

AN ITERATIVE METHOD FOR ELLIPTIC PROBLEMS WITH RAPIDLY OSCILLATING COEFFICIENTS

S. ARMSTRONG, A. HANNUKAINEN, T. KUUSI, AND J.-C. MOURRAT

ABSTRACT. We introduce a new iterative method for computing solutions of elliptic equations with random rapidly oscillating coefficients. Similarly to a multigrid method, each step of the iteration involves different computations meant to address different length scales. However, we use here the homogenized equation on all scales larger than a fixed multiple of the scale of oscillation of the coefficients. While the performance of standard multigrid methods degrades rapidly under the regime of large scale separation that we consider here, we show an explicit estimate on the contraction factor of our method which is independent of the size of the domain. We also present numerical experiments which confirm the effectiveness of the method.

1. INTRODUCTION

1.1. Informal summary of results. In this paper, we introduce a new iterative method for the numerical approximation of solutions of elliptic problems with rapidly oscillating coefficients. For definiteness, we consider the Dirichlet problem

$$(1.1) \quad \begin{cases} -\nabla \cdot (\mathbf{a}(x)\nabla u) = f & \text{in } U_r, \\ u = w & \text{on } \partial U_r, \end{cases}$$

where $r > 0$ is the length scale of the problem, which is typically very large ($r \gg 1$), and we write $U_r := rU$ where $U \subseteq \mathbb{R}^d$ is a bounded $C^{1,1}$ domain, in dimension $d \geq 2$. The boundary condition w belongs to $H^1(U_r)$, and the right-hand side f belongs to $H^{-1}(U_r)$. The coefficients $\mathbf{a}(x)$ are symmetric, uniformly elliptic and Hölder continuous. Moreover, in order to ensure that *quantitative homogenization* holds on large scales, we assume that the coefficients are sampled by a probability measure which is \mathbb{Z}^d -stationary and has a unit range of dependence (see below for the precise formulation of these assumptions). Our goal is to build a numerical method for the computation of u which remains efficient in the regime of fast oscillations of the coefficient field (which in our setting corresponds to the case in which the length scale is very large, $r \gg 1$) and does not rely on scale separation for convergence (the method computes the true solution for fixed r and not only in the limit $r \rightarrow \infty$).

In the absence of fast oscillations of the coefficient field, contemporary technology allows to access numerical approximations of elliptic problems involving billions of degrees of freedom. One of the most successful methods allowing to achieve such results is the *multigrid method* (see [15] for benchmarks). However, the performance of this method degrades as the coefficient field becomes more rapidly oscillating (see for instance [32, Table IV]).

Date: July 16, 2019.

2010 Mathematics Subject Classification. 65N55, 35B27.

We seek to remedy this problem by leveraging on *homogenization*. While standard multigrid methods use a decomposition of the elliptic problem into a series of scales, the difficulty in our context is that the slow eigenmodes of the heterogeneous operator still have fast oscillations, and are thus not easily captured through a coarse representation. We overcome this by introducing a suitable variant of the multigrid method that succeeds in replacing the heterogeneous operator by the homogenized one on length scales larger than a large but finite multiple of the correlation length scale. The result is a new iterative method that converges exponentially fast in the number of iterations, each of which is relatively inexpensive to compute—the memory and number of computations required scale linearly in the volume, and the computation is very amenable to parallelization. We give a rigorous proof of convergence and present numerical experiments which establish the efficiency of the method from a practical point of view.

1.2. Statement of the main result. We introduce some notation in order to state our main result. We begin with the precise assumptions on the coefficient field. We fix parameters $\Lambda > 1$ and $\alpha \in (0, 1]$ and require our coefficient fields $\mathbf{a}(x)$ to satisfy

$$(1.2) \quad \forall x, y \in \mathbb{R}^d, \quad |\mathbf{a}(y) - \mathbf{a}(x)| \leq \Lambda |x - y|^\alpha$$

and

$$(1.3) \quad \forall x \in \mathbb{R}^d, \quad \forall \xi \in \mathbb{R}^d, \quad \Lambda^{-1} |\xi|^2 \leq \xi \cdot \mathbf{a}(x) \xi \leq \Lambda |\xi|^2.$$

We denote by $\mathbb{R}_{\text{sym}}^{d \times d}$ the set of d -by- d real symmetric matrices and define

$$\Omega := \left\{ \mathbf{a} : \mathbb{R}^d \rightarrow \mathbb{R}_{\text{sym}}^{d \times d} \text{ satisfying (1.2) and (1.3)} \right\}.$$

For each Borel set $V \subseteq \mathbb{R}^d$, we denote by \mathcal{F}_V the Borel σ -algebra on Ω generated by the family of mappings

$$\mathbf{a} \mapsto \int_{\mathbb{R}^d} \chi \mathbf{a}_{ij}, \quad i, j \in \{1, \dots, d\}, \quad \chi \in C_c^\infty(V).$$

We also set $\mathcal{F} := \mathcal{F}_{\mathbb{R}^d}$. For each $y \in \mathbb{R}^d$, we denote by $T_y : \Omega \rightarrow \Omega$ the action of translation by y :

$$\forall x \in \mathbb{R}^d, \quad T_y \mathbf{a}(x) := \mathbf{a}(x + y).$$

We assume that \mathbb{P} is a probability measure on (Ω, \mathcal{F}) satisfying:

- stationarity with respect to \mathbb{Z}^d -translations: for every $y \in \mathbb{Z}^d$ and $A \in \mathcal{F}$,

$$\mathbb{P}[T_y A] = \mathbb{P}[A];$$

- unit range of dependence: for every Borel sets $V, W \subseteq \mathbb{R}^d$,

$$\text{dist}(V, W) \geq 1 \implies \mathcal{F}_V \text{ and } \mathcal{F}_W \text{ are } \mathbb{P}\text{-independent.}$$

The expectation associated with the probability measure \mathbb{P} is denoted by \mathbb{E} . We recall that, by classical homogenization theory (see [30, 3]), the heterogeneous operator $-\nabla \cdot \mathbf{a}(x) \nabla$ homogenizes to the homogeneous operator $-\nabla \cdot \bar{\mathbf{a}} \nabla$, where $\bar{\mathbf{a}} \in \mathbb{R}^{d \times d}$ is a deterministic, constant, positive definite matrix. For every $s, \theta > 0$ and random variable X , we write

$$(1.4) \quad X \leq \mathcal{O}_s(\theta) \quad \text{if and only if} \quad \mathbb{E} \left[\exp \left((\theta^{-1} \max(X, 0))^s \right) \right] \leq 2.$$

We also set, for every $\lambda \in (0, 1]$,

$$(1.5) \quad \ell(\lambda) := \begin{cases} (\log(1 + \lambda^{-1}))^{\frac{1}{2}} & \text{if } d = 2, \\ 1 & \text{if } d \geq 3. \end{cases}$$

For notational convenience, from now on we will suppress the explicit dependence on the spatial variable in the operator $-\nabla \cdot \mathbf{a}(x)\nabla$ and simply write $-\nabla \cdot \mathbf{a}\nabla$.

We now state the main result of the paper. We recall that \mathbb{P} is a probability measure on (Ω, \mathcal{F}) which specifies the law of the coefficient field $\mathbf{a}(x)$ and satisfies the assumptions stated above, that $\bar{\mathbf{a}}$ is the homogenized matrix associated to \mathbb{P} , and that $U \subseteq \mathbb{R}^d$ is a bounded domain with $C^{1,1}$ boundary.

Theorem 1.1 (*H^1 contraction*). *For each $s \in (0, 2)$, there exists a constant $C(s, U, \Lambda, \alpha, d) < \infty$ such that the following statement holds. Fix $r \geq 1$, $\lambda \in [r^{-1}, 1]$, $f \in H^{-1}(U_r)$, $w \in H^1(U_r)$ and let $u \in w + H_0^1(U_r)$ be the solution of (1.1). Also fix a function $v \in w + H_0^1(U_r)$ and define the functions $u_0, \bar{u}, \tilde{u} \in H_0^1(U_r)$ to be the solutions of the following equations (with null Dirichlet boundary condition on ∂U_r):*

$$\begin{aligned} (\lambda^2 - \nabla \cdot \mathbf{a}\nabla)u_0 &= f + \nabla \cdot \mathbf{a}\nabla v && \text{in } U_r, \\ -\nabla \cdot \bar{\mathbf{a}}\nabla \bar{u} &= \lambda^2 u_0 && \text{in } U_r, \\ (\lambda^2 - \nabla \cdot \mathbf{a}\nabla)\tilde{u} &= (\lambda^2 - \nabla \cdot \bar{\mathbf{a}}\nabla)\bar{u} && \text{in } U_r. \end{aligned}$$

For $\hat{v} \in w + H_0^1(U_r)$ defined by

$$(1.6) \quad \hat{v} := v + u_0 + \tilde{u},$$

we have the estimate

$$(1.7) \quad \|\nabla(\hat{v} - u)\|_{L^2(U_r)} \leq \mathcal{O}_s \left(C \ell(\lambda)^{\frac{1}{2}} \lambda^{\frac{1}{2}} \|\nabla(v - u)\|_{L^2(U_r)} \right).$$

The function $u \in H^1(U_r)$ appearing in Theorem 1.1 is the unknown we wish to approximate, and $v \in H^1(U_r)$ should be thought of as the current approximation to u . The function \hat{v} is then the new, updated approximation to u and the estimate (1.7) says that, if λ is chosen small enough, then the error in our approximation will be reduced by a multiplicative factor of $1/2$. As explained more precisely around (1.10) below, we can then iterate this procedure and obtain rapid convergence to the solution. The only assumption we make on v is that it satisfies the correct boundary condition, that is, $v \in w + H_0^1(U_r)$. In particular, we may begin the iteration with $v = w$ as the initial guess (or any other function with the correct boundary condition). The computation of \hat{v} reduces to solving the problems for u_0, \bar{u} , and \tilde{u} listed in the statement, and the point is that each of these problems is relatively inexpensive to compute, provided that λ is not too small. A fundamental aspect of the result is therefore that the required smallness of the parameter λ (so that (1.7) gives us a strict contraction in H^1) does not depend on the length scale r of the problem. In other words, we may need to take λ to be small, but it will still be of order one, no matter how large r is.

Similarly to standard multigrid methods, the equation for u_0 is meant to resolve the small-scale discrepancies between u and v . Note that the equation for u_0 can be rewritten as

$$(\lambda^2 - \nabla \cdot \mathbf{a}\nabla)u_0 = -\nabla \cdot \mathbf{a}\nabla(u - v) \quad \text{in } U_r.$$

The parameter λ^{-1} is the characteristic length scale of this problem, and in practice we will take it to be some fixed multiple of the scale of oscillations of the coefficients. The computation of u_0 can thus be decomposed into a large number of essentially unrelated elliptic problems posed on subdomains of side length of the order of λ^{-1} . In analogy with multigrid methods, we may also think of λ^{-2} as the number of elementary pre-smoothing steps performed during one global iteration.

As announced, we then use the homogenized operator on scales larger than λ^{-1} . This is what the problem for \bar{u} is meant to capture. Since the elliptic problem for \bar{u} involves the homogenized operator $-\nabla \cdot \bar{\mathbf{a}} \nabla$, it can be solved efficiently using the standard multigrid method. We note that the equation for \bar{u} can be rewritten, if desired, in the form

$$(1.8) \quad -\nabla \cdot \bar{\mathbf{a}} \nabla \bar{u} = -\nabla \cdot \mathbf{a} \nabla (u - v - u_0) \quad \text{in } U_r.$$

The final step of the iteration, involving the definition of \tilde{u} , is meant to add back some small-scale details to the function \bar{u} . It is analogous to the post-smoothing step in the standard V -cycle implementation of the multigrid method, and the parameter λ^{-2} represents the number of post-smoothing steps.

We next discuss the more probabilistic aspects involved in the statement of Theorem 1.1. Since the coefficient field is random, the statement of this theorem can only be valid with high probability, but not almost surely. Indeed, with non-zero probability, the coefficient field can be essentially arbitrary, and on such small-probability events, the idea of leveraging on homogenization can only perform badly (recall that we aim for a convergence result for large but fixed r , as opposed to asymptotic convergence). It may help the intuition to observe that, by Chebyshev's inequality, the assumption of (1.4) implies that

$$(1.9) \quad \forall x \geq 0, \quad \mathbb{P}[X \geq \theta x] \leq 2 \exp(-x^s),$$

and that conversely, the assumption of (1.9) implies that $X \leq \mathcal{O}_s(C\theta)$ for some constant $C(s) < \infty$ (see [3, Lemma A.1]).

We remark that Theorem 1.1 is new even when restricted to the subclass of periodic coefficient fields. In this case, both the probabilistic part of the estimate as well as the logarithmic factor of $\ell(\lambda)$ are not present, and (1.7) can be replaced with the simpler form

$$\|\nabla(\widehat{v} - u)\|_{L^2(U_r)} \leq C \lambda^{\frac{1}{2}} \|\nabla(v - u)\|_{L^2(U_r)}.$$

We stress that the probabilistic statement in (1.7) is valid for each fixed choice of $u, v \in H^1(U_r)$. In fact, the following stronger, uniform estimate is now proved in [19]. For each $s \in (0, 2)$, there exist a constant $C(s, p, U, \Lambda, d) < \infty$ and, for each $r \geq 1$ and $\lambda \in [r^{-1}, 1]$, a random variable $\mathcal{X}_{s,r,\lambda} : \Omega \rightarrow [0, +\infty]$ satisfying

$$\mathcal{X}_{s,r,\lambda} \leq \mathcal{O}_s(C)$$

such that, for every $u, v \in H^1(U_r)$ and \widehat{v} as in the statement of Theorem 1.1,

$$(1.10) \quad \|\nabla(\widehat{v} - u)\|_{L^2(U_r)} \leq \mathcal{X}_{s,r,\lambda} \ell(\lambda)^{\frac{1}{2}} \lambda^{\frac{1}{2}} (\log r)^{\frac{1}{s}} \|\nabla(v - u)\|_{L^2(U_r)}.$$

Moreover, the proof given in [19] does not require that the coefficient field be Hölder continuous. As is apparent in (1.10), the price one has to pay for the uniformity of this estimate in the functions u and v is a slight degradation of the

contraction factor, by a slowly diverging logarithmic factor of the domain size. Due to randomness, uniform estimates such as (1.10) must necessarily contain some logarithmic divergence in the domain size. Indeed, consider for instance the case of a coefficient field given by a random checkerboard in which we toss a fair coin, independently for each $z \in \mathbb{Z}^d$, the coefficient field in $z + [0, 1)^d$ to be either I_d or $2I_d$. Then, with probability tending to one as r tends to infinity, there will be in the domain U_r a region of space of side length of the order of $(\log r)^{\frac{1}{d}}$ where the coefficient field is constant equal to I_d . If the support of the solution we seek is concentrated in this region, then the iteration described in Theorem 1.1 will perform badly unless λ^{-1} is chosen larger than $(\log r)^{\frac{1}{d}}$.

The iteration proposed in Theorem 1.1 requires the user to make a judicious choice of the length scale λ^{-1} . Ideally, it would be preferable to devise an adaptive method which discovers a good choice for λ^{-1} automatically. The contraction of the iteration would then be guaranteed with probability one, and more subtle probabilistic quantifiers would instead enter into the complexity analysis of the method. A suitably designed adaptive algorithm would likely also work on more general coefficient fields than those considered here, allowing for instance to drop the assumption of stationarity. An assumption of approximate local stationarity would then also enter into the complexity analysis of the method. We leave the development of such adaptive methods to future work.

The method proposed here also requires that the user computes $\bar{\mathbf{a}}$ beforehand. An efficient method for doing so was presented in [26] in a discrete setting; see also [16, 11] and references therein for previous work on this problem. Moreover, one can check that in order to guarantee the contraction property of the iteration described in Theorem 1.1, say by a factor of 1/2, a coarse approximation of $\bar{\mathbf{a}}$, which may be off by a small but fixed positive amount, suffices.

The proof of Theorem 1.1 can be modified so that the L^2 norms in (1.7) are replaced by L^p norms, for any exponent $p < \infty$. Up to some additional logarithmic factors in λ , the contraction factor in the estimate would then be of order $\lambda^{\frac{1}{p}}$ rather than $\lambda^{\frac{1}{2}}$. This modification requires the application of large-scale Calderón-Zygmund-type L^p estimates which can be found in [3, Chapter 7]. The main required modification to the proof of Theorem 1.1 is simply to upgrade the two-scale expansion result of Theorem 3.1 from $p = 2$ to larger exponents by adapting the argument of [3, Theorem 7.10].

1.3. Previous works. There has been a lot of work on numerical algorithms that become sharp only in the limit of infinite scale separation (see for instance [27, 25, 6, 21, 10, 8, 1] and the references therein). That is, the error between the true solution u and its numerical approximation becomes small only as $r \rightarrow \infty$. Such algorithms typically have a computational complexity scaling sublinearly with the volume of the domain. An example of such a method in the context of the homogenization problem considered here is to compute an approximation of the solution to the homogenized equation. In addition to relying on scale separation, we note that such a sublinear method can only give an accurate global approximation in a weaker space such as L^2 , but not in stronger norms such as H^1 which are sensitive to small scale oscillations.

We now turn our attention to numerical algorithms that, like ours, converge to the true solution for each finite value of r . As pointed out in [18, 22], direct applications of standard multigrid methods result in coarse-scale systems that do not capture the relevant large-scale properties of the problem. Indeed, standard coarsening procedures produce effective coefficients that are simple arithmetic averages of the original coefficient field, instead of the homogenized coefficients. To remedy this problem, [18, 22] propose more subtle, matrix-dependent choices for the restriction and prolongation operators. The idea is to try to approximate a Schur complement calculation, while preserving some calculability constraints such as matrix sparsity. The method proposed there is shown numerically to perform better than simple averaging, but no theoretical guarantee is provided.

In [12, 13], the authors propose, in the periodic setting, to solve local problems for the correctors, deduce locally homogenized coefficients, and build coarsened operators from these. For the special two-dimensional case with $\mathbf{a}(x) = \tilde{\mathbf{a}}(x_1 - x_2)$ for some 1-periodic $\tilde{\mathbf{a}} \in C([0, 1]; \mathbb{R}_{\text{sym}}^{2 \times 2})$, they show (in our notation) that $O(r^{\frac{5}{3}} \log r)$ smoothing steps suffice to guarantee the contractivity of the two-step multigrid method (assuming that the chosen coarsening scale is a bounded multiple of the oscillation scale). For comparison, this roughly corresponds to the choice of $\lambda \simeq r^{-\frac{5}{6}}$ in our method. They also report better numerical performance than predicted by their theoretical arguments.

Beyond our current assumption of stationarity of the coefficient field, one can look for numerical methods for the resolution of general elliptic problems with rapidly oscillating coefficients. Possibly the simplest such method is to rely on the uniform ellipticity assumption (1.3) and appeal to a preconditioned conjugate gradient method, using the standard Laplacian as a preconditioner. However, the norm that is contracted at each iteration of this algorithm is the L^2 norm, as opposed to a contraction of the H^1 norm as obtained in the present paper¹. Moreover, the performance of this method degrades quickly if the ellipticity ratio Λ becomes large. In contrast, inspired by [2, 7], Chenlin Gu (at ENS Paris) has announced a result similar to that of Theorem 1.1 and [19] in the highly degenerate case of perforated media of percolation type, for which $\Lambda = \infty$.

Algebraic multigrid methods are intended to solve completely arbitrary linear systems of equations, by automatically discovering a hierarchy of coarsened problems [31]. In practice, it is however necessary to make some judicious choices of coarsening operators. In a sense, the present contribution as well as those of [18, 22, 12, 13] are descriptions of specific coarsening procedures which, under stronger assumptions such as stationarity, are shown to have fast convergence properties.

Many alternative approaches to the computation of elliptic problems with arbitrary coefficient fields have been developed. We mention in particular, without going into details, hierarchical matrices [5], generalized multiscale finite element methods [4, 9, 17], polyharmonic splines [29], local orthogonal decompositions [24], subspace correction methods [23] and gamblets [28]. While methods such as

¹Naturally, the L^2 and H^1 norms become equivalent after discretization; but a theoretical guarantee of contraction in H^1 ensures that the high frequencies can be resolved efficiently after only a few steps of the iteration, irrespectively of the size of the mesh refinement.

gamblets have been shown theoretically to have essentially linear complexity under weaker assumptions than those explored in the present paper, the construction and storage of the hierarchy of gamblets may actually be quite expensive in practice (we are not aware of large-scale computations that use gamblets; the main numerical example in [28] has $2^{24} \simeq 1.7 \cdot 10^7$ degrees of freedom). Methods such as local orthogonal decompositions introduce an intermediate scale, often denoted by H , inbetween the microscopic and the macroscopic scales, and an adapted basis of local functions is computed at this level. In this framework, the numerical error is bounded from below by a multiple of H . The method presented in Theorem 1.1 shares some aspects of this idea in that it also introduces an intermediate scale λ^{-1} ; however, the final numerical error is not constrained by this choice, and can be made arbitrarily low irrespectively of the value of λ .

1.4. Organization of the paper. We introduce some more notation in Section 2. Section 3 is devoted to the proof of Theorem 1.1. We report on our numerical results in Section 4. Finally, an appendix recalls some classical Sobolev and elliptic estimates for the reader's convenience.

2. NOTATION

In this section, we collect some notation used throughout the paper. Recall that the notation $\mathcal{O}_s(\cdot)$ was defined in (1.4). We will need the following fact, which says that \mathcal{O}_s is behaving like a norm: for each $s \in (0, \infty)$, there exists $C_s < \infty$ (with $C_s = 1$ for $s \geq 1$) such that the following triangle inequality for $\mathcal{O}_s(\cdot)$ holds: for any measure space (E, \mathcal{S}, μ) , measurable function $\theta : E \rightarrow (0, \infty)$ and jointly measurable family $\{X(z)\}_{z \in E}$ of random variables, we have (see [3, Lemma A.4])

$$(2.1) \quad \forall z \in E, X(z) \leq \mathcal{O}_s(\theta(z)) \implies \int_E X d\mu \leq \mathcal{O}_s\left(C_s \int_E \theta d\mu\right).$$

We denote by (e_1, \dots, e_d) the canonical basis of \mathbb{R}^d , and write $B(x, r) \subseteq \mathbb{R}^d$ for the Euclidean ball centered at $x \in \mathbb{R}^d$ and of radius $r > 0$. For a Borel set $V \subseteq \mathbb{R}^d$, we denote its Lebesgue measure by $|V|$. If $|V| < \infty$, then for every $p \in [1, \infty)$ and $f \in L^p(V)$ we write the scaled L^p norm of f by

$$\|f\|_{\underline{L}^p(V)} := \left(|V|^{-1} \int_V f^p\right)^{\frac{1}{p}} = |V|^{-\frac{1}{p}} \|f\|_{L^p(V)}.$$

For each $k \in \mathbb{N}$, we denote by $H^k(V)$ the classical Sobolev space on V , whose norm is given by

$$\|f\|_{H^k(V)} := \sum_{0 \leq |\beta| \leq k} \|\partial^\beta f\|_{L^2(V)}.$$

In the expression above, the parameter $\beta = (\beta_1, \dots, \beta_d)$ is a multi-index in \mathbb{N}^d , and we used the notation

$$|\beta| := \sum_{i=1}^d \beta_i \quad \text{and} \quad \partial^\beta f = \partial_{x_1}^{\beta_1} \dots \partial_{x_d}^{\beta_d} f.$$

Whenever $|V| < \infty$, we define the scaled Sobolev norm by

$$\|f\|_{\underline{H}^k(V)} := \sum_{0 \leq |\beta| \leq k} |V|^{\frac{|\beta|-k}{d}} \|\partial^\beta f\|_{\underline{L}^2(V)}.$$

We denote by $H_0^1(V)$ the completion in $H^1(V)$ of the space $C_c^\infty(V)$ of smooth functions with compact support in V . We write $H^{-1}(V)$ for the dual space to $H_0^1(V)$, which we endow with the (scaled) norm

$$\|f\|_{\underline{H}^{-1}(V)} := \sup \left\{ |V|^{-1} \int_V f g, \quad g \in H_0^1(V), \quad \|g\|_{\underline{H}^1(V)} \leq 1 \right\}.$$

The integral sign above is an abuse of notation and should be understood as the duality pairing between $H^{-1}(V)$ and $H_0^1(V)$. The spaces $H^{-1}(V)$ and $H_0^1(V)$ can be continuously embedded into the space of distributions, and we make sure that the duality pairing is consistent with the integral expression above whenever f and g are smooth functions. For every $r > 0$ and $x \in \mathbb{R}^d$, we denote the time-slice of the heat kernel which has length scale r by

$$(2.2) \quad \Phi_r(x) := (4\pi r^2)^{-\frac{d}{2}} \exp\left(-\frac{x^2}{4r^2}\right).$$

We denote by $\zeta \in C_c^\infty(\mathbb{R}^d)$ the standard mollifier

$$(2.3) \quad \zeta(x) := \begin{cases} c_d \exp(-(1 - |x|^2)^{-1}) & \text{if } |x| < 1, \\ 0 & \text{if } |x| \geq 1, \end{cases}$$

where the constant c_d is chosen so that $\int_{\mathbb{R}^d} \zeta = 1$. For $f \in L^p(\mathbb{R}^d)$ and $g \in L^{p'}(\mathbb{R}^d)$ with $\frac{1}{p} + \frac{1}{p'} = 1$, we denote the convolution of f and g by

$$f * g(x) := \int_{\mathbb{R}^d} f(y)g(x - y) dy.$$

3. PROOF OF THEOREM 1.1

This section is devoted to the proof of Theorem 1.1. We begin by introducing the notion of (*first-order*) *corrector*: for each $p \in \mathbb{R}^d$, the corrector in the direction of p is the function $\phi_p \in H_{\text{loc}}^1(\mathbb{R}^d)$ solving

$$-\nabla \cdot \mathbf{a}(p + \nabla \phi_p) = 0 \quad \text{in } \mathbb{R}^d,$$

and such that the mapping $x \mapsto \nabla \phi_p(x)$ is \mathbb{Z}^d -stationary and satisfies

$$\mathbb{E} \left[\int_{[0,1]^d} \nabla \phi_p \right] = 0.$$

The corrector ϕ_p is unique up to an additive constant (see [3, Definition 4.2] for instance). We also recall that one can define the homogenized matrix $\bar{\mathbf{a}} \in \mathbb{R}_{\text{sym}}^{d \times d}$ via the formula

$$\forall p \in \mathbb{R}^d, \quad \bar{\mathbf{a}}p = \mathbb{E} \left[\int_{[0,1]^d} \mathbf{a}(p + \nabla \phi_p) \right],$$

or equivalently,

$$\forall p \in \mathbb{R}^d, \quad p \cdot \bar{\mathbf{a}}p = \mathbb{E} \left[\int_{[0,1]^d} (p + \nabla \phi_p) \cdot \mathbf{a}(p + \nabla \phi_p) \right],$$

and in particular, as a consequence of (1.3), we have

$$(3.1) \quad \forall \xi \in \mathbb{R}^d, \quad \Lambda^{-1}|\xi|^2 \leq \xi \cdot \bar{\mathbf{a}}\xi \leq \Lambda|\xi|^2.$$

For each $k \in \{1, \dots, d\}$ and $\lambda > 0$, we denote

$$\phi_{e_k}^{(\lambda)} := \phi_{e_k} - \phi_{e_k} * \Phi_{\lambda^{-1}}.$$

A key ingredient in the proof of Theorem 1.1 is the following *quantitative two-scale expansion* for the operator $(\lambda^2 - \nabla \cdot \mathbf{a} \nabla)$. It is the only input from the quantitative theory of stochastic homogenization used in this paper and it follows from some estimates which can be found in [3].

Theorem 3.1 (Two-scale expansion and error estimate). *For each $s \in (0, 2)$, there exists a constant $C(s, U, \Lambda, \alpha, d) < \infty$ such that, for every $r \geq 1$, $\lambda \in [r^{-1}, 1]$, and $\bar{v} \in H_0^1(U_r) \cap H^2(U_r)$, defining*

$$(3.2) \quad w := \bar{v} + \sum_{k=1}^d \phi_{e_k}^{(\lambda)} \partial_{x_k} \bar{v},$$

we have the estimate

$$(3.3) \quad \|\nabla \cdot (\mathbf{a} \nabla w - \bar{\mathbf{a}} \nabla \bar{v})\|_{\underline{H}^{-1}(U_r)} \leq \mathcal{O}_s \left(C\ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} + C\lambda^{\frac{d}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)} \right).$$

Moreover, for every $\mu \in [0, \lambda]$ and $v \in H_0^1(U_r)$ such that

$$(3.4) \quad (\mu^2 - \nabla \cdot \mathbf{a} \nabla) v = (\mu^2 - \nabla \cdot \bar{\mathbf{a}} \nabla) \bar{v},$$

we have the estimate

$$(3.5) \quad \|v - w\|_{\underline{H}^1(U_r)} + (\mu + r^{-1}) \|v - \bar{v}\|_{\underline{L}^2(U_r)} + (\mu + r^{-1})^2 \|v - \bar{v}\|_{\underline{H}^{-1}(U_r)} \\ \leq \mathcal{O}_s \left(C \left(\mu\ell(\lambda) + \lambda^{\frac{d}{2}} \right) \|\bar{v}\|_{\underline{H}^1(U_r)} + C\ell(\lambda)^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} + C\ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} \right).$$

The proof of Theorem 3.1 follows that of a similar result from Chapter 6 of [3]. The main difference here is the presence of the zeroth order term with the factor of μ^2 , which presents no additional difficulty. We begin by recalling the concept of a flux corrector and stating some estimates on the correctors proved in [3].

For each $p \in \mathbb{R}^d$, we denote the (centered) flux of the corrector ϕ_p by

$$(\mathbf{g}_{p,i})_{1 \leq i \leq d} = \mathbf{g}_p := \mathbf{a}(p + \nabla \phi_p) - \bar{\mathbf{a}}p.$$

Since $\nabla \cdot \mathbf{g}_p = 0$, the flux of the corrector admits a representation as the “curl” of some vector potential, by Helmholtz’s theorem. This vector potential, the *flux corrector*, will be useful for the proof of Theorem 3.1. For each $p \in \mathbb{R}^d$, the vector potential $(\mathbf{S}_{p,ij})_{1 \leq i, j \leq d}$ is a matrix-valued random field with entries in $H_{\text{loc}}^1(\mathbb{R}^d)$ satisfying, for each $i, j \in \{1, \dots, d\}$,

$$(3.6) \quad \begin{aligned} \mathbf{S}_{p,ij} &= -\mathbf{S}_{p,ji}, \\ \nabla \cdot \mathbf{S}_p &= \mathbf{g}_p, \end{aligned}$$

and such that $x \mapsto \nabla \mathbf{S}_{p,ij}(x)$ is a stationary random field with mean zero. In (3.6), we used the shorthand notation

$$(\nabla \cdot \mathbf{S}_e)_i := \sum_{j=1}^d \partial_{x_j} \mathbf{S}_{e,ij}.$$

The conditions above do not specify the flux corrector uniquely. One way to “fix the gauge” is to enforce that, for each $i, j \in \{1, \dots, d\}$,

$$\Delta \mathbf{S}_{p,ij} = \partial_{x_j} \mathbf{g}_{p,i} - \partial_{x_i} \mathbf{g}_{p,j}.$$

This latter choice then defines $\mathbf{S}_{p,ij}$ uniquely, up to the addition of a constant. We refer to [3, Section 6.1] for more precision on this construction. We set

$$(3.7) \quad \mathbf{S}_e^{(\lambda)} := \mathbf{S}_e - \mathbf{S}_e * \Phi_{\lambda^{-1}}.$$

The fundamental ingredient for the proof of Theorem 3.1 is the following proposition, which quantifies the convergence to zero of the spatial averages of the gradients of the correctors.

Proposition 3.2 (Corrector estimates). *For each $s \in (0, 2)$, there exists a constant $C(s, U, \Lambda, \alpha, d) < \infty$ such that for every $\lambda \in (0, 1)$, $x \in \mathbb{R}^d$ and $i, j, k \in \{1, \dots, d\}$,*

$$(3.8) \quad |\nabla \phi_{e_k}(x)| \leq \mathcal{O}_s(C),$$

$$(3.9) \quad |(\nabla \phi_{e_k} * \Phi_{\lambda^{-1}})(x)| + |(\nabla \mathbf{S}_{e_k,ij} * \Phi_{\lambda^{-1}})(x)| \leq \mathcal{O}_s\left(C\lambda^{\frac{d}{2}}\right),$$

$$(3.10) \quad \left|\phi_{e_k}^{(\lambda)}(x)\right| + \left|\mathbf{S}_{e_k,ij}^{(\lambda)}(x)\right| = \mathcal{O}_s(C\ell(\lambda)).$$

Proof. By [3, Lemma 4.4], we have

$$\|\nabla \phi_{e_k}\|_{L^2(B(0,1))} \leq \mathcal{O}_s(C).$$

By the assumption of (1.2), we can apply standard Schauder estimates, see e.g. [20, Theorems 3.1 and 3.8], to deduce (3.8). The estimates in (3.9) are proved in [3, Theorem 4.9 and Proposition 6.2]. The estimates in (3.10) also follow from [3, Theorem 4.9 and Proposition 6.2], combined with the assumption of (1.2) and the Schauder estimate in [20, Corollary 3.2 and Theorem 3.8]. \square

In the next lemma, we provide a convenient representation of $\nabla \cdot \mathbf{a} \nabla w$ in terms of the correctors.

Lemma 3.3. *Let $\lambda > 0$, $\bar{v} \in H^1(U_r)$, and let $w \in H^1(U_r)$ be defined by (3.2). Then*

$$\nabla \cdot (\mathbf{a} \nabla w - \bar{\mathbf{a}} \nabla \bar{v}) = \nabla \cdot \mathbf{F},$$

where the i -th component of the vector field \mathbf{F} is given by

$$(3.11) \quad \mathbf{F}_i := \sum_{j,k=1}^d \left(\mathbf{a}_{ij} \phi_{e_k}^{(\lambda)} - \mathbf{S}_{e_k,ij}^{(\lambda)} \right) \partial_{x_j} \partial_{x_k} \bar{v} \\ + \sum_{j,k=1}^d \left(\mathbf{a}_{ij} (\partial_{x_j} \phi_{e_k} * \Phi_{\lambda^{-1}}) + \partial_{x_j} \mathbf{S}_{e_k,ij} * \Phi_{\lambda^{-1}} \right) \partial_{x_k} \bar{v}.$$

Proof. The argument is very similar to that for [3, Lemma 6.6], the main difference being that the definition of $\phi_{e_k}^{(\lambda)}$ is slightly different from that of $\phi_{e_k}^\varepsilon$ there. We recall the argument here for the reader's convenience. Observe that, for each $j \in \{1, \dots, d\}$,

$$(3.12) \quad \partial_{x_j} w = \sum_{k=1}^d \left((\delta_{jk} + \partial_{x_j} \phi_{e_k}) \partial_{x_k} \bar{v} - (\partial_{x_j} \phi_{e_k} * \Phi_{\lambda^{-1}}) \partial_{x_k} \bar{v} + \phi_{e_k}^{(\lambda)} \partial_{x_j} \partial_{x_k} \bar{v} \right).$$

We start by studying the contribution of the first summand. By (3.6) and (3.7), we have, for every $i, k \in \{1, \dots, d\}$,

$$\sum_{j=1}^d \partial_{x_j} \mathbf{S}_{e_k,ij}^{(\lambda)} = \sum_{j=1}^d \left(\mathbf{a}_{ij} (\delta_{jk} + \partial_{x_j} \phi_{e_k}) - \bar{\mathbf{a}}_{ij} \delta_{jk} - \partial_{x_j} \mathbf{S}_{e_k,ij} * \Phi_{\lambda^{-1}} \right).$$

We deduce that, for each $i \in \{1, \dots, d\}$,

$$(3.13) \quad \sum_{j,k=1}^d \mathbf{a}_{ij} (\delta_{jk} + \partial_{x_j} \phi_{e_k}) \partial_{x_k} \bar{v} = \sum_{j,k=1}^d \left(\bar{\mathbf{a}}_{ij} \delta_{jk} + \partial_{x_j} \mathbf{S}_{e_k, ij}^{(\lambda)} + \partial_{x_j} \mathbf{S}_{e_k, ij} * \Phi_{\lambda^{-1}} \right) \partial_{x_k} \bar{v},$$

and thus

$$\begin{aligned} \sum_{i,j,k=1}^d \partial_{x_i} \left(\mathbf{a}_{ij} (\delta_{jk} + \partial_{x_j} \phi_{e_k}) \partial_{x_k} \bar{v} \right) &= \nabla \cdot \bar{\mathbf{a}} \nabla \bar{v} \\ &+ \sum_{i,j,k=1}^d \partial_{x_i} \left(\partial_{x_j} \mathbf{S}_{e_k, ij}^{(\lambda)} \partial_{x_k} \bar{v} \right) + \sum_{i,j,k=1}^d \partial_{x_i} \left((\partial_{x_j} \mathbf{S}_{e_k, ij} * \Phi_{\lambda^{-1}}) \partial_{x_k} \bar{v} \right). \end{aligned}$$

By the skew-symmetry of $\mathbf{S}_e^{(\lambda)}$, we have

$$\begin{aligned} 0 &= \sum_{i,j,k=1}^d \partial_{x_i} \partial_{x_j} \left(\mathbf{S}_{e_k, ij}^{(\lambda)} \partial_{x_k} \bar{v} \right) \\ &= \sum_{i,j,k=1}^d \partial_{x_i} \left(\partial_{x_j} \mathbf{S}_{e_k, ij}^{(\lambda)} \partial_{x_k} \bar{v} \right) + \sum_{i,j,k=1}^d \partial_{x_i} \left(\mathbf{S}_{e_k, ij}^{(\lambda)} \partial_{x_j} \partial_{x_k} \bar{v} \right), \end{aligned}$$

and thus

$$\begin{aligned} \sum_{i,j,k=1}^d \partial_{x_i} \left(\mathbf{a}_{ij} (\delta_{jk} + \partial_{x_j} \phi_{e_k}) \partial_{x_k} \bar{v} \right) &= \nabla \cdot \bar{\mathbf{a}} \nabla \bar{v} \\ &- \sum_{i,j,k=1}^d \partial_{x_i} \left(\mathbf{S}_{e_k, ij}^{(\lambda)} \partial_{x_j} \partial_{x_k} \bar{v} \right) + \sum_{i,j,k=1}^d \partial_{x_i} \left((\partial_{x_j} \mathbf{S}_{e_k, ij} * \Phi_{\lambda^{-1}}) \partial_{x_k} \bar{v} \right). \end{aligned}$$

Recalling (3.12), we obtain the announced result. \square

We next present the proof of Theorem 3.1, which can be compared to the one of [3, Theorem 6.9].

Proof of Theorem 3.1. We will proceed by proving first (3.3), and then the H^1 , L^2 and H^{-1} estimates appearing in (3.5), in this order. We decompose the arguments into seven steps.

Step 1. We prove (3.3). In view of Lemma 3.3, it suffices to show that, for the vector field \mathbf{F} defined in (3.11),

$$(3.14) \quad \|\mathbf{F}\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C\ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} + C\lambda^{\frac{d}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)} \right).$$

We estimate each of the terms appearing in the definition of \mathbf{F} . By Proposition 3.2 and (2.1), we have, for every $i, j, k \in \{1, \dots, d\}$,

$$\left\| \left(\mathbf{a}_{ij} \phi_{e_k}^{(\lambda)} - \mathbf{S}_{e_k, ij}^{(\lambda)} \right) \partial_{x_j} \partial_{x_k} \bar{v} \right\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C\ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} \right),$$

as well as

$$\left\| \left(\mathbf{a}_{ij} (\partial_{x_j} \phi_{e_k} * \Phi_{\lambda^{-1}}) + \partial_{x_j} \mathbf{S}_{e_k, ij} * \Phi_{\lambda^{-1}} \right) \partial_{x_k} \bar{v} \right\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C\lambda^{\frac{d}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)} \right),$$

and thus (3.14) follows.

Step 2. In order to show (3.5), we first need to evaluate the contribution of a boundary layer. For every $\ell \geq 0$, we write $\zeta_\ell := \ell^{-d} \zeta(\ell^{-1} \cdot)$ (recall the definition of ζ in (2.3)) and

$$(3.15) \quad U_{r,\ell} := \{x \in U_r : \text{dist}(x, \partial U_r) > \ell\}.$$

With the definition of $\ell(\lambda)$ given in (1.5), we set

$$T := \left(\mathbf{1}_{\mathbb{R}^d \setminus U_{r,2\ell(\lambda)}} * \zeta_{\ell(\lambda)} \right) \sum_{k=1}^d \phi_{e_k}^{(\lambda)} \partial_{x_k} \bar{v}.$$

We will use the function T as a test function for an upper bound on the size of the actual boundary layer in the next step. In this step, we show that there exists $C(s, U, \Lambda, \alpha, d) < \infty$ such that

$$(3.16) \quad \|\nabla T\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C \ell(\lambda)^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} + C \ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} \right)$$

and

$$(3.17) \quad \|T\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C \ell(\lambda)^{\frac{3}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} \right).$$

By the chain rule,

$$\|\nabla T\|_{L^2(U_r)} \leq C \sum_{k=1}^d \left\| \left(\frac{|\nabla \bar{v}|}{\ell(\lambda)} + |\nabla^2 \bar{v}| \right) \left| \phi_{e_k}^{(\lambda)} \right| + |\nabla \bar{v}| \left| \nabla \phi_{e_k}^{(\lambda)} \right| \right\|_{L^2(U_r \setminus U_{r,3\ell(\lambda)})}.$$

By Proposition 3.2 and (2.1), we have

$$\left\| |\nabla^2 \bar{v}| \left| \phi_{e_k}^{(\lambda)} \right| \right\|_{L^2(U_r \setminus U_{r,3\ell(\lambda)})} \leq \mathcal{O}_s \left(C \ell(\lambda) \|\nabla^2 \bar{v}\|_{L^2(U_r)} \right).$$

Similarly,

$$(3.18) \quad \left\| \frac{|\nabla \bar{v}|}{\ell(\lambda)} \left| \phi_{e_k}^{(\lambda)} \right| \right\|_{L^2(U_r \setminus U_{r,3\ell(\lambda)})} \leq \mathcal{O}_s \left(C \|\nabla \bar{v}\|_{L^2(U_r \setminus U_{r,3\ell(\lambda)})} \right),$$

and by Proposition A.1,

$$(3.19) \quad \|\nabla \bar{v}\|_{L^2(U_r \setminus U_{r,3\ell(\lambda)})} \leq C \ell(\lambda)^{\frac{1}{2}} r^{\frac{d}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}}.$$

Finally, using again Proposition 3.2 and (2.1), we have

$$\left\| |\nabla \bar{v}| \left| \nabla \phi_{e_k}^{(\lambda)} \right| \right\|_{L^2(U_r \setminus U_{r,3\ell(\lambda)})} \leq \mathcal{O}_s \left(C \|\nabla \bar{v}\|_{L^2(U_r \setminus U_{r,3\ell(\lambda)})} \right),$$

and we can appeal once more to (3.19) to estimate the norm of $\nabla \bar{v}$ on the right side above. This completes the proof of (3.16). The estimate (3.17) follows from (3.18) and (3.19).

Step 3. We now evaluate the size of the boundary layer $b \in H^1(U_r)$ defined as the solution of

$$(3.20) \quad \begin{cases} (\mu^2 - \nabla \cdot \mathbf{a} \nabla) b = 0 & \text{in } U_r, \\ b = \sum_{k=1}^d \phi_{e_k}^{(\lambda)} \partial_{x_k} \bar{v} & \text{on } \partial U_r. \end{cases}$$

Since T and b share the same boundary condition on ∂U_r , by the variational formulation of (3.20), we have

$$\int_{U_r} (\mu^2 b^2 + \nabla b \cdot \mathbf{a} \nabla b) \leq \int_{U_r} (\mu^2 T^2 + \nabla T \cdot \mathbf{a} \nabla T).$$

By the result of the previous step, we thus obtain, for every $\mu \in [0, \lambda]$,

$$(3.21) \quad \mu \|b\|_{\underline{L}^2(U_r)} + \|\nabla b\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C \ell(\lambda)^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} + C \ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} \right).$$

Step 4. We are now prepared to prove that

$$(3.22) \quad \|\nabla(v-w)\|_{\underline{L}^2(U_r)} + \mu \|v-w\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C \left(\mu \ell(\lambda) + \lambda^{\frac{d}{2}} \right) \|\bar{v}\|_{\underline{H}^1(U_r)} + C \ell(\lambda)^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} + C \ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} \right).$$

For concision, we define

$$\mathcal{X}_1 := \|\nabla \cdot (\bar{\mathbf{a}} \nabla \bar{v} - \mathbf{a} \nabla w)\|_{\underline{H}^{-1}(U_r)},$$

and recall that, by (3.3),

$$(3.23) \quad \mathcal{X}_1 \leq \mathcal{O}_s \left(C \ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} + C \lambda^{\frac{d}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)} \right).$$

Moreover, by (3.4) and (3.20),

$$\begin{aligned} -\nabla \cdot (\bar{\mathbf{a}} \nabla \bar{v} - \mathbf{a} \nabla w) &= -\nabla \cdot \mathbf{a} \nabla(v-w) + \mu^2(v-\bar{v}) \\ &= -\nabla \cdot \mathbf{a} \nabla(v-w+b) + \mu^2(v-\bar{v}+b). \end{aligned}$$

Since $v-w+b \in H_0^1(U_r)$, we deduce that

$$\begin{aligned} |U_r|^{-1} \int_{U_r} (\nabla(v-w+b) \cdot \mathbf{a} \nabla(v-w+b) + \mu^2(v-w+b)(v-\bar{v}+b)) \\ \leq \mathcal{X}_1 \|\nabla(v-w+b)\|_{\underline{L}^2(U_r)}, \end{aligned}$$

and by the uniform ellipticity of \mathbf{a} and Hölder's inequality,

$$\begin{aligned} \|\nabla(v-w+b)\|_{\underline{L}^2(U_r)}^2 + \mu^2 \|v-w+b\|_{\underline{L}^2(U_r)}^2 \\ \leq C \mathcal{X}_1 \|\nabla(v-w+b)\|_{\underline{L}^2(U_r)} + \mu^2 \|w-\bar{v}\|_{\underline{L}^2(U_r)} \|v-w+b\|_{\underline{L}^2(U_r)}. \end{aligned}$$

Using Proposition 3.2 and (2.1), we verify that

$$(3.24) \quad \|w-\bar{v}\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C \ell(\lambda) \|\bar{v}\|_{\underline{H}^1(U_r)} \right).$$

Combining these two estimates with (3.23) and Young's inequality, we obtain that

$$\begin{aligned} \|\nabla(v-w+b)\|_{\underline{L}^2(U_r)} + \mu \|v-w+b\|_{\underline{L}^2(U_r)} \\ \leq \mathcal{O}_s \left(C \left(\mu \ell(\lambda) + \lambda^{\frac{d}{2}} \right) \|\bar{v}\|_{\underline{H}^1(U_r)} + C \ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} \right). \end{aligned}$$

An application of (3.21) then yields the announced estimate (3.22).

Step 5. In this step, we complete the proof of the fact that $\|v - w\|_{\underline{H}^1(U_r)}$ is bounded by the right side of (3.5). In view of (3.22), it suffices to show that

$$(3.25) \quad r^{-1}\|v - w\|_{\underline{L}^2(U_r)} \leq \mathcal{O}_s \left(C \left(\mu \ell(\lambda) + \lambda^{\frac{d}{2}} \right) \|\bar{v}\|_{\underline{H}^1(U_r)} + C \ell(\lambda)^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} + C \ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} \right).$$

By (3.16) and (3.22), we have

$$\begin{aligned} & \|\nabla(v - w + T)\|_{\underline{L}^2(U_r)} \\ & \leq \mathcal{O}_s \left(C \left(\mu \ell(\lambda) + \lambda^{\frac{d}{2}} \right) \|\bar{v}\|_{\underline{H}^1(U_r)} + C \ell(\lambda)^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{v}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} + C \ell(\lambda) \|\bar{v}\|_{\underline{H}^2(U_r)} \right). \end{aligned}$$

The estimate (3.25) then follows by the Poincaré inequality and (3.17).

Step 6. We now complete the proof that $(\mu + r^{-1})\|v - \bar{v}\|_{\underline{L}^2(U_r)}$ is bounded by the right side of (3.5). For $\mu \geq r^{-1}$, the result follows from (3.22) and (3.24), while $\mu \leq r^{-1}$, it follows from (3.5) and (3.24).

Step 7. We finally complete the proof of (3.5) by showing the estimate for the H^{-1} norm of $v - \bar{v}$. If $\mu \leq r^{-1}$, then the conclusion is immediate from the estimate on the L^2 norm of $v - \bar{v}$, by scaling. Otherwise, by the equations for v and \bar{v} , we have

$$\mu^2(v - \bar{v}) = \nabla \cdot (\mathbf{a}\nabla v - \bar{\mathbf{a}}\nabla \bar{v}),$$

and moreover,

$$\begin{aligned} & \|\nabla \cdot (\mathbf{a}\nabla v - \bar{\mathbf{a}}\nabla \bar{v})\|_{\underline{H}^{-1}(U_r)} \\ & \leq \|\nabla \cdot (\mathbf{a}\nabla w - \bar{\mathbf{a}}\nabla \bar{v})\|_{\underline{H}^{-1}(U_r)} + \|\nabla \cdot (\mathbf{a}\nabla v - \mathbf{a}\nabla w)\|_{\underline{H}^{-1}(U_r)} \\ & \leq \|\nabla \cdot (\mathbf{a}\nabla w - \bar{\mathbf{a}}\nabla \bar{v})\|_{\underline{H}^{-1}(U_r)} + C \|\nabla v - \nabla w\|_{\underline{L}^2(U_r)}. \end{aligned}$$

The terms on the right side above have been estimated in (3.3) and (3.22) respectively, so the proof is complete. \square

We next give the proof of the main result.

Proof of Theorem 1.1. Let $u, v, u_0, \bar{u}, \tilde{u} \in H^1(U_r)$ be as in the statement of Theorem 1.1. We first show the a priori estimates

$$(3.26) \quad \lambda \|u_0\|_{\underline{L}^2(U_r)} + \|\nabla u_0\|_{\underline{L}^2(U_r)} \leq C \|u - v\|_{\underline{H}^1(U_r)},$$

and

$$(3.27) \quad \|\bar{u}\|_{\underline{H}^1(U_r)} + \lambda^{-1} \|\bar{u}\|_{\underline{H}^2(U_r)} \leq C \|u - v\|_{\underline{H}^1(U_r)}.$$

By the variational formulation of the equation for $u_0 \in H_0^1(U_r)$, we have

$$\int_{U_r} (\lambda^2 u_0^2 + \nabla u_0 \cdot \mathbf{a}\nabla u_0) = \int_{U_r} \nabla u_0 \cdot \mathbf{a}\nabla(u - v).$$

By Hölder's and Young's inequalities and the uniform ellipticity of \mathbf{a} , we get (3.26). Using the equation (1.8) satisfied by $\bar{u} \in H_0^1(U_r)$ and the estimate (3.26), we deduce

$$\begin{aligned} \|\nabla \bar{u}\|_{\underline{L}^2(U_r)} & \leq C \|\nabla(u - v - u_0)\|_{\underline{L}^2(U_r)} \\ & \leq C \|\nabla(u - v)\|_{\underline{L}^2(U_r)}. \end{aligned}$$

By Proposition A.2 and the L^2 estimate in (3.26), we also have

$$\|\bar{u}\|_{\underline{H}^2(U_r)} \leq C\lambda^2 \|u_0\|_{\underline{L}^2(U_r)} \leq C\lambda \|u - v\|_{\underline{H}^1(U_r)},$$

as announced in (3.27).

We now introduce the two-scale expansion

$$w := \bar{u} + \sum_{k=1}^d \phi_{e_k}^{(\lambda)} \partial_{x_k} \bar{u}.$$

Using the equation for \bar{u} in (1.8) and Theorem 3.1 with $\mu = 0$, we obtain

$$\begin{aligned} & \|v + u_0 + w - u\|_{\underline{H}^1(U_r)} \\ & \leq \mathcal{O}_s \left(C\lambda^{\frac{d}{2}} \|\bar{u}\|_{\underline{H}^1(U_r)} + C\ell(\lambda)^{\frac{1}{2}} \|\bar{u}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{u}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} + C\ell(\lambda) \|\bar{u}\|_{\underline{H}^2(U_r)} \right), \end{aligned}$$

and thus, by (3.27),

$$(3.28) \quad \|v + u_0 + w - u\|_{\underline{H}^1(U_r)} \leq \mathcal{O}_s \left(C\ell(\lambda)^{\frac{1}{2}} \lambda^{\frac{1}{2}} \|u - v\|_{\underline{L}^2(U_r)} \right).$$

In order to complete the proof of Theorem 1.1, there remains to estimate the H^1 norm of $w - \tilde{u}$. By the equation for \tilde{u} , Theorem 3.1 and (3.27), we have

$$\begin{aligned} & \|w - \tilde{u}\|_{\underline{H}^1(U_r)} \\ & \leq \mathcal{O}_s \left(C \left(\lambda\ell(\lambda) + \lambda^{\frac{d}{2}} \right) \|\bar{u}\|_{\underline{H}^1(U_r)} \right. \\ & \quad \left. + C\ell(\lambda)^{\frac{1}{2}} \|\bar{u}\|_{\underline{H}^1(U_r)}^{\frac{1}{2}} \|\bar{u}\|_{\underline{H}^2(U_r)}^{\frac{1}{2}} + C\ell(\lambda) \|\bar{u}\|_{\underline{H}^2(U_r)} \right) \\ & \leq \mathcal{O}_s \left(C\ell(\lambda)^{\frac{1}{2}} \lambda^{\frac{1}{2}} \|u - v\|_{\underline{H}^1(U_r)} \right), \end{aligned}$$

as desired. \square

4. NUMERICAL RESULTS

In this section, we report on numerical tests demonstrating the performance of the iterative method described in Theorem 1.1. The code used in the tests can be consulted at

https://github.com/ahannuka/homo_mg

Throughout this section, we consider a two-dimensional *random checkerboard* coefficient field $x \mapsto \mathbf{a}(x)$, which is defined as follows: we give ourselves a family $(b(z))_{z \in \mathbb{Z}^2}$ of independent random variables such that for every $z \in \mathbb{Z}^2$,

$$\mathbb{P}[b(z) = 1] = \mathbb{P}[b(z) = 9] = \frac{1}{2}.$$

We then set, for every $x \in z + [0, 1)^2$,

$$\mathbf{a}(x) := b(z) I_2,$$

where I_2 denotes the 2-by-2 identity matrix. For this particular coefficient field, the homogenized matrix can be computed analytically as $\bar{\mathbf{a}} = 3I_2$ (see [3, Exercise 2.3]). When such an analytical expression does not exist, the homogenized coefficient can be approximated numerically, for example, by using the method presented in [26].

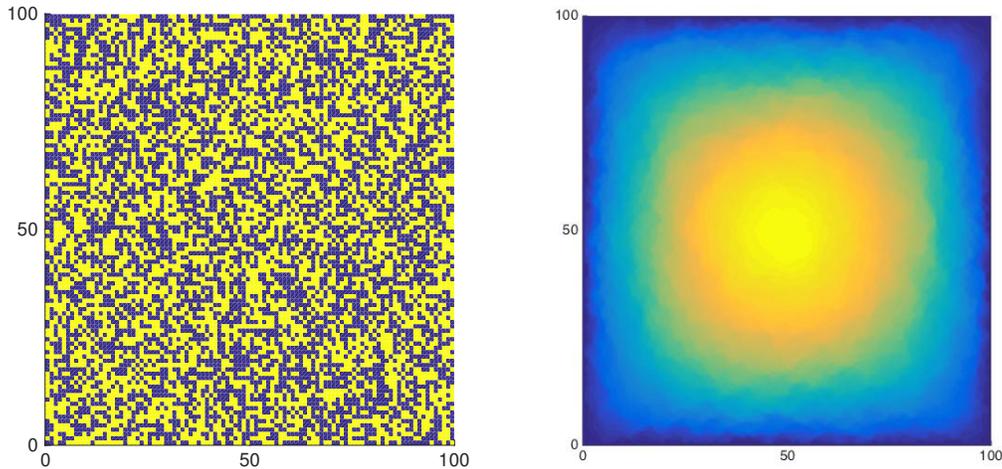


FIGURE 1. On the left, a typical realization of the coefficient field $\mathbf{a}(x)$, with $r = 100$ (yellow corresponds to the value 1 and blue to the value 9). On the right, the corresponding solution.

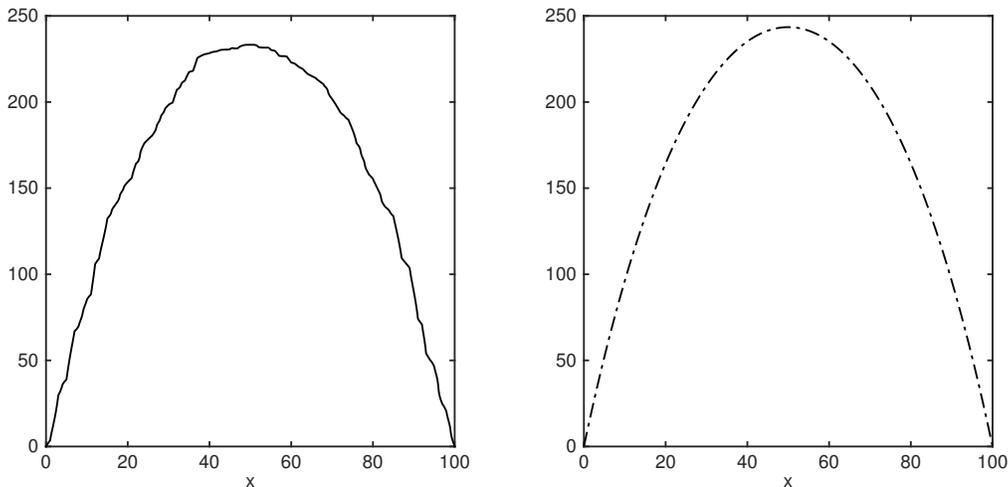


FIGURE 2. On the left, the FE-solution to the heterogeneous problem, and on the right, the FE-solution to the corresponding homogenized problem. Both solutions are plotted along the line $y = 55$. The fast oscillation in the left figure is clearly visible.

For each $r > 0$, we write $U_r := (0, r)^2$. We aim to compute the solution to the continuous partial differential equation in (1.1) with $w = 0$ (null Dirichlet boundary condition) and load function $f = 1$. We discretize this problem using a first-order finite element method. Let \mathcal{T} be a triangular mesh of the domain U_r constructed by first dividing each cell $z + [0, 1)^2$ ($z \in \mathbb{Z}^2$) into two triangles, and then using three levels of uniform mesh refinement. This results into a sufficiently fine mesh to capture the oscillations present in the exact solution u . The first order finite

element space

$$V_h := \{u \in H_0^1(U_r) \mid u|_K \in P^1(K) \quad \forall K \in \mathcal{T}\}$$

with standard nodal basis is used in all computations. The finite element solution $u_h \in V_h$ satisfies

$$(4.1) \quad \forall v_h \in V_h, \quad (\mathbf{a} \nabla u_h, \nabla v_h) = (f, v_h).$$

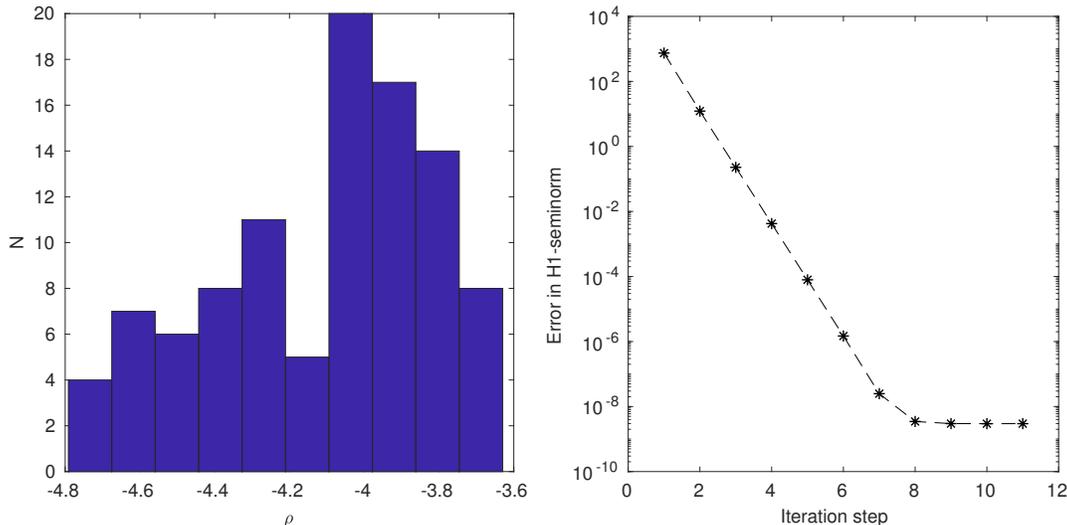


FIGURE 3. On the left, the empirical distribution of the factor ρ for $\lambda = 0.1$ and $r = 100$, based on 100 runs. On the right, the error in the H^1 seminorm for $r = 100$ and $\lambda = 0.1$, after each iteration. The method converges after 8 iterations.

A typical realization of the coefficient field $\mathbf{a}(x)$ and of the corresponding exact solution u are visualized in Figure 1, with the choice of $r = 100$. The high-frequency oscillations in the solution are clearly visible in Figure 2, where the solution is visualized along the line $y = 55$.

Our interest lies in the contraction factor of the iterative procedure. The contraction factor is studied by first solving the finite dimensional problem (4.1) exactly using a direct solver. Then a sequence of approximate solutions $\{u_h^{(i)}\}_{i=1}^N$ is generated by starting from $u_h^{(1)} = 0$ and applying the iterative procedure described in Theorem 1.1. The logarithm of the error $\|\nabla(u - u_h^{(i)})\|_{L^2(U_r)}$ is computed for each $i \in \{1, \dots, 10\}$, a regression line is fitted, and the slope of this line is denoted by ρ . It is our numerical estimate of the logarithm of the contraction factor; roughly speaking,

$$\rho \approx \log \left(\frac{\|\nabla(u_h - u_h^{(i+1)})\|_{L^2(U_r)}}{\|\nabla(u_h - u_h^{(i)})\|_{L^2(U_r)}} \right)$$

(“log” denotes the natural logarithm.) The iteration is said to converge when the relative error is smaller than 10^{-9} . Past this threshold, the error between the exact and the iterative solutions is smaller than the accuracy of the discretization itself, and thus cannot be measured.

Since the coefficient field is random, the contraction factor will vary for different realizations of \mathbf{a} . For the choice of $\lambda = 0.1$ and $r = 100$, the empirical distribution of

the contraction factor is given in Figure 3, based on one hundred samples of the coefficient field. Apart from the purposes of displaying this histogram, each of our estimates for ρ is an average over ten realizations of the coefficient field.

In our first test, the parameter λ is fixed to $\lambda = 0.1, 0.2$, and then 0.4 . The size of the domain r is varied between 10 and 200. The averaged contraction factor is visualized on the left side of Figure 4. The results are in excellent agreement with Theorem 1.1. After a pre-asymptotic region, the contraction factor becomes independent of the size of the domain r . The pre-asymptotic region is due to the fact that for small values of r , the pre- and post-smoothing steps are essentially sufficient to solve the equation. We emphasize that the contraction factor remains very good, of the order of 0.1, even for the relatively large value of $\lambda = 0.4$.

In the second test, the size of the domain r takes values $r = 100, 200$, and 300 , while λ is varied between 0.01 and 0.5. For each λ , the exponent of the averaged contraction factor is computed based on ten simulation runs. The results are presented on the right side of Figure 4. After a pre-asymptotic region, the exponential of the contraction factor behaves like $\lambda^{1/2}$, as predicted by Theorem 1.1. The pre-asymptotic region is roughly characterized by the scaling $r \lesssim 10\lambda^{-1}$.

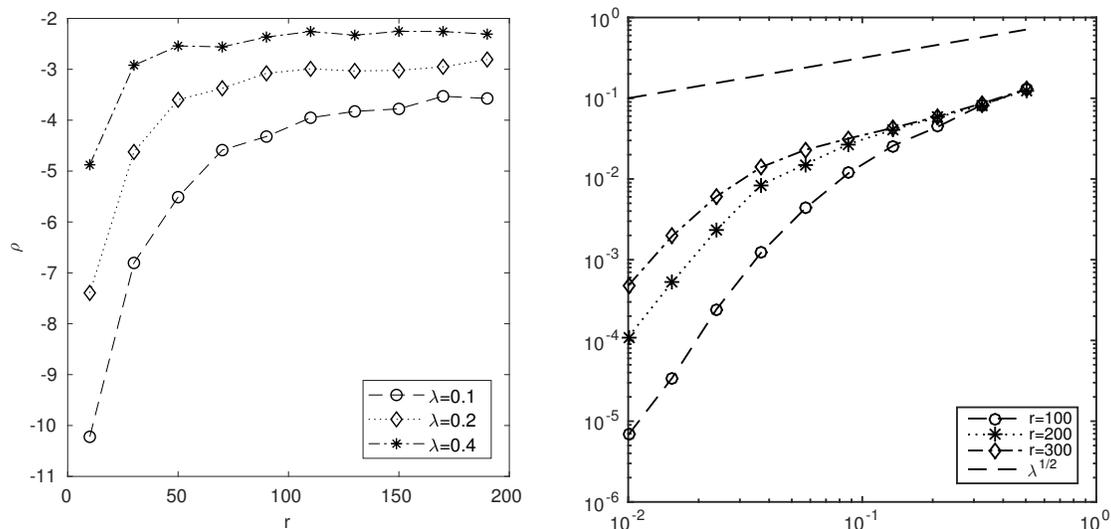


FIGURE 4. On the left, averaged factor ρ as a function of r , for $\lambda = 0.1, 0.2$, and 0.4 . On the right, the exponential of the averaged factor ρ as a function of λ for $r = 100, 200$, and 300 . In all cases, the average is computed from ten simulation runs.

APPENDIX A. SOBOLEV ESTIMATES

In this appendix, we prove an estimate for the norm of a function restricted to a layer close to the boundary of a domain. The estimate is an integrated version of a trace theorem. For convenience, we will also recall a standard H^2 estimate for homogeneous elliptic equations. As in (3.15), for every $\ell \geq 0$, we write

$$U_{r,\ell} := \{x \in U_r : \text{dist}(x, \partial U_r) > \ell\}.$$

Proposition A.1 (Trace estimate). *There exists $C(U, d) < \infty$ such that for every $r \geq 1$, $\ell \in (0, r]$ and $f \in H^1(U_r)$,*

$$r^{-d} \|f\|_{L^2(U_r \setminus U_{r,\ell})}^2 \leq C \ell \|f\|_{\underline{L}^2(U_r)} \|f\|_{\underline{H}^1(U_r)}.$$

Proof. Denote by $\mathbf{n}_{r,t}$ the unit normal vector to $\partial U_{r,\ell}$, which we extend to $U_{r,\ell}$ harmonic continuation. Since U is $C^{1,1}$, there exists $C(U, d) < \infty$ such that for every $t \in (0, r/C]$, we have $\|\nabla \mathbf{n}_{r,t}\|_{L^\infty(U_{r,t})} \leq Cr^{-1}$. It thus follows that

$$\begin{aligned} \int_{\partial U_{r,t}} f^2 &= \int_{U_{r,t}} \nabla \cdot (f^2 \mathbf{n}_{r,t}) \\ &\leq Cr^{-1} \|f\|_{L^2(U_r)}^2 + C \|f\|_{L^2(U_r)} \|\nabla f\|_{L^2(U_r)} \\ &\leq Cr^d \|f\|_{\underline{L}^2(U_r)} \|f\|_{\underline{H}^1(U_r)}. \end{aligned}$$

By the coarea formula, for every $\ell \leq r/C$, we have

$$\|f\|_{L^2(U_r \setminus U_{r,\ell})}^2 = \int_0^\ell \|f\|_{L^2(\partial U_{r,t})}^2 dt.$$

Combining the previous two displays yields

$$\|f\|_{L^2(U_r \setminus U_{r,\ell})}^2 \leq C \ell r^d \|f\|_{\underline{L}^2(U_r)} \|f\|_{\underline{H}^1(U_r)},$$

which is the announced result. The case $\ell > r/C$ is immediate. \square

Proposition A.2 (H^2 estimate). *Let $\bar{\mathbf{a}} \in \mathbb{R}_{\text{sym}}^{d \times d}$ satisfy (3.1). There exists a constant $C(\Lambda, U, d) < \infty$ such that for every $u \in H_0^1(U_r)$ and $f \in L^2(U_r)$, if*

$$-\nabla \cdot \bar{\mathbf{a}} \nabla u = f,$$

then $u \in H^2(U_r)$ and

$$\|u\|_{\underline{H}^2(U_r)} \leq C \|f\|_{\underline{L}^2(U_r)}.$$

Proof. See [14, Theorem 6.3.2.4]. \square

Acknowledgments. SA was partially supported by the NSF Grant DMS-1700329. AH was partially supported by the Stenbäck stiftelse. TK was supported by the Academy of Finland. JCM was partially supported by the ANR Grant LSD (ANR-15-CE40-0020-03).

REFERENCES

- [1] A. Abdulle, W. E, B. Engquist, and E. Vanden-Eijnden. The heterogeneous multiscale method. *Acta Numer.*, 21:1–87, 2012.
- [2] S. Armstrong and P. Dario. Elliptic regularity and quantitative homogenization on percolation clusters. *Comm. Pure Appl. Math.*, 71(9):1717–1849, 2018.
- [3] S. Armstrong, T. Kuusi, and J.-C. Mourrat. *Quantitative stochastic homogenization and large-scale regularity*. Preliminary version available at arXiv:1705.05300.
- [4] I. Babuska and R. Lipton. Optimal local approximation spaces for generalized finite element methods with application to multiscale problems. *Multiscale Model. Simul.*, 9(1):373–406, 2011.
- [5] M. Bebendorf and W. Hackbusch. Existence of \mathcal{H} -matrix approximants to the inverse FE-matrix of elliptic operators with L^∞ -coefficients. *Numer. Math.*, 95(1):1–28, 2003.
- [6] A. Brandt. Multiscale scientific computation: review 2001. In *Multiscale and multiresolution methods*, volume 20 of *Lect. Notes Comput. Sci. Eng.*, pages 3–95. Springer, Berlin, 2002.
- [7] P. Dario. Optimal corrector estimates on percolation clusters, preprint, arXiv:1805.00902.
- [8] W. E. *Principles of multiscale modeling*. Cambridge University Press, Cambridge, 2011.

- [9] Y. Efendiev, J. Galvis, and T. Y. Hou. Generalized multiscale finite element methods (GMsFEM). *J. Comput. Phys.*, 251:116–135, 2013.
- [10] Y. Efendiev and T. Y. Hou. *Multiscale finite element methods*, volume 4 of *Surveys and Tutorials in the Applied Mathematical Sciences*. Springer, New York, 2009.
- [11] A.-C. Egloffé, A. Gloria, J.-C. Mourrat, and T. N. Nguyen. Random walk in random environment, corrector equation and homogenized coefficients: from theory to numerics, back and forth. *IMA J. Numer. Anal.*, 35(2):499–545, 2015.
- [12] B. Engquist and E. Luo. New coarse grid operators for highly oscillatory coefficient elliptic problems. *J. Comput. Phys.*, 129(2):296–306, 1996.
- [13] B. Engquist and E. Luo. Convergence of a multigrid method for elliptic equations with highly oscillatory coefficients. *SIAM J. Numer. Anal.*, 34(6):2254–2273, 1997.
- [14] L. C. Evans. *Partial differential equations*, volume 19 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, second edition, 2010.
- [15] A. Gholami, D. Malhotra, H. Sundar, and G. Biros. FFT, FMM, or multigrid? A comparative study of state-of-the-art Poisson solvers for uniform and nonuniform grids in the unit cube. *SIAM J. Sci. Comput.*, 38(3):C280–C306, 2016.
- [16] A. Gloria. Numerical approximation of effective coefficients in stochastic homogenization of discrete elliptic equations. *ESAIM Math. Model. Numer. Anal.*, 46(1):1–38, 2012.
- [17] L. Grasedyck, I. Greff, and S. Sauter. The AL basis for the solution of elliptic problems in heterogeneous media. *Multiscale Model. Simul.*, 10(1):245–258, 2012.
- [18] M. Griebel and S. Knapek. A multigrid-homogenization method. In *Modeling and computation in environmental sciences (Stuttgart, 1995)*, volume 59 of *Notes Numer. Fluid Mech.*, pages 187–202. Friedr. Vieweg, Braunschweig, 1997.
- [19] C. Gu. Uniform estimate of an iterative method for elliptic problems with rapidly oscillating coefficients, preprint, arXiv:1807.06565.
- [20] Q. Han and F. Lin. *Elliptic partial differential equations*, volume 1 of *Courant Lecture Notes in Mathematics*. American Mathematical Society, Providence, RI, 1997.
- [21] I. G. Kevrekidis, C. W. Gear, J. M. Hyman, P. G. Kevrekidis, O. Runborg, and C. Theodoropoulos. Equation-free, coarse-grained multiscale computation: enabling microscopic simulators to perform system-level analysis. *Commun. Math. Sci.*, 1(4):715–762, 2003.
- [22] S. Knapek. Matrix-dependent multigrid homogenization for diffusion problems. *SIAM J. Sci. Comput.*, 20(2):515–533, 1998.
- [23] R. Kornhuber and H. Yserentant. Numerical homogenization of elliptic multiscale problems by subspace decomposition. *Multiscale Model. Simul.*, 14(3):1017–1036, 2016.
- [24] A. Målqvist and D. Peterseim. Localization of elliptic multiscale problems. *Math. Comp.*, 83(290):2583–2603, 2014.
- [25] A.-M. Matache and C. Schwab. Two-scale FEM for homogenization problems. *M2AN Math. Model. Numer. Anal.*, 36(4):537–572, 2002.
- [26] J.-C. Mourrat. Efficient methods for the estimation of homogenized coefficients, preprint, arXiv:1609.06674.
- [27] N. Neuss, W. Jäger, and G. Wittum. Homogenization and multigrid. *Computing*, 66(1):1–26, 2001.
- [28] H. Owhadi. Multigrid with rough coefficients and multiresolution operator decomposition from hierarchical information games. *SIAM Rev.*, 59(1):99–149, 2017.
- [29] H. Owhadi, L. Zhang, and L. Berlyand. Polyharmonic homogenization, rough polyharmonic splines and sparse super-localization. *ESAIM Math. Model. Numer. Anal.*, 48(2):517–552, 2014.
- [30] G. C. Papanicolaou and S. R. S. Varadhan. Boundary value problems with rapidly oscillating random coefficients. In *Random fields, Vol. I, II (Esztergom, 1979)*, volume 27 of *Colloq. Math. Soc. János Bolyai*, pages 835–873. North-Holland, Amsterdam, 1981.
- [31] K. Stüben. A review of algebraic multigrid. *J. Comput. Appl. Math.*, 128(1-2):281–309, 2001.
- [32] H. Sundar, G. Biros, C. Burstedde, J. Rudi, O. Ghattas, and G. Stadler. Parallel geometric-algebraic multigrid on unstructured forests of octrees. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. IEEE Computer Society Press, article No. 43, 2012.

(S. N. Armstrong) COURANT INSTITUTE OF MATHEMATICAL SCIENCES, NEW YORK UNIVERSITY, USA

E-mail address: `scotta@cims.nyu.edu`

(A. Hannukainen) DEPARTMENT OF MATHEMATICS AND SYSTEMS ANALYSIS, AALTO UNIVERSITY, FINLAND

E-mail address: `antti.hannukainen@aalto.fi`

(T. Kuusi) DEPARTMENT OF MATHEMATICS AND STATISTICS, UNIVERSITY OF HELSINKI, FINLAND

E-mail address: `tuomo.kuusi@helsinki.fi`

(J.-C. Mourrat) DMA, ECOLE NORMALE SUPÉRIEURE, CNRS, PSL RESEARCH UNIVERSITY, PARIS, FRANCE

E-mail address: `mourrat@dma.ens.fr`