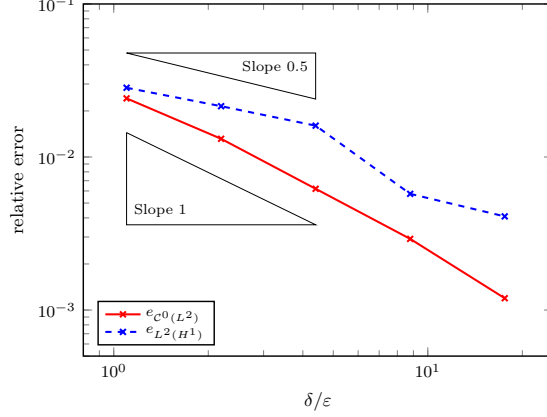


Figure 5: Degenerated parabolic multiscale problem of Section 6.2. Effect of sampling domain size  $\delta$  for linearized multiscale method (15) using Dirichlet coupling  $W(K_\delta) = H_0^1(K_\delta)$ . Sampling domain size  $\delta$  taken as  $\delta = 2^k \log(3)\varepsilon$  for  $k = 0, \dots, 4$ . Temporal and spatial macro and micro discretization errors kept constant. Comparison to solution obtained by the linearized multiscale scheme (15) with optimal periodic coupling.

and space. The main feature of the proposed approach is that the upscaling strategy is based only on *linear* micro problems, which makes the implementation of the method efficient and straightforward, as demonstrated in the numerical experiments.

Since the nonlinearity of the studied problem possibly leads to a low regularity of the exact solution, a combination of the proposed method with adaptivity techniques in time and space would be of practical interest. In view of the recent work [22], where a posteriori estimates for linearization errors in nonlinear solvers have been derived, one might aim to control the linearization error  $e_{n,K}$ , see (31), by some a posteriori error indicators. Thus, the hypothesis (28) in Theorem 4.4, assuming that  $e_{n,K}$  is small enough, could possibly be ensured by using an appropriate adaptive refinement.

**Acknowledgements.** The research of A. A., M. H., and G. V. is partially supported by the Swiss National Foundation, Grant: No 200021\_134716/1, and No 200020\_144313/1, respectively.

## A Appendix

*Proof of Remark 2.1.* For  $\xi \in \mathbb{R}^d$  and a.e.  $x \in \Omega$ , the derivative  $D_\xi \mathcal{A}^\varepsilon$  can be represented by

$$D_\xi \mathcal{A}^\varepsilon(x, \xi) = a^\varepsilon(x, \xi) + \partial_\xi a^\varepsilon(x, \xi)[\xi], \quad \text{with } (\partial_\xi a^\varepsilon(x, \xi)[\xi])_{ik} = \sum_{j=1}^d \frac{\partial a_{ij}^\varepsilon}{\partial \xi_k}(x, \xi) \xi_j, \quad 1 \leq i, k \leq d. \quad (81)$$

Using the Cauchy-Schwarz inequality and condition (4), we derive that (a.e.  $x \in \Omega$ ,  $\xi \in \mathbb{R}^d$ )

$$\|\partial_\xi a^\varepsilon(x, \xi)[\xi]\|_{\mathcal{F}}^2 \leq \sum_{i,j,k=1}^d \left| \frac{\partial a_{ij}^\varepsilon}{\partial \xi_k}(x, \xi) \right|^2 |\xi|^2 \leq L_a^2 \frac{|\xi|^2}{(1+|\xi|)^2} \leq L_a^2, \quad \text{i.e., } \|\partial_\xi a^\varepsilon(x, \xi)[\xi]\|_{\mathcal{F}} \leq L_a. \quad (82)$$

First, we show that  $\mathcal{A}^\varepsilon$  satisfies  $(\mathcal{A}_1)$ . For a.e.  $x \in \Omega$ ,  $\xi_1, \xi_2 \in \mathbb{R}^d$ , the representation (81) yields

$$\begin{aligned} \mathcal{A}^\varepsilon(x, \xi_1) - \mathcal{A}^\varepsilon(x, \xi_2) &= \int_0^1 D_\xi \mathcal{A}^\varepsilon(x, \xi_2 + t(\xi_1 - \xi_2))(\xi_1 - \xi_2) dt \\ &= \int_0^1 a^\varepsilon(x, \xi_2 + t(\xi_1 - \xi_2)) + \partial_\xi a^\varepsilon(x, \xi_2 + t(\xi_1 - \xi_2))[\xi_2 + t(\xi_1 - \xi_2)] dt (\xi_1 - \xi_2) \\ &\leq (\Lambda_a + L_a) |\xi_1 - \xi_2|, \end{aligned}$$

where we used the boundedness (3) of  $a^\varepsilon$  and bound (82) for  $\partial_\xi a^\varepsilon$ . Thus, the map  $\mathcal{A}^\varepsilon$  satisfies  $(\mathcal{A}_1)$ . Similarly, using the ellipticity of  $a^\varepsilon$  stated in (3) we obtain that (for a.e.  $x \in \Omega$ )

$$[\mathcal{A}^\varepsilon(x, \xi_1) - \mathcal{A}^\varepsilon(x, \xi_2)] \cdot (\xi_1 - \xi_2) \geq (\lambda_a - L_a) |\xi_1 - \xi_2|^2,$$

i.e., the map  $\mathcal{A}^\varepsilon$  is indeed strongly monotone if  $L_a < \lambda_a$ .  $\square$

## References

- [1] A. ABDULLE, *On a priori error analysis of fully discrete heterogeneous multiscale FEM*, Multiscale Model. Simul., 4 (2005), pp. 447–459.
- [2] ———, *The finite element heterogeneous multiscale method: a computational strategy for multiscale PDEs*, in Multiple scales problems in biomathematics, mechanics, physics and numerics, vol. 31 of GAKUTO Internat. Ser. Math. Sci. Appl., Gakkōtoshō, Tokyo, 2009, pp. 133–181.
- [3] ———, *A priori and a posteriori error analysis for numerical homogenization: a unified framework*, Ser. Contemp. Appl. Math. CAM, 16 (2011), pp. 280–305.
- [4] ———, *Discontinuous Galerkin finite element heterogeneous multiscale method for elliptic problems with multiple scales*, Math. Comp., 81 (2012), pp. 687–713.
- [5] A. ABDULLE, W. E, B. ENGQUIST, AND E. VANDEN-EIJNDEN, *The heterogeneous multiscale method*, Acta Numer., 21 (2012), pp. 1–87.
- [6] A. ABDULLE AND M. E. HUBER, *Finite element heterogeneous multiscale method for nonlinear monotone parabolic homogenization problems*. MATHICSE Technical Report July 2014, École Polytechnique Fédérale de Lausanne.
- [7] A. ABDULLE AND G. VILMART, *A priori error estimates for finite element methods with numerical quadrature for nonmonotone nonlinear elliptic problems*, Numer. Math., 121 (2012), pp. 397–431.
- [8] ———, *Analysis of the finite element heterogeneous multiscale method for quasilinear elliptic homogenization problems*, Math. Comp., 83 (2014), pp. 513–536.
- [9] I. BABUŠKA, *Solution of interface problems by homogenization. III*, SIAM J. Math. Anal., 8 (1977), pp. 923–937.
- [10] J. W. BARRETT AND W. B. LIU, *Finite element approximation of the parabolic  $p$ -Laplacian*, SIAM J. Numer. Anal., 31 (1994), pp. 413–428.
- [11] S. C. BRENNER AND L. R. SCOTT, *The mathematical theory of finite element methods*, vol. 15 of Texts in Applied Mathematics, Springer, New York, third ed., 2008.
- [12] P. G. CIARLET, *The finite element method for elliptic problems*, vol. 4 of Studies in Mathematics and its Applications, North-Holland, 1978.
- [13] J. E. DENDY, JR., *Galerkin’s method for some highly nonlinear problems*, SIAM J. Numer. Anal., 14 (1977), pp. 327–347.
- [14] L. DIENING, C. EBMAYER, AND M. RŮŽIČKA, *Optimal convergence for the implicit space-time discretization of parabolic systems with  $p$ -structure*, SIAM J. Numer. Anal., 45 (2007), pp. 457–472 (electronic).
- [15] M. DOBROWOLSKI,  *$L^\infty$ -convergence of linear finite element approximation to nonlinear parabolic problems*, SIAM J. Numer. Anal., 17 (1980), pp. 663–674.
- [16] J. DOUGLAS, JR. AND T. DUPONT, *Galerkin methods for parabolic equations*, SIAM J. Numer. Anal., 7 (1970), pp. 575–626.
- [17] R. DU AND P. MING, *Heterogeneous multiscale finite element method with novel numerical integration schemes*, Commun. Math. Sci., 8 (2010), pp. 863–885.
- [18] W. E AND B. ENGQUIST, *The heterogeneous multiscale methods*, Commun. Math. Sci., 1 (2003), pp. 87–132.
- [19] W. E, P. MING, AND P. ZHANG, *Analysis of the heterogeneous multiscale method for elliptic homogenization problems*, J. Amer. Math. Soc., 18 (2005), pp. 121–156.
- [20] Y. EFENDIEV AND T. Y. HOU, *Multiscale finite element methods. Theory and applications*, vol. 4 of Surveys and Tutorials in the Applied Mathematical Sciences, Springer, New York, 2009.

- [21] Y. EFENDIEV AND A. PANKOV, *Numerical homogenization of nonlinear random parabolic operators*, Multiscale Model. Simul., 2 (2004), pp. 237–268.
- [22] A. ERN AND M. VOHRALÍK, *Adaptive inexact Newton methods with a posteriori stopping criteria for nonlinear diffusion PDEs*, SIAM J. Sci. Comput., 35 (2013), pp. A1761–A1791.
- [23] M. FEISTAUER AND A. ŽENÍŠEK, *Finite element solution of nonlinear elliptic problems*, Numer. Math., 50 (1987), pp. 451–475.
- [24] J. FREHSE AND R. RANNACHER, *Asymptotic  $L^\infty$ -error estimates for linear finite element approximations of quasilinear boundary value problems*, SIAM J. Numer. Anal., 15 (1978), pp. 418–431.
- [25] E. HAIRER AND G. WANNER, *Solving ordinary differential equations II. Stiff and differential-algebraic problems*, Springer-Verlag, Berlin and Heidelberg, 1996.
- [26] P. HENNING AND M. OHLBERGER, *A Newton-scheme framework for multiscale methods for nonlinear elliptic homogenization problems*, in Proceedings of the ALGORITMY 2012, 19th Conference on Scientific Computing, Vysoké Tatry, Podbanské, 2012, pp. 65–74.
- [27] V. H. HOANG, *Sparse finite element method for periodic multiscale nonlinear monotone problems*, Multiscale Model. Simul., 7 (2008), pp. 1042–1072.
- [28] P. HOUSTON, J. ROBSON, AND E. SÜLI, *Discontinuous Galerkin finite element approximation of quasilinear elliptic boundary value problems. I. The scalar case*, IMA J. Numer. Anal., 25 (2005), pp. 726–749.
- [29] O. A. LADYZHENSKAYA AND N. N. URAL'TSEVA, *Linear and quasilinear elliptic equations*, Translated from the Russian by Scripta Technica, Inc. Translation editor: Leon Ehrenpreis, Academic Press, New York-London, 1968.
- [30] W. B. LIU AND J. W. BARRETT, *Finite element approximation of some degenerate monotone quasilinear elliptic systems*, SIAM J. Numer. Anal., 33 (1996), pp. 88–106.
- [31] C. LUBICH AND A. OSTERMANN, *Linearly implicit time discretization of non-linear parabolic equations*, IMA J. Numer. Anal., 15 (1995), pp. 555–583.
- [32] C. G. MAKRIDAKIS, *Finite element approximations of nonlinear elastic waves*, Math. Comp., 61 (1993), pp. 569–594.
- [33] Y.-Y. NIE AND V. THOMÉE, *A lumped mass finite-element method with quadrature for a non-linear parabolic problem*, IMA J. Numer. Anal., 5 (1985), pp. 371–396.
- [34] I. NIYONZIMA, R. V. SABARIEGO, P. DULAR, AND C. GEUZAINÉ, *Finite element computational homogenization of nonlinear multiscale materials in magnetostatics*, IEEE Transactions on Magnetics, 48 (2012), pp. 587–590.
- [35] I. NIYONZIMA, R. V. SABARIEGO, P. DULAR, F. HENROTTE, AND C. GEUZAINÉ, *Computational homogenization for laminated ferromagnetic cores in magnetodynamics*, IEEE Transactions on Magnetics, 49 (2013), pp. 2049–2052.
- [36] A. PANKOV, *G-convergence and homogenization of nonlinear partial differential operators*, vol. 422 of Mathematics and its Applications, Kluwer Academic Publishers, Dordrecht, 1997.
- [37] N. SVANSTEDT, *G-convergence of parabolic operators*, Nonlinear Anal., 36 (1999), pp. 807–842.
- [38] N. SVANSTEDT, N. WELLANDER, AND J. WYLLER, *A numerical algorithm for nonlinear parabolic equations with highly oscillating coefficients*, Numer. Methods Partial Differential Equations, 12 (1996), pp. 423–440.
- [39] V. THOMÉE, *Galerkin finite element methods for parabolic problems*, vol. 25 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, second ed., 2006.
- [40] E. ZEIDLER, *Nonlinear functional analysis and its applications. II/B*, Springer-Verlag, New York, 1990. Nonlinear monotone operators, Translated from the German by the author and Leo F. Boron.