

## SPARSE COMPRESSION OF EXPECTED SOLUTION OPERATORS\*

MICHAEL FEISCHL<sup>†</sup> AND DANIEL PETERSEIM<sup>‡</sup>

**Abstract.** We show that the expected solution operator of prototypical linear elliptic PDEs with random coefficients is well approximated by a computable sparse matrix. This result is based on a random localized orthogonal multiresolution decomposition of the solution space that allows both the sparse approximate inversion of the random operator represented in this basis as well as its stochastic averaging. The approximate expected solution operator can be interpreted in terms of classical Haar wavelets. When combined with a suitable sampling approach for the expectation, this construction leads to an efficient method for computing a sparse representation of the expected solution operator.

**Key words.** Monte Carlo, random PDEs, expected solution, sparse approximation, uncertainty quantification

**AMS subject classifications.** 65N30, 65C05

**DOI.** 10.1137/20M132571X

**1. Introduction.** For a random (or parameterized) family of prototypical linear elliptic partial differential operators  $\mathcal{A}(\omega) = -\operatorname{div}(\mathbf{A}(\omega)\nabla\bullet)$  and a given deterministic right-hand side  $f$ , we consider the family of solutions

$$\mathbf{u}(\omega) := \mathcal{A}(\omega)^{-1}f$$

with events  $\omega \in \Omega$  in some probability space  $\Omega$ . We define the harmonically averaged operator

$$\mathcal{A} := (\mathbb{E}[\mathcal{A}(\omega)^{-1}])^{-1}.$$

The idea behind this definition is that  $\mathbb{E}(\mathbf{u})$  satisfies

$$\mathbb{E}[\mathbf{u}] = \mathcal{A}^{-1}f.$$

In this sense,  $\mathcal{A}$  may be understood as a stochastically homogenized operator, and  $\mathcal{A}^{-1}$  is the effective solution operator. Note that this definition does not rely on probabilistic structures of the random diffusion coefficient  $\mathbf{A}$ , such as stationarity, ergodicity, or any characteristic length of correlation. However, we shall emphasize that  $\mathcal{A}$  does not coincide with the partial differential operator that would result from the standard theory of stochastic homogenization (under stationarity and ergodicity) [31, 36, 45]; see, e.g., [5, 21, 12, 22, 1] for quantitative results. Recent works on discrete random problems on  $\mathbb{Z}^d$  with independent and identically distributed (i.i.d.) edge conductivities indicate that  $\mathcal{A}$  is rather a nonlocal perturbation of the Laplacian

---

\*Received by the editors March 16, 2020; accepted for publication (in revised form) July 31, 2020; published electronically November 4, 2020.

<https://doi.org/10.1137/20M132571X>

**Funding:** The work of the first author was supported by the Deutsche Forschungsgemeinschaft (DFG) grant CRC 1173. The work of the second author was supported by the DFG in the Priority Program 1748 “Reliable Simulation Techniques in Solid Mechanics” grant PE2143/2-2.

<sup>†</sup>Institute for Analysis and Scientific Computing, Vienna University of Technology, Vienna 1040, Austria (Michael.Feischl@tuwien.ac.at).

<sup>‡</sup>Institut für Mathematik, Universität Augsburg, Augsburg 86159, Germany (daniel.peterseim@math.uni-augsburg.de).

by a convolution-type operator [4, 28, 11], which indicates that there might exist an efficient algorithm to approximate the operator.

The goal of the present work is to show the following result, even in the more general PDE setup of this paper without any assumptions on the distribution of the random coefficient.

**MAIN RESULT** (see Theorem 10 for the precise statement). *We can compute a sparse matrix  $R^\delta$  in almost linear cost which approximates the expected solution operator  $\mathcal{A}^{-1}$  in the sense that the corresponding operator  $\mathcal{R}^\delta$  satisfies*

$$\|\mathcal{A}^{-1} - \mathcal{R}^\delta\|_{L^2(D) \rightarrow L^2(D)} \leq \delta$$

for any  $\delta > 0$  the number of nonzero entries of  $R^\delta$  as well as the computational cost to evaluate  $\mathcal{R}^\delta$  scale like  $\delta^{-d}$  up to logarithmic-in- $\delta$  terms.

The sparse matrix representation of  $\mathcal{A}^{-1}$  is based on multiresolution decompositions of the energy space in the spirit of numerical homogenization by localized orthogonal decomposition (LOD) [32, 25, 37, 18, 29, 19] and, in particular, its multiscale generalization that is popularized under the name gamblets [33]. In this paper, a one-to-one correspondence of a gamblet decomposition and classical Haar wavelets is established via  $L^2$ -orthogonal projections and conversely by corrections involving the solution operator (see section 2). The resulting problem-dependent multiresolution decompositions block-diagonalize the random operator  $\mathcal{A}$  for any event in the probability space (see section 3). The block-diagonal representations (with sparse blocks) are well conditioned and, hence, easily inverted to high accuracy using a few steps of standard linear iterative solvers. The sparsity of the inverted blocks is preserved to the degree that it deteriorates only logarithmically with higher accuracy.

While the sparsity pattern of the inverted block-diagonal operator is independent of the stochastic parameter and, hence, not affected when taking the expectation (or any sample mean), the resulting object cannot be interpreted in a known basis. This issue is circumvented by reinterpreting the approximate inverse stiffness matrices in terms of the deterministic Haar basis before stochastic averaging (see section 4). This leads to an accurate representation of  $\mathcal{A}^{-1}$  in terms of piecewise constant functions. Sparsity is not directly preserved by this transformation but can be retained by some appropriate hyperbolic cross truncation which is justified by scaling properties of the multiresolution decomposition (see section 5).

While the mathematical question of sparse approximability of the expected operator can be answered positively with simpler techniques, the above construction leads to a computationally efficient method for approximating  $\mathcal{A}^{-1}$  when combined with any sampling approach for the approximation of the expectation (see section 6). This new sparse compression algorithm for the direct discretization of  $\mathcal{A}^{-1}$  may be beneficial if we want to compute  $\mathbb{E}[\mathbf{u}]$  for multiple right-hand sides  $f$ . This, for example, is the case if we have an independent probability space  $\xi \in \Xi$  influencing  $f = \mathbf{f}(\xi)$  as well as the corresponding solution  $\mathbf{U}(\omega, \xi) := \mathcal{A}(\omega)^{-1} \mathbf{f}(\xi)$ . Then we might be interested in the average behavior  $\mathbb{E}_{\Omega \times \Xi}[\mathbf{U}]$ , which is the solution of

$$(1.1) \quad \mathbb{E}_{\Omega \times \Xi}[\mathbf{U}] = \mathbb{E}_{\Xi}[\mathcal{A}^{-1} \mathbf{f}] = \mathcal{A}^{-1} \mathbb{E}_{\Xi}[\mathbf{f}].$$

While this can be computed efficiently with sparse approximations of the random parameter (see, e.g., [2, 3]) or multilevel algorithms (see, e.g., [9, 20]) under the regularity assumption on the random parameter, the present approach does not assume any smoothness apart from integrability. A practical example for the problem might serve the Darcy flow as a model of groundwater flow. Here,  $\mathcal{A}$  is a random diffusion process modeling the unknown diffusion coefficient of the ground material. The

right-hand side  $\mathbf{f}$  would be the random (unknown) injection of pollutants into the groundwater. Ultimately, the user would be interested in the average distribution of pollutants in the ground. Obviously, computing the right-hand side of (1.1) requires the user to sample  $\Omega$  and  $\Xi$  successively, whereas computing the left-hand side of (1.1) forces the user to sample the much larger product space  $\Omega \times \Xi$ . While for plain Monte Carlo sampling only the possibly increased variance of the product random variable affects the convergence, higher-order sampling methods such as sparse grids and quasi-Monte Carlo will directly and, in case of lack of regularity on the random parameter, quite drastically (see, e.g., exponential dependence on dimension in [7]) benefit from the reduction of dimension of the probability space. Therefore, an accurate discretization of  $\mathcal{A}$  can help save significant computational cost.

We consider some prototypical linear second-order elliptic PDEs with a random diffusion coefficient. Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space with set of events  $\Omega$ ,  $\sigma$ -algebra  $\mathcal{F} \subseteq 2^\Omega$ , and probability measure  $\mathbb{P}$ . The expectation operator is denoted by  $\mathbb{E}$ . Let  $D \subseteq \mathbb{R}^d$  for  $d \in \{1, 2, 3\}$  be a bounded Lipschitz polytope with diameter of order 1. The set of admissible coefficients reads

$$\mathcal{M}(D, \gamma_{\min}, \gamma_{\max}) = \left\{ A \in L^\infty(D; \mathbb{R}_{\text{sym}}^{d \times d}) \text{ s.t. } \gamma_{\min} |\xi|^2 \leq (A(x)\xi) \cdot \xi \leq \gamma_{\max} |\xi|^2 \right. \\ \left. \text{for a.e. } x \in D \text{ and all } \xi \in \mathbb{R}^d \right\}$$

for given uniform spectral bounds  $0 < \gamma_{\min} \leq \gamma_{\max} < \infty$ . Here,  $\mathbb{R}_{\text{sym}}^{d \times d}$  denotes the set of symmetric  $d \times d$  matrices. Let  $\mathbf{A}$  be a Bochner-measurable  $\mathcal{M}(D, \gamma_{\min}, \gamma_{\max})$ -valued random field with  $\gamma_{\max} > \gamma_{\min} > 0$ . Note that we do not make any structural assumptions regarding the distribution of  $\mathbf{A}$ . Moreover, realizations in  $\mathcal{M}(D, \gamma_{\min}, \gamma_{\max})$  are fairly free to vary within the bounds  $\gamma_{\min}$  and  $\gamma_{\max}$  without any conditions on frequencies of variation or smoothness.

Denote the energy space by  $V := H_0^1(D)$ , and let  $f \in V^* = H^{-1}(D)$  be deterministic. The prototypical second-order elliptic variational problem seeks a  $V$ -valued random field  $\mathbf{u}$  such that, for almost all  $\omega \in \Omega$ ,

$$(1.2) \quad \mathbf{a}_\omega(\mathbf{u}(\omega), v) := \int_D (\mathbf{A}(\omega)(x) \nabla \mathbf{u}(\omega)(x)) \cdot \nabla v(x) dx = f(v) \quad \text{for all } v \in V.$$

The bilinear form  $\mathbf{a}_\omega$  depends continuously on the coefficient  $\mathbf{A}(\omega) \in \mathcal{M}(D, \gamma_{\min}, \gamma_{\max})$  and, particularly, is measurable as a function of  $\omega$ . Hence, the reformulation of this problem in the Hilbert space  $L^2(\Omega; V)$  of  $V$ -valued random fields with finite second moments shows well-posedness in the sense that there exists a unique solution  $\mathbf{u} \in L^2(\Omega; V)$  with

$$\|\nabla \mathbf{u}\|_{L^2(\Omega; V)} := \int_\Omega \int_D |\nabla(\mathbf{u}(\omega))(x)|^2 dx d\mathbb{P}(\omega)^{1/2} \leq \gamma_{\min}^{-1} \|f\|_{V^*}.$$

To connect the model problem to the operator setting of the introduction, we shall define the random operator  $\mathcal{A}: \Omega \rightarrow \mathcal{L}(V, V^*)$  by

$$\langle \mathcal{A}(\omega)u, v \rangle_{V^*, V} := \mathbf{a}_\omega(u, v)$$

for functions  $u, v \in V$  and  $\omega \in \Omega$ . Then the model problem (1.2) can be rephrased as

$$\mathcal{A}(\omega)\mathbf{u}(\omega) = f \quad \text{for almost all } \omega \in \Omega.$$

For convenience, we define the sample-dependent energy norm  $\|\cdot\|_\omega^2 := \mathbf{a}_\omega(\cdot, \cdot)$ .

**2. Coefficient-adapted hierarchical bases.** Let  $\mathcal{T}_\ell$ ,  $\ell = 0, \dots, L$  denote a sequence of uniform refinements with mesh size  $h_\ell$  of some initial mesh  $\mathcal{T}_0$  of  $D$ , and let  $\mathcal{N}(\mathcal{T}_\ell)$  denote the nodes of the meshes. We allow fairly general meshes in the sense that we only require a reference element  $T_{\text{ref}}$  together with a family of uniformly bi-Lipschitz maps  $\Psi_T: T_{\text{ref}} \rightarrow T$  for all elements  $T \in \mathcal{T}_\ell$ ,  $\ell = 0, \dots, L$ . Straightforward examples are simplicial meshes generated from an initial triangulation by red refinement (or newest vertex bisection) or quadrilateral meshes generated by subdividing the elements into  $2^d$  new elements. Particularly, hanging nodes do not pose problems as long as the other properties are observed.

The number of levels (or scales)  $L$  will typically be chosen proportional to the modulus of some logarithm of the desired accuracy  $1 \gtrsim \delta > 0$ . We assume  $h_{\ell+1} \leq h_\ell/2$ . Note that any other fixed factor of mesh width reduction strictly smaller than one would do the job. Define the set of descendants of an element  $T \in \mathcal{T}_\ell$  by  $\text{ref}(T) := \{T' \in \mathcal{T}_{\ell+1} : T' \subseteq T\}$ . For each  $T \in \bigcup_{\ell=0}^{L-1} \mathcal{T}_\ell$ , we pick piecewise constant functions  $\phi_{T,1}, \phi_{T,2}, \dots, \phi_{T, \#\text{ref}(T)} \in P^0(\text{ref}(T))$  such that they are pairwise  $L^2(T)$ -orthogonal and  $\int_T \phi_{T,j} dx = 0$  for all  $j = 1, \dots, \#\text{ref}(T)$ . With the indicator functions  $\chi_{(\cdot)}$ , we then define  $\mathcal{H}_0 := \{\chi_T : T \in \mathcal{T}_0\}$  and for  $\ell \geq 1$

$$(2.1) \quad \mathcal{H}_\ell := \bigcup_{T \in \mathcal{T}_{\ell-1}} \{\phi_{T,j} : j = 1, \dots, \#\text{ref}(T)\}.$$

We define a Haar basis via

$$\mathcal{H} := \bigcup_{\ell=0}^L \mathcal{H}_\ell.$$

LEMMA 1. *The basis  $\mathcal{H}$  is  $L^2$ -orthogonal and local in the sense that  $\phi \in \mathcal{H}_\ell$  satisfies  $\text{supp}(\phi) = T$  for some  $T \in \mathcal{T}_{\ell-1}$  or  $T \in \mathcal{T}_0$  for  $\ell = 0$ .*

*Proof.* Let  $\phi_\ell \in \mathcal{H}_\ell$  and  $\phi_k \in \mathcal{H}_k$ . If  $k = \ell$ , then the interiors of the supports of any  $\phi_k = \phi_\ell \in \mathcal{H}_k$  are disjoint, which implies  $L^2(D)$ -orthogonality. If  $k < \ell$ , we have that  $\phi_k$  is constant on  $\text{supp}(\phi_\ell)$ . Since  $\int_D \phi_\ell dx = 0$  by definition, this concludes the proof of  $L^2$ -orthogonality. Locality follows readily from the construction.  $\square$

*Remark 2.* For uniform Cartesian meshes,  $\mathcal{H}$  is the Haar basis. The choice of the  $2^d - 1$  generating functions follows the standard procedure for Haar wavelets (see, e.g., [42]). The construction is applicable to general meshes that are not based on tensor-product structures.

Due to the lack of  $V$ -conformity, the basis  $\mathcal{H}$  is not suited for approximating the solution of model problem (1.2) in a Galerkin approach. It will, however, serve as a companion of certain regularized hierarchical bases  $\mathcal{B}(\omega) = \bigcup_{\ell=0}^L \mathcal{B}_\ell(\omega) \subset V$  to be defined below. The new bases are connected to  $\mathcal{H}$  (and to each other) via  $L^2$ -orthogonal projections  $\Pi_\ell: V \rightarrow P^0(\mathcal{T}_\ell)$  onto  $\mathcal{T}_\ell$ -piecewise constant functions by

$$(2.2) \quad \Pi_\ell \mathcal{B}_\ell(\omega) = \mathcal{H}_\ell$$

for all  $\ell = 0, 1, \dots, L$  and  $\omega \in \Omega$ . Among the infinitely many possible choices, we define the elements of  $\mathcal{B}_\ell(\omega)$  by minimizing the energies  $\frac{1}{2} \mathbf{a}_\omega(\bullet, \bullet)$  in the closed affine space of preimages of  $\Pi_\ell$  restricted to  $V$ ; i.e., given  $\phi \in \mathcal{H}_\ell$  and  $\omega \in \Omega$ , we define  $\mathbf{b}_\phi(\omega) \in \mathcal{B}_\ell(\omega)$  by

$$(2.3) \quad \mathbf{b}_\phi(\omega) := \underset{v \in V}{\text{argmin}} \frac{1}{2} \mathbf{a}_\omega(v, v) \quad \text{subject to} \quad \Pi_\ell v = \phi.$$

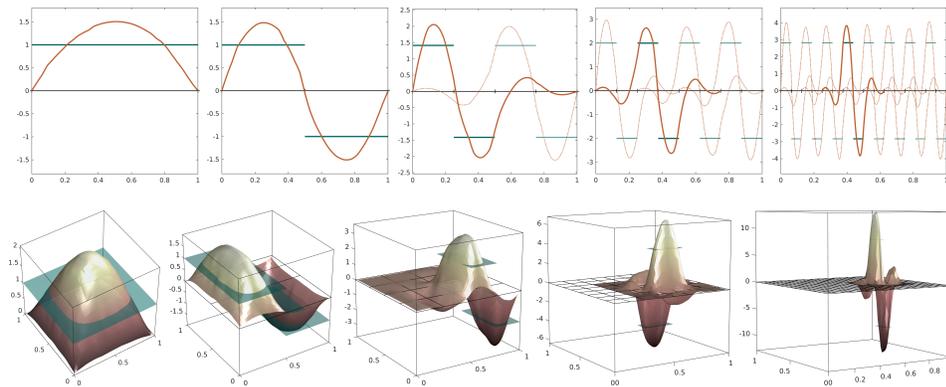


FIG. 1. Realization of coefficient-adapted hierarchical decomposition in one dimension (top row) and two dimensions (bottom row) based on an elliptic partial differential operator with random coefficient *i.i.d.* with respect to Cartesian grid of width  $\varepsilon = 2^{-6}$ . Five levels from coarse (left) to fine (right). Green lines/surfaces represent classical Haar wavelets.

This construction (visualized in Figure 1) is strongly inspired by numerical homogenization, where this sort of orthogonalization of scales in the energy space paved the way to a scheme that works with arbitrary rough coefficients beyond periodicity or scale separation [32, 25, 37, 29]. While most results in the context of this so-called LOD are based on a conforming companion (the Faber basis), early works also addressed the possibility of using discontinuous companions [13, 14]. This dG version of LOD is very useful when taking the step from two levels or scales in numerical homogenization to actual multilevel decomposition. This was first shown in [33], where so-called gamblets are introduced; see also [41, 26, 27, 35, 34]. In particular, piecewise constants induce a natural hierarchical structure with nested kernels of local projection operators (here the  $\Pi_\ell$ ) that is not easily achieved with  $H^1$ -conforming functions. The construction of the present paper coincides with the gamblet decomposition of [33] in the sense that the approximation spaces on all levels coincide in some idealized deterministic setting. However, our particular choice of basis is connected to the Haar wavelets which decouple the definition and computation of the basis across levels. More importantly, our particular choice of basis is crucial in the context of the random problem at hand because it is exactly the link to the deterministic Haar basis that allows a meaningful interpretation of the averaged approximate solution operator.

We shall express the mapping of bases encoded in (2.2)–(2.3) in terms of two concatenated linear operators. This will be useful for both analysis and actual computations. First, let  $\tilde{\Pi}_\ell : L^2(D) \rightarrow V$  be such that

$$(2.4) \quad \Pi_\ell \circ \tilde{\Pi}_\ell = \text{id} \quad \text{on } \text{span}\mathcal{H}_\ell.$$

In particular, this means that  $\tilde{\Pi}_\ell$  maps any  $\phi \in \mathcal{H}_\ell$  to some function that is admissible in the sense of the minimization problem (2.3). The operators  $\tilde{\Pi}_\ell$  are easily constructed using nonnegative bubble functions  $\tilde{\chi}_T$  supported on an element  $T \in \mathcal{T}_\ell$  with  $\Pi_\ell \tilde{\chi}_T = \chi_T$ . Then

$$\tilde{\Pi}_\ell v = \sum_{T \in \mathcal{T}_\ell} (\Pi_\ell v)|_T \tilde{\chi}_T.$$

There is even locality in the sense of

$$(2.5) \quad \text{supp } \tilde{\Pi}_\ell \phi \subset \text{supp } \phi$$

for all  $\phi \in \text{span}\mathcal{H}_\ell$ . The bubbles can be chosen such that, for some  $C > 0$ ,

$$(2.6) \quad \|\tilde{\Pi}_\ell \phi\|_{H^m(D)} \leq Ch^{-m} \|\phi\|_{L^2(D)}$$

holds for  $m \in \{0, 1\}$ .

The second step involves  $\mathbf{a}_\omega$ -orthogonal projections  $\mathcal{C}_\ell(\omega)$  onto the closed subspaces

$$(2.7) \quad W_\ell := \text{kernel}(\Pi_\ell|_V) = \text{kernel}(\tilde{\Pi}_\ell|_V)$$

of  $V$ . Given any  $u \in V$ , define  $\mathcal{C}_\ell(\omega)u \in W_\ell$  as the unique solution of the variational problem

$$(2.8) \quad \mathbf{a}_\omega(\mathcal{C}_\ell(\omega)u, v) = \mathbf{a}_\omega(u, v) \quad \text{for all } v \in W_\ell.$$

With the two operators  $\tilde{\Pi}_\ell$  and  $\mathcal{C}_\ell$  we rewrite (2.3) as

$$\mathbf{b}_\phi = (\text{id} - \mathcal{C}_\ell)\tilde{\Pi}_\ell \phi$$

for all  $\phi \in \mathcal{H}_\ell$  and  $\ell = 0, 1, \dots, L$ . Actually, for any  $\omega \in \Omega$ ,  $(\text{id} - \mathcal{C}_\ell(\omega))\tilde{\Pi}_\ell$  defines a bijection from  $\mathcal{H}$  to  $\mathcal{B}(\omega)$  with left inverse  $\Pi_\ell$ .

While the  $L^2$ -orthogonality of the Haar basis is not preserved under these mappings, we have achieved  $\mathbf{a}$ -orthogonality between the levels of the hierarchies as shown in the following lemma.

**LEMMA 3** ( *$\mathbf{a}$ -orthogonality and scaling of  $\mathcal{B}$* ). *Any two functions  $b_k \in \mathcal{B}_k(\omega)$  and  $b_\ell \in \mathcal{B}_\ell(\omega)$  with  $k \neq \ell$  satisfy*

$$\mathbf{a}_\omega(b_k, b_\ell) = 0.$$

Moreover,

$$(2.9) \quad C^{-1} \|\phi_k\|_{L^2(D)} \leq C^{-1} \|b_k\|_{L^2(D)} \leq h_k \|b_k\|_\omega \leq C \|\phi_k\|_{L^2(D)}$$

with some generic constant  $C > 0$  independent of the mesh sizes and the event.

*Proof.* Since

$$(2.10) \quad \Pi_k(\text{id} - \mathcal{C}_\ell(\omega))\tilde{\Pi}_\ell \mathcal{H}_\ell = \Pi_k \Pi_\ell(\text{id} - \mathcal{C}_\ell(\omega))\tilde{\Pi}_\ell \mathcal{H}_\ell = \Pi_k \mathcal{H}_\ell = \{0\}$$

whenever  $k < \ell$ , we have that

$$\mathcal{B}_\ell(\omega) \subset W_k.$$

This and the orthogonality

$$\mathbf{a}_\omega(\mathcal{B}_k(\omega), W_k) = 0$$

from (2.8) prove the (block-)orthogonality of the bases. The scaling follows from  $\Pi_{k-1}\mathcal{B}_k = \{0\}$  (which is a special instance of (2.10)), the Poincaré inequality, (2.6), and the construction. More precisely,

$$(2.11) \quad \begin{aligned} \|\phi_k\|_{L^2(D)} &= \|\Pi_k b_k\|_{L^2(D)} \leq \|b_k\|_{L^2(D)} = \|(1 - \Pi_{k-1})b_k\|_{L^2(D)} \lesssim h_k \|b_k\|_\omega \\ &= h_k \|(\text{id} - \mathcal{C}_k(\omega))\tilde{\Pi}_k \phi_k\|_\omega \leq h_k \|\tilde{\Pi}_k \phi_k\|_\omega \lesssim \|\phi_k\|_{L^2(D)}. \end{aligned}$$

This concludes the proof. □

We shall emphasize that, in general, the basis elements  $\mathbf{b}_\phi(\omega)$  have global support in  $D$ . However, their moduli decay exponentially away from  $\text{supp } \phi$  in scales of  $h_\ell$ ,

$$(2.12) \quad \|\mathbf{b}_\phi(\omega)\|_{H^1(D \setminus B_R(\text{supp } \phi))} \leq C e^{-cR/h_\ell} \|\mathbf{b}_\phi(\omega)\|_{H^1(D)},$$

with some generic constants  $c, C > 0$  that solely depend on the contrast  $\gamma_{\max}/\gamma_{\min}$  and the shape regularity of the mesh  $\mathcal{T}_\ell$  (and thus on  $\mathcal{T}_0$ ) but not on the mesh size. This is a well-established result of numerical homogenization since [32] and valid in many different settings (see [37] and references therein). Here, we will provide some elements of a more recent constructive proof of the decay that provides local approximations by the theory of preconditioned iterative solvers [29], which in turn is based on [30].

We start with introducing an overlapping decomposition of  $D$  that we will later use to define the local preconditioner. Let the level  $\ell \in \{0, 1, \dots, L\}$  and the event  $\omega \in \Omega$  be arbitrary but fixed. For any element of the mesh, define the patch

$$D_T := \bigcup \{K \in \mathcal{T}_\ell \mid \bar{K} \cap \bar{T} = \emptyset\}$$

and a corresponding local subspace

$$V_T := \{v \in V \mid v = 0 \text{ in } D \setminus D_T\} \subset V.$$

Note that  $V_T$  is equal to  $H_0^1(D_T)$  up to extension by zero outside of  $D_T$ . Let  $\lambda_T$ ,  $T \in \mathcal{T}_\ell$ , be a partition of unity with  $\text{supp } \lambda_T \subset D_T$  and  $\|\lambda_T\|_{W^{m,\infty}(D)} \lesssim h_\ell^{-m}$ ,  $m = 0, 1$ . Under the complementary projection  $(\text{id} - \tilde{\Pi}_\ell)$ , these subspaces are turned into subspaces

$$W_T := (\text{id} - \tilde{\Pi}_\ell)V_T = \{v \in W_\ell \mid v = 0 \text{ in } D \setminus D_T\}$$

of  $W_\ell$ . For each  $T \in \mathcal{T}_\ell$ , we define the corresponding  $\mathbf{a}_\omega$ -orthogonal projection  $\mathcal{P}_T(\omega): V \rightarrow W_T \subset W_\ell \subset V$  by the variational problem

$$\mathbf{a}_\omega(\mathcal{P}_T(\omega)u, w) = \mathbf{a}_\omega(u, w) \quad \text{for all } w \in W_T.$$

The sum of these local Ritz projections

$$(2.13) \quad \mathcal{P}_\ell(\omega) := \sum_{T \in \mathcal{T}_\ell} \mathcal{P}_T(\omega)$$

defines a bounded linear operator from  $V$  to  $W_\ell$  that can be seen as a preconditioned version of the correction operator  $\mathcal{C}_\ell(\omega)$ . The operator  $\mathcal{P}_\ell(\omega)$  is quasi-local with respect to the mesh  $\mathcal{T}_\ell$  since information can only propagate over distances of order  $h_\ell$  each time  $\mathcal{P}_\ell(\omega)$  is applied.

The remaining part of this section aims to show that the preconditioned operators  $\mathcal{P}_\ell(\omega)$  serve well within iterative solvers for linear equations. Following the abstract theory for subspace correction or additive Schwarz methods for operator equations [30] (see also [43, 44] for the matrix case), we need to verify that the energy norm of a function  $u \in V$  can be bounded in terms of the sum of local contributions from  $V_T$ , and, for one specific decomposition, we need a reverse estimate.

LEMMA 4. *For every decomposition  $u = \sum_{T \in \mathcal{T}_\ell} u_T$  of  $u \in W_\ell$  with  $u_T \in W_T$ , we have*

$$\|\nabla u\|_{L^2(D)}^2 \leq K_2 \sum_{z \in \mathcal{T}_\ell} \|\nabla u_T\|_{L^2(D)}^2$$

with constant  $K_2 > 0$  depending only on the shape regularity of  $\mathcal{T}_\ell$  (and thus on  $\mathcal{T}_0$ ). With the partition of unity functions  $\lambda_T$  associated with the elements  $T \in \mathcal{T}_\ell$ , the one decomposition  $\sum_{T \in \mathcal{T}_\ell} u_T = u$  with  $u_T := (1 - \tilde{\Pi}_\ell)(\lambda_T u) \in W_T$  for  $T \in \mathcal{T}_\ell$  satisfies

$$\sum_{T \in \mathcal{T}_\ell} \|\nabla u_T\|_{L^2(D)}^2 \leq K_1 \|\nabla u\|_{L^2(D)}^2$$

with constant  $K_1 > 0$  that only depends on the shape regularity of  $\mathcal{T}_\ell$  and the contrast  $\gamma_{\max}/\gamma_{\min}$ .

*Proof.* With  $K_2$  the maximum number of elements of  $\mathcal{T}_\ell$  covered by one patch  $D_T$  for  $T \in \mathcal{T}_\ell$ , we can estimate on a single element  $T'$

$$\|\nabla u\|_{L^2(T')}^2 = \sum_{T \in \mathcal{T}_\ell} \sum_{L^2(T')}^2 \|\nabla u_T\|_{L^2(T')}^2 \leq K_2 \sum_{T \in \mathcal{T}_\ell} \|\nabla u_T\|_{L^2(T')}^2.$$

Due to shape regularity of  $\mathcal{T}_\ell$ ,  $K_2$  is independent of  $h_\ell$ . A summation over all  $T'$  yields the first inequality. The second one follows from the  $H^1$ -stability of  $\tilde{\Pi}_\ell$  on  $W_\ell$ , the product rule, (2.6), and the Poincaré inequality. For further details, we refer the reader to [29, Lemma 3.1], where these results are proved in detail in a very similar setting. □

Lemma 4 implies that

$$(2.14) \quad 1/K_1 \mathbf{a}_\omega(v, v) \leq \mathbf{a}_\omega(\mathcal{P}_\ell(\omega)v, v) \leq K_2 \mathbf{a}_\omega(v, v)$$

holds for functions  $v$  in the kernel  $W_\ell$  of  $\Pi_\ell|_V$  and any  $\omega \in \Omega$  (cf. [29, equation (3.11)]). Following the construction of [30, 29], there exists a localized linear approximation  $\mathcal{C}_\ell^\delta(\omega)$  based on  $\mathcal{O}(\log(1/\delta))$  steps of some linear iterative solver applied to the preconditioned corrector problems [29, Equations (3.8) or equation (3.18)] such that

$$(2.15) \quad \|\nabla(\mathcal{C}_\ell(\omega)u - \mathcal{C}_\ell^\delta(\omega)u)\|_{L^2(D)} \leq \delta \|\nabla \mathcal{C}_\ell(\omega)u\|_{L^2(D)};$$

see [29, Lemma 3.2]. With the approximate correctors, we can define modified (localized) bases

$$\mathcal{B}^\delta(\omega) := \bigcup_{\ell=0}^L \mathcal{B}_\ell^\delta(\omega) := \bigcup_{\ell=0}^L \{\mathbf{b}_\phi^\delta(\omega) : \phi \in \mathcal{H}_\ell\},$$

where

$$\mathbf{b}_\phi^\delta(\omega) := (\text{id} - \mathcal{C}_\ell^\delta(\omega))\tilde{\Pi}_\ell\phi$$

for  $\phi \in \mathcal{H}_\ell$ . The previous discussion shows that there exist constants  $C_1, C_2 > 0$  that only depend on the shape regularity of the meshes  $\mathcal{T}_\ell$  and the contrast  $\gamma_{\max}/\gamma_{\min}$  of the coefficients such that

$$(2.16) \quad \|\mathbf{b}_\phi(\omega) - \mathbf{b}_\phi^\delta(\omega)\|_\omega \leq C_1 \delta \|\mathbf{b}_\phi(\omega)\|_\omega,$$

while

$$(2.17) \quad \text{supp } \mathbf{b}_\phi^\delta(\omega) \subset \{x \in D : \text{dist}(x, \text{supp } \phi) \leq C_2 |\log(\delta)| h_\ell\}.$$

Later, we will typically use normalized bases. Since

$$(2.18) \quad (1 - C_1\delta)\|\mathbf{b}_\phi(\omega)\|_\omega \leq \|\mathbf{b}_\phi^\delta(\omega)\|_\omega \leq (1 + C_1\delta)\|\mathbf{b}_\phi(\omega)\|_\omega$$

by (2.16), the normalization of the localized bases is meaningful whenever  $\delta < 1/C_1$ . Normalization does not affect the local supports (2.17), and the approximation property (2.16) is preserved in the sense of

$$(2.19) \quad \begin{aligned} \left\| \frac{\mathbf{b}_\phi(\omega)}{\|\mathbf{b}_\phi(\omega)\|_\omega} - \frac{\mathbf{b}_\phi^\delta(\omega)}{\|\mathbf{b}_\phi^\delta(\omega)\|_\omega} \right\|_\omega &\leq \frac{\|\mathbf{b}_\phi(\omega) - \mathbf{b}_\phi^\delta(\omega)\|_\omega}{\|\mathbf{b}_\phi(\omega)\|_\omega} + \frac{\|\mathbf{b}_\phi(\omega)\|_\omega - \|\mathbf{b}_\phi^\delta(\omega)\|_\omega}{\|\mathbf{b}_\phi^\delta(\omega)\|_\omega} \\ &\leq \delta C_1 + \frac{C_1\delta}{1 - C_1\delta} \leq 3\delta C_1 \end{aligned}$$

for any  $\delta \leq 1/(2C_1)$ .

**3. Sparse stiffness matrices.** With the localized bases of the previous section, we can now study the sparsity of corresponding stiffness matrices and their inverses. We define the level function  $\text{lev}(\cdot)$  according to the Haar basis by  $\text{lev}(\mathbf{b}) = \text{lev}(\mathbf{b}^\delta) = \text{lev}(\phi) = \ell$  for  $\mathbf{b} = \mathbf{b}_\phi \in \mathcal{B}_\ell(\omega)$ ,  $\mathbf{b}^\delta = \mathbf{b}_\phi^\delta \in \mathcal{B}_\ell^\delta(\omega)$ , and  $\phi \in \mathcal{H}_\ell$ . We order the basis functions in  $\mathcal{B}$ ,  $\mathcal{B}^\delta$ , and  $\mathcal{H}$  such that  $\text{lev}$  is monotonically increasing in the index running from 1 to  $N := \#\mathcal{B} = \#\mathcal{B}^\delta = \#\mathcal{H}$ . With this convention, we may also write  $\text{lev}(i) := \text{lev}(\mathbf{b}_i) = \text{lev}(\mathbf{b}_i^\delta) = \text{lev}(\phi_i)$  for all  $i = 1, \dots, N$ . Moreover, we define a (semi-)metric  $d(\cdot, \cdot)$  on  $\{1, \dots, N\}$  by

$$d(i, j) := \frac{\text{dist}(\text{mid}(\phi_i), \text{mid}(\phi_j))}{h_{\min\{\text{lev}(i), \text{lev}(j)\}}},$$

where  $\text{mid}(w)$  defines the barycenter of  $\text{supp}(w)$ .

Define the stiffness matrices  $\mathbf{S}(\omega) \in \mathbb{R}^{N \times N}$  associated with the bases  $\mathcal{B}(\omega)$  by

$$\mathbf{S}(\omega)_{ij} := \mathbf{a}_\omega \frac{\mathbf{b}_j(\omega)}{\|\mathbf{b}_j(\omega)\|_\omega}, \frac{\mathbf{b}_i}{\|\mathbf{b}_i(\omega)\|_\omega}.$$

The orthogonality of the bases  $\mathcal{B}$  motivates the approximation of the stiffness matrices by block-diagonal ones even after localization. Given  $1/C_1 > \delta > 0$ , define the block-diagonal stiffness matrices  $\mathbf{S}^\delta(\omega) \in \mathbb{R}^{N \times N}$  by

$$\mathbf{S}^\delta(\omega)_{ij} := \begin{cases} \mathbf{a}_\omega \left( \frac{\mathbf{b}_j^\delta(\omega)}{\|\mathbf{b}_j^\delta(\omega)\|_\omega}, \frac{\mathbf{b}_i^\delta(\omega)}{\|\mathbf{b}_i^\delta(\omega)\|_\omega} \right) & \text{for } \text{lev}(i) = \text{lev}(j), \\ 0 & \text{else.} \end{cases}$$

In the following, we use the spectral norm  $\|\cdot\|_2$ , i.e., the matrix norm induced by the Euclidean norm.

LEMMA 5. *There exists a constant  $C > 0$  that depends only on  $D$ , the shape regularity of  $\mathcal{T}_0$ , and the contrast  $\gamma_{\max}/\gamma_{\min}$  such that, for any  $\omega \in \Omega$  and for all  $\delta \leq 1/(2C_1)$ ,*

$$\|\mathbf{S}(\omega) - \mathbf{S}^\delta(\omega)\|_2 \leq C\delta.$$

Moreover, there exists a constant  $\zeta > 0$  which depends only on  $D$  such that

$$(3.1) \quad d(i, j) > \zeta(|\log(\delta)| + 1) \text{ or } \text{lev}(i) = \text{lev}(j) \implies \mathbf{S}_{ij}^\delta(\omega) = 0;$$

in particular, the number of nonzero entries  $\text{nnz}(\mathbf{S}^\delta(\omega)) \lesssim N(1 + |\log \delta|)^d$  is bounded uniformly in  $\omega$ .

*Proof.* The sparsity of the diagonal blocks follows from (2.17). For the proof of the error bound, define

$$\tilde{\mathbf{S}}^\delta(\omega)_{ij} := \begin{cases} \mathbf{a}_\omega \left( \frac{\mathbf{b}_j(\omega)}{\|\mathbf{b}_j(\omega)\|_\omega}, \frac{\mathbf{b}_i^\delta(\omega)}{\|\mathbf{b}_i^\delta(\omega)\|_\omega} \right) & \text{for } \text{lev}(i) = \text{lev}(j), \\ 0 & \text{else.} \end{cases}$$

Since  $|\mathbf{S}_{ij}(\omega) - \tilde{\mathbf{S}}_{ij}^\delta(\omega)| = 0$  whenever  $\text{lev}(i) = \text{lev}(j)$ , it suffices to bound the errors related to the diagonal blocks indexed by  $\ell = 1, 2, \dots, L$ . We have for any vectors  $x, y \in \mathbb{R}^{\#\mathcal{B}_\ell^\delta}$  that

$$\begin{aligned} & |x \cdot (\mathbf{S}_\ell(\omega) - \tilde{\mathbf{S}}_\ell^\delta(\omega))y| \\ &= \sum_{\text{lev}(i)=\ell} \sum_{\text{lev}(j)=\ell} x_i y_j \mathbf{a}_\omega \left( \frac{\mathbf{b}_i(\omega)}{\|\mathbf{b}_i(\omega)\|_\omega}, \frac{\mathbf{b}_j^\delta(\omega)}{\|\mathbf{b}_j^\delta(\omega)\|_\omega} \right) - \frac{\mathbf{b}_j(\omega)}{\|\mathbf{b}_j(\omega)\|_\omega} \lesssim \delta \|x\|_{\ell_2} \|y\|_{\ell_2} \end{aligned}$$

by (2.18). The same arguments show  $|x \cdot (\mathbf{S}_\ell^\delta(\omega) - \tilde{\mathbf{S}}_\ell^\delta(\omega))y| \lesssim \delta \|x\|_{\ell_2} \|y\|_{\ell_2}$ , and the triangle inequality readily proves the assertion.  $\square$

LEMMA 6. For any  $\omega \in \Omega$ , the normalized set  $\mathcal{B} = \mathcal{B}(\omega)$  or  $\mathcal{B} = \mathcal{B}^\delta(\omega)$  (with  $\delta \lesssim 1/L$  sufficiently small) is a Riesz basis in the sense that

$$(3.2) \quad C^{-1} \sum_{b \in \mathcal{B}} \alpha_b^2 \leq \sum_{b \in \mathcal{B}} \alpha_b \frac{b}{\|\mathbf{b}\|_\omega} \Big|_{H^1(D)}^2 \leq C \sum_{b \in \mathcal{B}} \alpha_b^2$$

holds with some constant  $C > 0$  which depends only on  $D$ , the shape regularity of  $\mathcal{T}_0$ , and the contrast  $\gamma_{\max}/\gamma_{\min}$ . This immediately implies that  $\mathbf{S}(\omega)$  and  $\mathbf{S}^\delta(\omega)$  are uniformly well conditioned.

*Proof.* Since  $\|\cdot\|_{H^1(D)}$  and  $\|\cdot\|_\omega$  are equivalent uniformly in  $\omega$  and the basis  $\mathcal{B}(\omega)$  is  $\mathbf{a}_\omega(\cdot, \cdot)$ -orthogonal across the levels, it suffices to consider one level  $k \in \{1, \dots, L\}$  in the case  $\mathcal{B} = \mathcal{B}(\omega)$ . The  $L^2(D)$ -orthogonality of the Haar basis and the construction of  $\mathcal{B}$  implies

$$\begin{aligned} \sum_{b \in \mathcal{B}_k} \alpha_b^2 &= \sum_{b \in \mathcal{B}_k} \alpha_b \frac{\phi_b}{\|\phi_b\|_{L^2(D)}} \Big|_{L^2(D)}^2 = \Pi_k \sum_{b \in \mathcal{B}_k} \alpha_b \frac{b}{\|\phi_b\|_{L^2(D)}} \Big|_{L^2(D)}^2 \\ &\leq (1 - \Pi_{k-1}) \sum_{b \in \mathcal{B}_k} \alpha_b \frac{b}{\|\phi_b\|_{L^2(D)}} \Big|_{L^2(D)}^2 \lesssim \sum_{b \in \mathcal{B}_k} \alpha_b \frac{b}{\|\mathbf{b}\|_\omega} \Big|_{H^1(D)}^2, \end{aligned}$$

where the last estimate follows from the Poincaré inequality and (2.9). For the proof of the converse direction, the construction of  $\mathcal{B}$  and boundedness of  $\mathcal{C}_k$  show

$$\begin{aligned} \sum_{b \in \mathcal{B}_k} \alpha_b \frac{b}{\|\mathbf{b}\|_\omega} \Big|_{H^1(D)}^2 &= (\text{id} - \mathcal{C}_k(\omega)) \tilde{\Pi}_k \sum_{b \in \mathcal{B}_k} \alpha_b \frac{b}{\|\mathbf{b}\|_\omega} \Big|_{H^1(D)}^2 \\ &\lesssim \tilde{\Pi}_k \sum_{b \in \mathcal{B}_k} \alpha_b \frac{\phi_b}{\|\phi_b\|_\omega} \Big|_{H^1(D)}^2 \lesssim \sum_{b \in \mathcal{B}_k} \alpha_b \frac{\phi_b}{\|\phi_b\|_{L^2(D)}} \Big|_{L^2(D)}^2 = \sum_{b \in \mathcal{B}_k} \alpha_b^2, \end{aligned}$$

where the second inequality follows from the inverse inequality (2.6) and (2.9).

The result for  $\mathcal{B} = \mathcal{B}^\delta(\omega)$  is slightly more involved as the  $\mathbf{a}_\omega(\cdot, \cdot)$ -orthogonal across the levels is lost. In a first step, Lemma 5 and the equivalence of  $\|\cdot\|_\omega$  and  $\|\cdot\|_{H^1(D)}$  imply

$$\sum_{b^\delta \in \mathcal{B}_\ell^\delta(\omega)} \alpha_{b^\delta} \frac{b^\delta}{\|\!|b^\delta\!\|_\omega} \Big|_{H^1(D)}^2 \simeq \mathcal{S}_\ell^\delta(\omega) \alpha \cdot \alpha \simeq \mathcal{S}_\ell(\omega) \alpha \cdot \alpha \pm C\delta \|\alpha\|_{\ell_2}^2 \simeq (1 \pm C\delta) \|\alpha\|_{\ell_2}^2$$

with the constant  $C$  from Lemma 5 and  $\delta \leq 1/(2C)$ . The second step concerns the quantification of nonorthogonality. The estimate (2.15) and the norm equivalence  $\|\cdot\|_\omega \simeq \|\cdot\|_{H^1(D)}$  imply

$$\|\!(\mathcal{C}_\ell^\delta - \mathcal{C}_\ell)v\!\|_\omega \lesssim \delta \|v\|_\omega \quad \text{for all } v \in H_0^1(D).$$

Consequently, we obtain as in the proof of Lemma 5, for some  $1 \leq k = \ell \leq L$ ,

$$\begin{aligned} & \mathbf{a}_\omega \left( \sum_{b^\delta \in \mathcal{B}_\ell^\delta(\omega)} \alpha_{b^\delta} \frac{b^\delta}{\|\!|b^\delta\!\|_\omega}, \sum_{b \in \mathcal{B}_k(\omega)} \beta_b \frac{b}{\|\!|b\!\|_\omega} \right) \\ &= \mathbf{a}_\omega \left( (\mathcal{C}_\ell^\delta - \mathcal{C}_\ell) \tilde{\Pi}_\ell \sum_{b_\phi^\delta \in \mathcal{B}_\ell^\delta(\omega)} \alpha_b \frac{\phi}{\|\!|b_\phi^\delta\!\|_\omega}, \sum_{b \in \mathcal{B}_k(\omega)} \beta_{b_\phi^\delta} \frac{b}{\|\!|b\!\|_\omega} \right) \\ &\lesssim \delta \|\alpha\|_{\ell_2} \|\beta\|_{\ell_2}, \end{aligned}$$

where we used the orthogonality across levels and the stability of  $\mathcal{B}$ . Symmetry of the argument concludes  $\mathbf{a}_\omega(\sum_{b^\delta \in \mathcal{B}_\ell^\delta(\omega)} \alpha_{b^\delta} \frac{b^\delta}{\|\!|b^\delta\!\|_\omega}, \sum_{b \in \mathcal{B}_k(\omega)} \beta_b \frac{b}{\|\!|b\!\|_\omega}) \lesssim \delta \|\alpha\|_{\ell_2} \|\beta\|_{\ell_2}$ , and we find

$$\begin{aligned} \sum_{b^\delta \in \mathcal{B}^\delta(\omega)} \alpha_{b^\delta} \frac{b^\delta}{\|\!|b^\delta\!\|_\omega} \Big|_{H^1(D)}^2 &\simeq \sum_{\ell=1}^L (1 \pm \delta) \|\alpha|_{\mathcal{B}_\ell^\delta(\omega)}\|_{\ell_2}^2 \pm \delta \sum_{\substack{i,j=1 \\ i=j}}^L \|\alpha|_{\mathcal{B}_i^\delta(\omega)}\|_{\ell_2} \|\alpha|_{\mathcal{B}_j^\delta(\omega)}\|_{\ell_2} \\ &\simeq (1 \pm C\delta L) \|\alpha\|_{\ell_2}^2 \end{aligned}$$

for sufficiently small  $\delta \leq 1/(2CL)$ . This concludes the proof. □

**4. Basis transformations.** This section analyzes the properties of a certain matrix representation of the  $L^2(D)$ -orthogonal projections  $\Pi_\ell: L^2(D) \rightarrow \text{span}(\bigcup_{j=0}^\ell \mathcal{H}_j)$  for  $\ell = 1, \dots, L$ . Given  $\omega \in \Omega$ , define the matrix  $\mathbf{T}(\omega) \in \mathbb{R}^{N \times N}$  by

$$\mathbf{T}_{ij}(\omega) := \frac{(\mathbf{b}_j(\omega), \phi_i)_{L^2(D)}}{\|\!|\mathbf{b}_j(\omega)\!\|_\omega \|\phi_i\|_{L^2(D)}^2}.$$

Given some  $v = \sum_{i=1}^N \alpha_i \frac{\mathbf{b}_i}{\|\!|\mathbf{b}_i\!\|_\omega}$  with  $\Pi_L v = \sum_{i=1}^N \beta_i \frac{\phi_i}{\|\phi_i\|_{L^2(D)}}$ , then, by definition,

$$\beta_i = \frac{(\Pi_L v, \phi_i)}{\|\phi_i\|_{L^2(D)}^2} = \sum_{j=1}^N \alpha_j \mathbf{T}_{ij} = (\mathbf{T}\alpha)_i,$$

i.e.,  $\beta = \mathbf{T}(\omega)\alpha$ . Given  $\delta > 0$ , a truncated approximation  $\mathbf{T}^\delta(\omega)$  of  $\mathbf{T}(\omega)$  is defined by

$$\mathbf{T}_{ij}^\delta(\omega) := \begin{cases} \frac{(\mathbf{b}_j^\delta(\omega), \phi_i)_{L^2(D)}}{\|\!|\mathbf{b}_j^\delta(\omega)\!\|_\omega \|\phi_i\|_{L^2(D)}^2} & \text{if } \text{lev}(j) \leq \text{lev}(i), \\ 0 & \text{else} \end{cases}$$

for any  $i, j \in \{1, \dots, N\}$ .  $\mathbf{T}^\delta(\omega)$  is a sparse lower block-triangular matrix, and the next lemma shows that the error of truncation is at most proportional to  $\delta$ . To explore the block structure of matrices, we shall introduce the following notation first. For any matrix  $K \in \mathbb{R}^{N \times N}$ , we define subblocks  $K_{(k,\ell)} \in \mathbb{R}^{\#\mathcal{H}_\ell \times \#\mathcal{H}_k}$  according to the level structure by

$$K_{(k,\ell)} := K|_{\{(i,j) : \text{lev}(i)=k, \text{lev}(j)=\ell\}}.$$

Thus, we may write

$$K = \begin{pmatrix} K_{(0,0)} & K_{(0,1)} & \cdots & K_{(0,L)} \\ K_{(1,0)} & K_{(1,1)} & \cdots & K_{(1,L)} \\ \vdots & \vdots & \ddots & \vdots \\ K_{(L,0)} & K_{(L,1)} & \cdots & K_{(L,L)} \end{pmatrix}.$$

LEMMA 7. For  $\delta > 0$  as in Lemma 6, there holds

$$\|\mathbf{T}(\omega) - \mathbf{T}^\delta(\omega)\|_2 \leq CL\delta,$$

and, for  $0 \leq \ell \leq k \leq L$ , there holds

$$(4.1) \quad \|\mathbf{T}^\delta(\omega)_{(k,\ell)}\|_2 \leq Ch_k.$$

Moreover,  $\mathbf{T}^\delta$  is lower block-triangular with sparse blocks; more precisely,

$$(\text{lev}(j) \geq \text{lev}(i) \text{ and } i = j) \text{ or } d(i, j) > \zeta(1 + |\log(\delta)|) \implies \mathbf{T}_{ij}^\delta = 0,$$

where  $\zeta > 0$  is the bandwidth from Lemma 5. The number of nonzero entries per block is bounded by  $\text{nnz}(\mathbf{T}^\delta(\omega)_{(k,\ell)}) \lesssim \#\mathcal{H}_k(1 + |\log \delta|)^d$ . The constant  $C > 0$  depends only on  $D$ , the shape regularity of  $\mathcal{T}_0$ , and the contrast  $\gamma_{\max}/\gamma_{\min}$ .

*Proof.* We see immediately that  $\mathbf{T}_{ij}(\omega) = 0$  for all  $\text{lev}(j) \geq \text{lev}(i)$  and  $i = j$  since

$$(\mathbf{b}_j(\omega), \phi_i)_{L^2(D)} = (\Pi_{\text{lev}(\phi_i)} \mathbf{b}_j(\omega), \phi_i)_{L^2(D)} = (\phi_j, \phi_i)_{L^2(D)} = 0.$$

Since  $\text{supp}(\mathbf{b}_i^\delta(\omega)) \cap \text{supp}(\phi_j) = \emptyset$  as soon as  $d(i, j) \gtrsim |\log(\delta)|$ , there is some  $\zeta > 0$  which depends only on  $D$  such that  $\mathbf{T}_{ij}^\delta(\omega) = 0$  for all  $d(i, j) > \zeta(1 + |\log(\delta)|)$ .

For any vectors  $x \in \mathbb{R}^{\#\mathcal{B}_k^\delta}$  and  $y \in \mathbb{R}^{\#\mathcal{B}_\ell^\delta}$ , we have

$$\begin{aligned} & x \cdot (\mathbf{T}_{(k,\ell)}(\omega) - \mathbf{T}_{(k,\ell)}^\delta(\omega))y \\ &= \sum_{\text{lev}(i)=k} \sum_{\text{lev}(j)=\ell} x_i y_j \frac{\phi_i}{\|\phi_i\|_{L^2(D)}^2}, \frac{\mathbf{b}_j^\delta(\omega)}{\|\mathbf{b}_j^\delta(\omega)\|_\omega} - \frac{\mathbf{b}_j(\omega)}{\|\mathbf{b}_j(\omega)\|_\omega} \Big|_{L^2(D)} \lesssim \delta \|x\|_{\ell_2} \|y\|_{\ell_2} \end{aligned}$$

by the Friedrichs inequality and (2.19). This implies  $\|\mathbf{T}(\omega)_{(\ell,k)} - \mathbf{T}^\delta(\omega)_{(\ell,k)}\|_2 \lesssim \delta$ . Summing up over the levels proves  $\|\mathbf{T}(\omega) - \mathbf{T}^\delta(\omega)\|_2 \lesssim L\delta$ .

To see (4.1), note that  $w := \sum_{\phi_i \in \mathcal{H}_k} \alpha_i \phi_i$  and  $b := \sum_{\mathbf{b}_j(\omega) \in \mathcal{B}_\ell} \beta_j \mathbf{b}_j^\delta(\omega)$  satisfy

$$\begin{aligned} \alpha^T \mathbf{T}^\delta(\omega) \beta &= (w, b)_{L^2(D)} = ((1 - \Pi_k)w, b)_{L^2(D)} = (w, (1 - \Pi_k)b)_{L^2(D)} \\ &\lesssim h_k \|w\|_{L^2(D)} \|b\|_\omega \lesssim h_k \|\alpha\|_{\ell_2} \|\beta\|_{\ell_2} \end{aligned}$$

by Lemma 6. This concludes the proof. □

**5. Inverse stiffness matrices and averaging.** This section proves that the inverse of the stiffness matrix  $\mathbf{S}^\delta(\omega)$  (w.r.t. the coefficient adapted bases  $\mathbf{B}^\delta(\omega)$ ) defined in the previous section can be efficiently approximated by a sparse matrix. One possibility to compute an approximate inverse of the matrix  $\mathbf{S}^\delta(\omega)$  is to apply the conjugate gradient (CG) method to the matrix with unit vectors  $e_i \in \mathbb{R}^N$  as right-hand sides. The sparsity pattern from Lemma 5 shows that one matrix-vector product with  $\mathbf{S}^\delta e_i$  increases the number of nonzero entries to  $\#\{1 \leq j \leq N : d(i, j) \lesssim 1 + |\log(\delta)|\}$ . Thus, after  $k \in \mathbb{N}$  iterations of the CG method, the resulting vector has about  $\#\{1 \leq j \leq N : d(i, j) \lesssim k(1 + |\log(\delta)|)\}$  nonzero entries. Since the condition number  $\kappa(\mathbf{S}^\delta)$  is uniformly bounded due to Lemma 6, the number of iterations grows only logarithmically in the desired accuracy  $\delta$ . Thus, the cost of  $k \simeq 1 + |\log(\delta)|$  iterations of the CG method to reach the accuracy can be bounded roughly by  $(1 + |\log(\delta)|)^2$ .

LEMMA 8. For  $\delta > 0$  as in Lemma 6, there exists a matrix  $\mathbf{R}^\delta(\omega)$  such that  $\|\mathbf{S}(\omega)^{-1} - \mathbf{R}^\delta(\omega)\|_2 \leq \delta$ . Moreover,  $\mathbf{R}^\delta(\omega)$  satisfies

$$(5.1) \quad d(i, j) > C_{\text{inv}}\zeta(|\log(\delta)|^2 + 1) \text{ or } \text{lev}(i) = \text{lev}(j) \implies \mathbf{R}_{ij}^\delta(\omega) = 0$$

for  $\zeta$  from Lemma 5 and  $C_{\text{inv}} > 0$  depending only on  $D$ , the shape regularity of  $\mathcal{T}_0$ , and the contrast  $\gamma_{\text{max}}/\gamma_{\text{min}}$ . The number of nonzero entries is bounded by  $\text{nnz}(\mathbf{R}^\delta) \lesssim N(1 + |\log(\delta)|)^d$ , and the cost to compute  $\mathbf{R}^\delta$  is bounded by the same number.

*Proof.* Due to Lemma 5 and the fact that  $\mathbf{B}(\omega)$  is a Riesz basis (Lemma 6), we observe that all eigenvalues of  $\mathbf{S}^\delta(\omega)$  are of order  $\mathcal{O}(1)$  as long as  $\delta \lesssim 1$ . Therefore, we can obtain  $\mathbf{R}^\delta(\omega)$  by application of CG steps to  $\mathbf{S}^{\tilde{\delta}}(\omega)$  (we chose  $\tilde{\delta} > 0$  later; see, e.g., [40, Chapter 6]). The convergence properties of CG show

$$\|\mathbf{S}^{\tilde{\delta}}(\omega)^{-1} - \mathbf{R}^\delta(\omega)\|_2 \leq \delta$$

if we perform  $k = \mathcal{O}(|\log(\delta)| + 1)$  CG steps. This follows since

$$\|\text{res}_k\|_{\ell_2} \simeq \sqrt{\mathbf{S}^{\tilde{\delta}}(\omega)\text{res}_k \cdot \text{res}_k}$$

for the residual  $\text{res}_k$  of the CG method. From Lemma 5, we see that  $\mathbf{R}^\delta(\omega)$  satisfies

$$d(i, j) > \zeta(|\log(\delta)| + 1)^2 \text{ or } \text{lev}(i) = \text{lev}(j) \implies \mathbf{R}_{ij}^\delta(\omega) = 0$$

since each CG step increases the bandwidth by the original bandwidth. With Lemma 5, we conclude the proof by choosing  $k \simeq 1 + |\log(\delta)|$  and  $\tilde{\delta} \simeq \delta$ . The constructive nature of this proof immediately reveals the cost estimate.  $\square$

LEMMA 9. We define a discrete approximation to  $\mathcal{A}^{-1}$  by

$$R := \mathbb{E}[(\mathbf{T}^{-T}(\omega)\mathbf{S}(\omega)\mathbf{T}^{-1}(\omega))^{-1}] = \mathbb{E}[\mathbf{T}(\omega)\mathbf{S}(\omega)^{-1}\mathbf{T}(\omega)^T].$$

For  $\delta > 0$  as in Lemma 6, we define a perturbed and truncated version of  $R$  by  $R^\delta \in \mathbb{R}^{N \times N}$

$$(5.2) \quad (R^\delta)_{(\ell,k)} := \begin{cases} \mathbb{E}[\mathbf{T}^\delta(\omega)\mathbf{R}^\delta(\omega)\mathbf{T}^\delta(\omega)^T]_{(\ell,k)} & \ell + k \leq |\log(\delta)|, \\ 0 & \text{else,} \end{cases}$$

which satisfies  $\|R - R^\delta\|_2 \leq CL^2\delta$ . The number of nonzero entries in  $R^\delta$  is bounded by

$\text{nnz}(R^\delta) \lesssim L/\delta^d$ , and up to the computation of the expectation, the cost to produce  $R^\delta$  is bounded by  $\mathcal{O}(L^d/\delta^d)$ . The constant  $C > 0$  depends only on  $D$ , the shape regularity of  $\mathcal{T}_0$ , and the contrast  $\gamma_{\max}/\gamma_{\min}$ .

*Proof.* We define the auxiliary operator

$$\tilde{R}^\delta := \mathbb{E}[\mathbf{T}^\delta(\omega)\mathbf{R}^\delta(\omega)\mathbf{T}^\delta(\omega)^T].$$

Analogously to matrix subblocks, we may partition vectors  $x \in \mathbb{R}^N$  by

$$x = (x_{(1)}, \dots, x_{(L)})$$

with  $x_{(\ell)} \in \mathbb{R}^{\#\mathcal{H}_\ell}$ . Using this notation and following the proofs of Lemmas 5, 7, and 8, we show

$$\|(R - \tilde{R}^\delta)x\|_{\ell_2}^2 \lesssim \sum_{\ell=1}^L \sum_{k=1}^L (R_{(\ell,k)} - \tilde{R}_{(\ell,k)}^\delta)x_{(k)} \Big|_{\ell_2}^2 \lesssim \delta^2 L^4 \sum_{\ell=1}^L \sum_{k=1}^L \|x_{(k)}\|_{\ell_2}^2 = \delta^2 L^4 \|x\|_{\ell_2}^2$$

and, hence,  $\|R - \tilde{R}^\delta\|_2 \lesssim \delta L^2$ . The estimate (4.1) implies, for  $\ell + k > |\log(\delta)|$ ,

$$\begin{aligned} \|(\tilde{R}^\delta - R^\delta)_{(\ell,k)}\|_2 &\leq \sum_{j=0}^L \|\mathbf{T}^\delta(\omega)_{(\ell,j)}\|_2 \|\mathbf{R}^\delta(\omega)_{(j,j)}\|_2 \|(\mathbf{T}^\delta(\omega)_{(k,j)})\|_2 \\ &\lesssim \sum_{j=0}^L h_\ell(1 + \delta)h_k \\ &\lesssim L2^{-\ell-k}. \end{aligned}$$

This implies, for  $x \in \mathbb{R}^N$ ,

$$\begin{aligned} \|(\tilde{R}^\delta - R^\delta)x\|_{\ell_2}^2 &\leq \sum_{i,j=0}^L \|(\tilde{R}^\delta - R^\delta)_{(i,j)}x|_{(j)}\|_{\ell_2}^2 \lesssim \sum_{j=0}^L \|x|_{(j)}\|_{\ell_2}^2 \sum_{i=|\log(\delta)|-j}^L L2^{-i-j} \\ &\lesssim L\delta \|x\|_{\ell_2}^2. \end{aligned}$$

The number of nonzero entries in  $R^\delta$  can be bounded sufficiently by ignoring the sparsity within the blocks and just summing up the entries

$$\sum_{0 \leq i+j \leq |\log(\delta)|} \#(R^\delta)_{(i,j)} \lesssim \sum_{0 \leq i+j \leq |\log(\delta)|} 2^{d(i+j)} \lesssim L\delta^{-d},$$

where we used that  $(R^\delta)_{i,j} \in \mathbb{R}^{\#\mathcal{H}_i \times \#\mathcal{H}_j}$  and  $\#\mathcal{H}_i \simeq 2^{di}$ . The cost to compute  $R^\delta$  (up to the computation of the expectation) is bounded by the linear cost in the number of nonzeros to set up  $\mathbf{T}^\delta$  (see Lemma 7) as well as the cost to set up  $\mathbf{R}^\delta$ . The latter is bounded by  $\mathcal{O}(N(1 + |\log(\delta)|^d))$  according to Lemma 8, where  $N = \#\mathcal{H} \simeq \delta^{-d}$ . This concludes the proof.  $\square$

To formulate the following main theorem, we identify the matrix  $R^\delta$  with an operator  $\mathcal{R}^\delta: L^2(D) \rightarrow L^2(D)$  via the natural embedding  $\iota: \mathbb{R}^N \rightarrow \text{span}(\mathcal{H})$ ,  $\iota(\alpha) = \sum_{i=1}^N \alpha_i \phi_i \in L^2(D)$ . There holds  $\mathcal{R}^\delta := \iota R^\delta \iota^*$ .

**THEOREM 10.** *For a given accuracy  $\delta > 0$  with  $\delta \lesssim 1/L$  sufficiently small, there exists a finite-dimensional operator  $\mathcal{R}^\delta: L^2(D) \rightarrow L^2(D)$  which depends only on  $\delta$  such that*

$$\|\mathcal{A}^{-1} - \mathcal{R}^\delta\|_{\mathcal{L}(L^2(D), L^2(D))} \leq \delta.$$

*The corresponding operator matrix  $R^\delta$  from Lemma 9 has at most  $\mathcal{O}(|\log(\delta)|^{2d+1}\delta^{-d})$  nonzero entries, and up to the computation of the expectation, the cost of producing  $R^\delta$  is bounded by  $\mathcal{O}(|\log(\delta)|^{3d}\delta^{-d})$ . The hidden constant depends only on  $D$ , the shape regularity of  $\mathcal{T}_0$ , and the contrast  $\gamma_{\max}/\gamma_{\min}$ .*

*Constructive proof.* We use the operator matrix  $R^\delta \in \mathbb{R}^{N \times N}$  from Lemma 9. Given  $f \in L^2(D)$ , define  $\mathbf{F}(\omega) \in \mathbb{R}^N$  by  $\mathbf{F}_i(\omega) := (f, \mathbf{b}_i / \|\mathbf{b}_i\|_\omega)$ . By definition, there holds  $\mathbf{S}(\omega)\boldsymbol{\alpha}(\omega) = \mathbf{F}(\omega)$  with  $\mathbf{u}_L(\omega) := \sum_{i=1}^N \boldsymbol{\alpha}_i(\omega)\mathbf{b}_i / \|\mathbf{b}_i\|_\omega \in \text{span}(\mathcal{B}(\omega))$  being the Galerkin approximation to  $\mathbf{u}(\omega) \in H_0^1(\Omega)$ . Galerkin orthogonality

$$\mathbf{a}_\omega(\mathbf{u}(\omega) - \mathbf{u}_L(\omega), \text{span}(\mathcal{B}(\omega))) = 0$$

implies  $\mathbf{u}(\omega) - \mathbf{u}_L(\omega) \in W_L$  and, hence,  $\Pi_L(\mathbf{u}(\omega) - \mathbf{u}_L(\omega)) = 0$ . Thus,

$$\|\mathbf{u}(\omega) - \Pi_L \mathbf{u}_L(\omega)\|_{L^2(D)} \leq \|(1 - \Pi_L)\mathbf{u}(\omega)\|_{L^2(D)} \lesssim h_L \|f\|_{L^2(D)}$$

using standard approximation properties of piecewise constants (Poincaré inequality) and a standard energy bound. With the transfer matrices  $\mathbf{T}(\omega)$  from Lemma 7, we obtain

$$\tilde{F} := \iota^* f = \mathbf{T}^{-T}(\omega)\mathbf{F}(\omega),$$

and, hence,  $\boldsymbol{\beta} \in \mathbb{R}^N$  with  $\mathbf{T}^{-T}(\omega)\mathbf{S}(\omega)\mathbf{T}^{-1}(\omega)\boldsymbol{\beta}(\omega) = \tilde{F}$  satisfies  $\mathbf{T}(\omega)\boldsymbol{\alpha}(\omega) = \boldsymbol{\beta}(\omega)$ . Together with Lemma 7, this shows that  $\Pi_L \mathbf{u}_L(\omega) = \sum_{i=1}^N \boldsymbol{\beta}_i(\omega)\phi_i / \|\phi_i\|_{L^2(D)}$ . The approximate solution  $\mathcal{R}^\delta f = \sum_{i=1}^N \gamma_i \phi_i / \|\phi_i\|_{L^2(D)}$  with  $\gamma := R^\delta \tilde{F}$  satisfies

$$|\gamma - \mathbb{E}[\boldsymbol{\beta}]| \lesssim L^2 \delta \|f\|_{L^2(D)}$$

by use of Lemma 9 and since  $\mathbb{E}_M[\boldsymbol{\beta}] = R\tilde{F}$ . Since  $\mathcal{H}$  is an orthogonal basis, we obtain immediately  $\|\mathcal{R}^\delta f - \mathcal{R}f\|_{L^2(D)} \lesssim L^2 \delta \|f\|_{L^2(D)}$ , where

$$\mathcal{R}^\delta f = \mathbb{E}[\mathbf{u}_L].$$

Combining the above error bounds, we conclude

$$\|\mathbb{E}[\mathbf{u}] - \mathcal{R}^\delta f\|_{L^2(D)} \lesssim (L^2 \delta + h_L) \|f\|_{L^2(D)}.$$

With  $L \simeq |\log \delta|$  and  $h_L \simeq \delta$ , there holds  $\|\mathbb{E}[\mathbf{u}] - \mathcal{R}^\delta f\|_{L^2(D)} \lesssim (1 + |\log(\delta)|^2)\delta \|f\|_{L^2(D)}$ , and Lemma 9 shows  $\text{nnz}(R^\delta) \simeq L/\delta^d \simeq |\log \delta|/\delta^d$ . To clean up the result, we replace  $\delta$  with  $\tilde{\delta} := (1 + |\log(\delta)|^2)\delta$  and obtain  $\text{nnz}(R^\delta) \lesssim (1 + |\log(\tilde{\delta})|^{2d+1})/\tilde{\delta}^d$  (note that we used  $\delta \gtrsim \tilde{\delta}/(1 + |\log \tilde{\delta}|^2)$  and  $\delta \leq \tilde{\delta}$  for sufficiently small  $\delta > 0$ ). The same calculation with  $L^d/\delta^d$  instead of  $L/\delta^d$  (see Lemma 9) proves the cost estimate for producing  $R^\delta$  and, hence, concludes the proof.  $\square$

**6. Sparse operator compression.** Theorem 10 shows that the expected operator can indeed be compressed to a sparse matrix. The constructive proof motivates a compression algorithm by simply replacing the expectation by a suitable sample mean. For this purpose, let  $\Omega_M \subset \Omega$  be a finite set of sampling points with  $|\Omega_M| = M \in \mathbb{N}$ , and define the sample mean  $\mathbb{E}_M[\mathbf{X}] := M^{-1} \sum_{\omega \in \Omega_M} \mathbf{X}(\omega)$  for a random field  $\mathbf{X}$ . It is readily seen that Lemma 9 remains valid when  $\mathbb{E}$  is replaced by  $\mathbb{E}_M$ . More precisely, define

$$R_M := \mathbb{E}_M[(\mathbf{T}^{-T}(\omega)\mathbf{S}(\omega)\mathbf{T}^{-1}(\omega))^{-1}] = \mathbb{E}_M[\mathbf{T}(\omega)\mathbf{S}(\omega)^{-1}\mathbf{T}(\omega)^T]$$

and a perturbed and truncated version of  $R_M$  by  $R_M^\delta \in \mathbb{R}^{N \times N}$ ,

$$(6.1) \quad (R_M^\delta)_{(\ell,k)} := \begin{cases} \mathbb{E}_M[\mathbf{T}^\delta(\omega)\mathbf{R}^\delta(\omega)\mathbf{T}^{\delta T}(\omega)]_{(\ell,k)} & \ell + k \leq |\log(\delta)|, \\ 0 & \text{else.} \end{cases}$$

Then

$$(6.2) \quad \|R_M - R_M^\delta\|_2 \leq CL^2\delta,$$

and the number of nonzero entries in  $R_M^\delta$  is bounded by  $\mathcal{O}(L/\delta^d)$ .

*Remark 11.* The truncation condition  $\ell + k \leq |\log(\delta)|$  in (6.1) can be relaxed to  $\ell + k \leq C|\log(\delta)|$  for some  $C \simeq 1$  without any harm. In practice, when  $L \simeq |\log \delta|$  is chosen, a natural choice would be  $\ell + k \leq L$ . In the numerical experiment of section 7, we will see that sometimes it can be advantageous to include a few more blocks of the lower right part of the matrix (see (7.1)) to recover gradient information.

The analogue of Theorem 10 in this discrete stochastic setting then as follows.

**COROLLARY 12.** *For a given accuracy  $\delta > 0$  as in Theorem 10 and a set of  $M$  samples  $\Omega_M \subset \Omega$ ,  $M \in \mathbb{N}$ , there exists a finite-dimensional operator  $\mathcal{R}_M^\delta: L^2(D) \rightarrow L^2(D)$  which depends only on the sample coefficients  $\mathbf{A}(\omega)$ ,  $\omega \in \Omega_M$ ,  $\delta$ , and  $D$  such that*

$$\|\mathcal{A}^{-1} - \mathcal{R}_M^\delta\|_{\mathcal{L}(L^2(D),L^2(D))} \leq \delta + \|(\mathbb{E} - \mathbb{E}_M)[\mathcal{A}^{-1}]\|_{\mathcal{L}(L^2(D),L^2(D))}.$$

*The corresponding operator matrix  $R_M^\delta$  has  $\mathcal{O}(|\log(\delta)|^{2d+1}\delta^{-d})$  nonzero entries and can be computed with cost bounded by  $\mathcal{O}(|\log(\delta)|^{3d}\delta^{-d})$ . The hidden constant depends only on  $D$ , the shape regularity of  $\mathcal{T}_0$ , and the contrast  $\gamma_{\max}/\gamma_{\min}$ .*

When using a plain Monte Carlo sampling, the mean squared sampling error scales like  $M^{-1}$ , meaning that  $M \simeq \delta^{-2}$  samples suffice to ensure that the sampling error is not dominating the error bound. This is optimal in the present setting with no assumptions on the distribution of the random diffusion coefficient. More advanced sampling techniques such as quasi-Monte Carlo methods are certainly possible under additional assumptions, such as a rapid decay of eigenvalues of a given Karhunen–Loève expansion of the random parameter (see [10] for a discussion in terms of PDEs with random parameters). Even more promising is the possible intertwining of the hierarchical decomposition and the sampling procedure in the spirit of multilevel/multi-index Monte Carlo (see, e.g., [20, 23] for the seminal works as well as [8]). At least in the regime where stochastic homogenization applies, the computation of basis functions is likely to be essentially independent of the parameter  $\omega$  for levels that are much coarser than the characteristic length scale of random oscillation

(or correlation) [19]. This has been made rigorous in a two-level setting in [16]. For the increasing variance for the levels approaching the scale of correlation, stationarity could be exploited to improve the overall complexity.

Another interesting case is the use of lognormal coefficients  $\mathbf{A}(\omega) = \exp(\mathbf{Z}(\omega))$  for a normal random field  $\mathbf{Z}$ . As shown in [15], such random fields can be efficiently generated for general covariance functions and nonuniform grids. The present analysis, however, breaks down since the assumption of bounded contrast in (1) is violated. The authors are confident, however, that the arguments can be modified in the sense that the extreme contrast samples will only appear with very low probability (the tails of the Gaussian density). Thus, a polynomial dependence on the contrast (as is observed for the present construction) will not perturb the final result.

We shall finally mention that so far the construction relies on the exact solution of the (infinite-dimensional) corrector problems (2.8) and their preconditioned variant, respectively. The elegant way to transfer all results to a fully discrete setting is to consider a space-discrete problem from the very beginning. It is readily seen that all constructions and results remain valid if we replace the space  $V = H_0^1(D)$  by a suitable finite-dimensional subspace  $V_h \subset V$  throughout the paper. We have in mind some standard  $V$ -conforming finite element space  $V_h$  that is based on some regular mesh of width  $h$  which turns the preconditioned corrector problems into finite element problems on the mesh  $h$  restricted to local subdomains of diameter  $h_\ell |\log \delta|$ . The only restriction that comes with this discretization step is that the mesh size  $h$  limits the number of possible levels  $L$  in the hierarchical decomposition and, hence, the possible accuracy  $\delta \lesssim h$  when the sparse approximation is compared with the reference solution  $\mathbb{E}[\mathbf{u}_h]$ , where  $\mathbf{u}_h$  solves (1.2) with  $V$  replaced with  $V_h$ . Clearly, the overall accuracy of the fully discrete method depends on the error  $\|\mathbb{E}[\mathbf{u} - \mathbf{u}_h]\|_{L^2(D)}$ , which is a standard finite element error that depends on the spatial regularity of  $\mathbf{A}$  and also its possible frequencies of oscillations. All this is well understood and implies the usual conditions on the smallness of  $h$  so that  $\mathbf{A}$  is properly resolved (see, e.g., [39]).

**7. Numerical experiment.** This section presents some simple numerical experiments to illustrate the performance of the method. We consider the domain  $D = [0, 1]^d$  for  $d = 1, 2$ , and the coefficient  $\mathbf{A}$  is scalar i.i.d., and, on each cell of the uniform Cartesian mesh  $\mathcal{T}_\varepsilon$ , it is uniformly distributed in the interval  $[\gamma_{\min}, \gamma_{\max}] = [0.5, 10]$ . The mesh width (scale of oscillation/correlation length) is  $\varepsilon = 2^{-8}$  ( $d = 1$ ) and  $\varepsilon = 2^{-5}$  ( $d = 2$ ).

The approximations of the solution operator are based on sequences of uniform Cartesian meshes  $\mathcal{T}_\ell$  ( $\ell = 0, 1, 2, \dots, L$ ) of mesh width  $h_\ell = 2^{-\ell}$  that do not necessarily resolve  $\varepsilon$ . We compute approximations  $\mathcal{R}^L = \mathcal{R}_{M_L}^\delta$  of the expected solution operator depending on the maximal level  $L$ , which means that we expect  $L^2(D)$  errors of order  $\delta \approx 2^{-L}$ . The truncation of blocks is performed based on the criterion  $k + \ell \leq L$  as indicated in Remark 11. For the solution of the corrector problems and the reference solution  $\mathbf{u}_h$ , we use  $d$ -linear finite elements on the mesh  $\mathcal{T}_h$ , where  $h = 2^{-14}$  ( $d = 1$ ) and  $h = 2^{-9}$  ( $d = 2$ ). To achieve accuracy of order  $\delta$  (w.r.t. to the reference solution), we perform  $\lceil L/2 \rceil$  CG iterations for both computing the correctors  $\mathcal{C}^\delta(\omega)$  and inverting the block-diagonal stiffness matrices  $\mathcal{S}^\delta(\omega)$ . For the approximation of the expected values, we use a quasi-Monte Carlo method (particularly a Sobol sequence) with appropriate numbers of sampling points  $M_h := h^{-1}$  for the reference solutions and  $M_L := 2^L$  for the approximations. While we did not show that the problem is smooth enough to justify the use of quasi-Monte Carlo sampling, we still

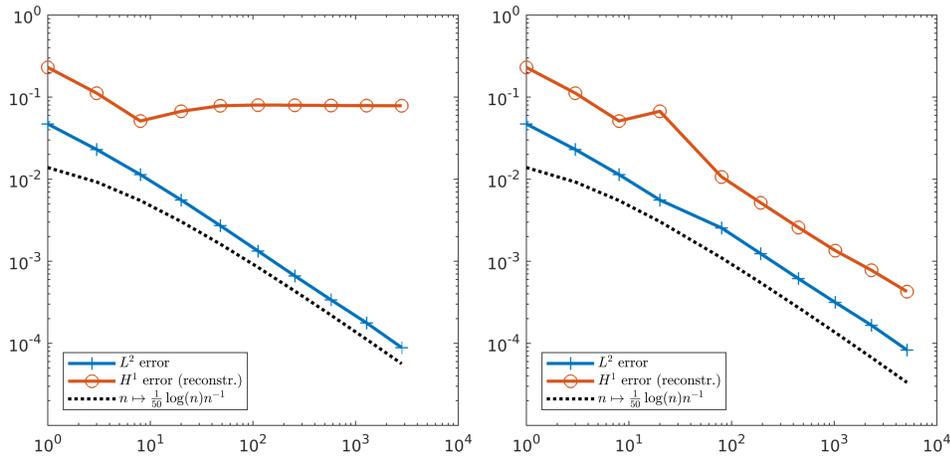


FIG. 2. Numerical results in one dimension:  $L^2(D)$ -errors  $\|\mathbb{E}_{M_h}[\mathbf{u}_h] - \mathcal{R}^L f\|_{L^2(D)}$  and  $H^1(D)$ -errors  $\|\nabla(\mathbb{E}_{M_h}[\mathbf{u}_h] - u_L^1)\|_{L^2(D)}$  of postprocessed approximation for  $L = 1, 2, \dots, 10$ . Left: Errors versus  $\text{nnz}(R^L)$  using original approach (6.1). Right: Errors versus  $\text{nnz}(\tilde{R}^L)$  using modified approach (7.1).

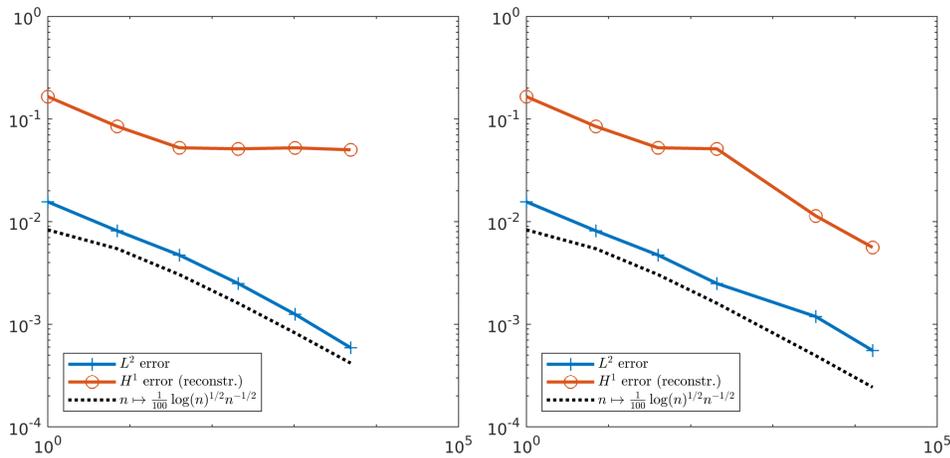


FIG. 3. Numerical results in two dimensions:  $L^2(D)$ -errors  $\|\mathbb{E}_{M_h}[\mathbf{u}_h] - \mathcal{R}^L f\|_{L^2(D)}$  and  $H^1(D)$ -errors  $\|\nabla(\mathbb{E}_{M_h}[\mathbf{u}_h] - u_L^1)\|_{L^2(D)}$  of postprocessed approximation for  $L = 1, 2, \dots, 6$ . Left: Errors versus  $\text{nnz}(R^L)$  using original approach (6.1). Right: Errors versus  $\text{nnz}(\tilde{R}^L)$  using modified approach (7.1).

observe the expected higher convergence rate compared to plain Monte Carlo sampling and thus save significant computing time.

Since the computation of a reference expected operator is hardly feasible, we only compute the error for one nonsmooth deterministic right-hand side  $f = \chi_{[.5,1] \times [0,1]^{d-1}} \in L^2(D) \setminus H^1(D)$ . Figures 2 and 3 (left plots) depict the errors  $\|\mathbb{E}_{M_h}[\mathbf{u}_h] - \mathcal{R}^L f\|_{L^2(D)}$  versus the number of nonzero entries of  $\text{nnz}(R^L)$  for  $L = 1, 2, \dots$ . The results are very well in agreement (up to a possibly pessimistic logarithmic factor) with the prediction that

$$\|\mathbb{E}_{M_h}[\mathbf{u}_h] - \mathcal{R}^L f\|_{L^2(D)} \lesssim M_L^{-1} + \frac{|\log(\text{nnz}(R^L))|^{2+1/d}}{\text{nnz}(R^L)^{1/d}}$$

for  $d = 1, 2$ . This is the optimal rate of convergence (up to a logarithmic factor) given a piecewise constant approximation.

In this setting, where the expected solution  $\mathbb{E}[\mathbf{u}]$  is even  $H^2(D)$  regular, it would be desirable to recover gradient information from the piecewise constant approximation by suitable postprocessing, e.g., in the hierarchical basis associated with a constant coefficient. Figures 2 and 3 (left plots) indicate that this is not automatically achieved for nonsmooth right-hand sides with the present choice of parameters. However, when the truncation in (6.1) is slightly relaxed in the form

$$(7.1) \quad (\tilde{R}^L)_{(\ell,k)} := \begin{cases} \mathbb{E}_M[\mathbf{T}^\delta(\omega)\mathbf{R}^\delta(\omega)\mathbf{T}^\delta(\omega)^T]_{(\ell,k)} & \ell + k \leq L + \max(1, \lceil \log_2 L \rceil), \\ 0 & \text{else,} \end{cases}$$

accurate reconstruction of gradients seems possible. From this slightly more accurate but slightly more dense approximation  $\tilde{R}^L$ , we can reconstruct the coefficients of a smooth approximation  $u_L^1 \in \text{span } \mathcal{B}(\omega_\Delta)$  (with  $\omega_\Delta \in \Omega$  such that  $\mathbf{A}(\omega_\Delta) = 1$ ) in the hierarchical basis that corresponds to the Laplacian by simply applying  $T^\delta(\omega_\Delta)^{-1}$  to  $\tilde{R}^L f$ . The errors of this smooth postprocessing  $\|\nabla(\mathbb{E}_{M_h}[\mathbf{u}_h] - u_L^1)\|_{L^2(D)}$  are plotted in Figures 2 and 3 (right plots) against the number of nonzero entries  $\text{nnz}(\tilde{R}^L)$ . The observed rate of convergence for the  $H^1$ -error is  $\text{nnz}(\tilde{R}^L)^{-1/d}$  (up to a logarithmic factor) which is nearly optimal. See also the plots on the left of Figures 2 and 3, which indicate that the step from (6.1) to (7.1) is essential for meaningful gradient reconstruction.

These first numerical results support the theoretical findings and indicate the potential of the approach. Since the techniques that were used in the construction of the method and its analysis, in particular the LOD, generalize in a straightforward way to other classes of operators, such as linear elasticity [24] or Helmholtz problems [38, 17, 6], we believe that the sparse compression algorithm for the approximation of expected solution operators is applicable beyond the prototypical model problem of this paper.

#### REFERENCES

- [1] S. ARMSTRONG, T. KUUSI, AND J.-C. MOURRAT, *The additive structure of elliptic homogenization*, *Invent. Math.*, 208 (2017), pp. 999–1154, <https://doi.org/10.1007/s00222-016-0702-4>.
- [2] M. BACHMAYR, A. COHEN, R. DEVORE, AND G. MIGLIORATI, *Sparse polynomial approximation of parametric elliptic PDEs. Part II: Lognormal coefficients*, *ESAIM Math. Model. Numer. Anal.*, 51 (2017), pp. 341–363, <https://doi.org/10.1051/m2an/2016051>.
- [3] M. BACHMAYR, A. COHEN, AND G. MIGLIORATI, *Sparse polynomial approximation of parametric elliptic PDEs. Part I: Affine coefficients*, *ESAIM Math. Model. Numer. Anal.*, 51 (2017), pp. 321–339, <https://doi.org/10.1051/m2an/2016045>.
- [4] J. BOURGAIN, *On a homogenization problem*, *J. Stat. Phys.*, (2018), <https://doi.org/10.1007/s10955-018-1981-5>.
- [5] A. BOURGEAT AND A. PIATNITSKI, *Approximations of effective coefficients in stochastic homogenization*, *Ann. Inst. H. Poincaré Probab. Statist.*, 40 (2004), pp. 153–165, [https://doi.org/10.1016/S0246-0203\(03\)00065-7](https://doi.org/10.1016/S0246-0203(03)00065-7).
- [6] D. L. BROWN, D. GALLISTL, AND D. PETERSEIM, *Multiscale Petrov-Galerkin method for high-frequency heterogeneous Helmholtz equations*, in *Meshfree Methods for Partial Differential equations VIII*, *Lect. Notes Comput. Sci. Eng.* 115, Springer, Cham, Switzerland, 2017, pp. 85–115.
- [7] J. DICK, *Walsh spaces containing smooth functions and quasi-Monte Carlo rules of arbitrary high order*, *SIAM J. Numer. Anal.*, 46 (2008), pp. 1519–1553, <https://doi.org/10.1137/060666639>.

- [8] J. DICK, M. FEISCHL, AND C. SCHWAB, *Improved Efficiency of a Multi-Index FEM*, preprint, <https://arxiv.org/abs/1806.04159>, 2018.
- [9] J. DICK, R. N. GANTNER, Q. T. LE GIA, AND C. SCHWAB, *Multilevel higher-order quasi-Monte Carlo Bayesian estimation*, *Math. Models Methods Appl. Sci.*, 27 (2017), pp. 953–995, <https://doi.org/10.1142/S021820251750021X>.
- [10] J. DICK, F. Y. KUO, Q. T. LE GIA, D. NUYENS, AND C. SCHWAB, *Higher order QMC Petrov-Galerkin discretization for affine parametric operator equations with random field inputs*, *SIAM J. Numer. Anal.*, 52 (2014), pp. 2676–2702, <https://doi.org/10.1137/130943984>.
- [11] M. DUERINCKX, A. GLORIA, AND M. LEMM, *A remark on a surprising result by Bourgain in homogenization*, *Comm. Partial Differential Equations*, 44 (2019), pp. 1345–1357, <https://doi.org/10.1080/03605302.2019.1638934>.
- [12] M. DUERINCKX, A. GLORIA, AND F. OTTO, *The Structure of Fluctuations in Stochastic Homogenization*, preprint, arXiv 1602.01717 [math.AP] (2016).
- [13] D. ELFVERSON, E. GEORGIOULIS, AND A. MÅLQVIST, *An adaptive discontinuous Galerkin multiscale method for elliptic problems*, *Multiscale Model. Simul.*, 11 (2013), pp. 747–765, <https://doi.org/10.1137/120863162>.
- [14] D. ELFVERSON, E. H. GEORGIOULIS, A. MÅLQVIST, AND D. PETERSEIM, *Convergence of a discontinuous Galerkin multiscale method*, *SIAM J. Numer. Anal.*, 51 (2013), pp. 3351–3372, <https://doi.org/10.1137/120900113>.
- [15] M. FEISCHL, F. Y. KUO, AND I. H. SLOAN, *Fast random field generation with  $H$ -matrices*, *Numer. Math.*, 140 (2018), pp. 639–676, <https://doi.org/10.1007/s00211-018-0974-2>.
- [16] J. FISCHER, D. GALLISTL, AND D. PETERSEIM, *A Priori Error Analysis of a Numerical Stochastic Homogenization Method*, preprint, <https://arxiv.org/abs/1912.11646>, (2019).
- [17] D. GALLISTL AND D. PETERSEIM, *Stable multiscale Petrov-Galerkin finite element method for high frequency acoustic scattering*, *Comput. Methods Appl. Mech. Engrg.*, 295 (2015), pp. 1–17.
- [18] D. GALLISTL AND D. PETERSEIM, *Computation of quasi-local effective diffusion tensors and connections to the mathematical theory of homogenization*, *Multiscale Model. Simul.*, 15 (2017), pp. 1530–1552.
- [19] D. GALLISTL AND D. PETERSEIM, *Numerical stochastic homogenization by quasi-local effective diffusion tensors*, *Commun. Math. Sci.*, 17 (2019), pp. 637–651, <https://doi.org/10.4310/CMS.2019.v17.n3.a3>.
- [20] M. B. GILES, *Multilevel Monte Carlo path simulation*, *Oper. Res.*, 56 (2008), pp. 607–617, <https://doi.org/10.1287/opre.1070.0496>.
- [21] A. GLORIA AND F. OTTO, *The Corrector in Stochastic Homogenization: Optimal Rates, Stochastic Integrability, and Fluctuations*, preprint, arXiv, 1510.08290 [math.AP] (2015).
- [22] A. GLORIA AND F. OTTO, *Quantitative results on the corrector equation in stochastic homogenization*, *J. Eur. Math. Soc. (JEMS)*, 19 (2017), pp. 3489–3548.
- [23] A.-L. HAJI-ALI, F. NOBILE, AND R. TEMPONE, *Multi-index Monte Carlo: When sparsity meets sampling*, *Numer. Math.*, 132 (2016), pp. 767–806, <https://doi.org/10.1007/s00211-015-0734-5>.
- [24] P. HENNING AND A. PERSSON, *A multiscale method for linear elasticity reducing Poisson locking*, *Comput. Methods Appl. Mech. Engrg.*, 310 (2016), pp. 156–171, <https://doi.org/10.1016/j.cma.2016.06.034>.
- [25] P. HENNING AND D. PETERSEIM, *Oversampling for the multiscale finite element method*, *Multiscale Model. Simul.*, 11 (2013), pp. 1149–1175, <https://doi.org/10.1137/120900332>.
- [26] T. Y. HOU, D. HUANG, K. C. LAM, AND P. ZHANG, *An adaptive fast solver for a general class of positive definite matrices via energy decomposition*, *Multiscale Model. Simul.*, 16 (2018), pp. 615–678, <https://doi.org/10.1137/17M1140686>.
- [27] T. Y. HOU AND P. ZHANG, *Sparse operator compression of higher-order elliptic operators with rough coefficients*, *Res. Math. Sci.*, 4 (2017), Paper No. 24, 49, <https://doi.org/10.1186/s40687-017-0113-1>.
- [28] J. KIM AND M. LEMM, *On the averaged Green’s function of an elliptic equation with random coefficients*, *Arch. Ration. Mech. Anal.*, 234 (2019), pp. 1121–1166, <https://doi.org/10.1007/s00205-019-01409-1>.
- [29] R. KORNUBER, D. PETERSEIM, AND H. YSERENTANT, *An analysis of a class of variational multiscale methods based on subspace decomposition*, *Math. Comp.*, 87 (2018), pp. 2765–2774, <https://doi.org/10.1090/mcom/3302>.
- [30] R. KORNUBER AND H. YSERENTANT, *Numerical homogenization of elliptic multiscale problems by subspace decomposition*, *Multiscale Model. Simul.*, 14 (2016), pp. 1017–1036, <https://doi.org/10.1137/15M1028510>.

- [31] S. M. KOZLOV, *The averaging of random operators*, Mat. Sb. (N.S.), 109 (1979), pp. 188–202, 327.
- [32] A. MÅLQVIST AND D. PETERSEIM, *Localization of elliptic multiscale problems*, Math. Comp., 83 (2014), pp. 2583–2603.
- [33] H. OWHADI, *Multigrid with rough coefficients and multiresolution operator decomposition from hierarchical information games*, SIAM Rev., 59 (2017), pp. 99–149, <https://doi.org/10.1137/15M1013894>.
- [34] H. OWHADI AND C. SCOVEL, *Operator-Adapted Wavelets, Fast Solvers, and Numerical Homogenization*, Cambridge Monographs on Applied and Computational Mathematics 35, Cambridge University Press, Cambridge, 2019.
- [35] H. OWHADI AND L. ZHANG, *Gamblets for opening the complexity-bottleneck of implicit schemes for hyperbolic and parabolic ODEs/PDEs with rough coefficients*, J. Comput. Phys., 347 (2017), pp. 99–128, <https://doi.org/10.1016/j.jcp.2017.06.037>.
- [36] G. C. PAPANICOLAOU AND S. R. S. VARADHAN, *Boundary Value Problems with Rapidly Oscillating Random Coefficients*, in Random Fields, Vols. I, II (Esztergom, 1979), Colloq. Math. Soc. János Bolyai 27, North-Holland, Amsterdam, 1981, pp. 835–873.
- [37] D. PETERSEIM, *Variational multiscale stabilization and the exponential decay of fine-scale correctors*, in Building Bridges: Connections and Challenges in Modern Approaches to Numerical Partial Differential Equations, G. R. Barrenea, F. Brezzi, A. Cangiani, and E. H. Georgoulis, eds., Springer International Publishing, Cham, Switzerland, 2016, pp. 343–369, [https://doi.org/10.1007/978-3-319-41640-3\\_11](https://doi.org/10.1007/978-3-319-41640-3_11).
- [38] D. PETERSEIM, *Eliminating the pollution effect in Helmholtz problems by local subscale correction*, Math. Comp., 86 (2017), pp. 1005–1036, <https://doi.org/10.1090/mcom/3156>.
- [39] D. PETERSEIM AND S. SAUTER, *Finite elements for elliptic problems with highly varying, nonperiodic diffusion matrix*, Multiscale Model. Simul., 10 (2012), pp. 665–695, <https://doi.org/10.1137/10081839X>.
- [40] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, 2nd ed., SIAM, Philadelphia, 2003, <https://doi.org/10.1137/1.9780898718003>.
- [41] F. SCHÄFER, T. J. SULLIVAN, AND H. OWHADI, *Compression, Inversion, and Approximate PCA of Dense Kernel Matrices at Near-Linear Computational Complexity*, preprint, <https://arxiv.org/abs/1706.02205>, (2017).
- [42] P. WOJTASZCZYK, *A Mathematical Introduction to Wavelets*, London Mathematical Society Student Texts 37, Cambridge University Press, Cambridge, 1997, <https://doi.org/10.1017/CBO9780511623790>.
- [43] J. XU, *Iterative methods by space decomposition and subspace correction*, SIAM Rev., 34 (1992), pp. 581–613, <https://doi.org/10.1137/1034116>.
- [44] H. YSERENTANT, *Old and new convergence proofs for multigrid methods*, Acta Numer., 2 (1993), pp. 285–326, <https://doi.org/10.1017/S0962492900002385>.
- [45] V. V. YURINSKIĬ, *Averaging of symmetric diffusion in a random medium*, Sib. Mat. Zh., 27 (1986), pp. 167–180, 215.