# Variationally correct operator learning: Reduced basis neural operator with a posteriori error estimation

Yuan Qiu[a], Wolfgang Dahmen[b], Peng Chen[a,*]

[a]*School of Computational Science and Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA*
[b]*Department of Mathematics, University of South Carolina, Columbia, SC 29208, USA*

**Abstract**

Minimizing PDE-residual losses is a common strategy to promote physical consistency in neural operators. However, standard formulations often lack *variational correctness*, meaning that small residuals do not guarantee small solution errors due to the use of non-compliant norms or ad hoc penalty terms for boundary conditions. This work develops a variationally correct operator learning framework by constructing first-order system least-squares (FOSLS) objectives whose values are provably equivalent to the solution error in PDE-induced norms. We demonstrate this framework on stationary diffusion and linear elasticity, incorporating mixed Dirichlet-Neumann boundary conditions via *variational lifts* to preserve norm equivalence without inconsistent penalties. To ensure the function space conformity required by the FOSLS loss, we propose a Reduced Basis Neural Operator (RBNO). The RBNO predicts coefficients for a pre-computed, conforming reduced basis, thereby ensuring variational stability by design while enabling efficient training. We provide a rigorous convergence analysis that bounds the total error by the sum of finite element discretization bias, reduced basis truncation error, neural network approximation error, and statistical estimation errors arising from finite sampling and optimization. Numerical benchmarks validate these theoretical bounds and demonstrate that the proposed approach achieves superior accuracy in PDE-compliant norms compared to standard baselines, while the residual loss serves as a reliable, computable *a posteriori* error estimator.

*Keywords:* parametric partial differential equations, neural operator, reduced basis methods, first-order system least squares, variational correctness, a posteriori error estimate

## 1   Introduction

Determining physical states of interest solely from observational data is generally intractable without invoking governing physical laws. These laws are typically modeled by systems of partial differential equations (PDEs) that depend on unspecified problem data, such as boundary conditions, source terms, or coefficient fields. The collection of solutions generated by varying these inputs forms the *solution manifold*. In many-query tasks such as uncertainty quantification (UQ) [1, 2], Bayesian inversion [3, 4] and optimal experimental design [5, 6], exploring this manifold requires evaluating the *data-to-solution map*—the solution operator—thousands or millions of times.

The computational burden of traditional high-fidelity discretizations (e.g., finite elements) has driven a surge of interest in *operator learning* for efficient surrogate modeling. Architectures such as Deep Operator Network (DeepONet) [7] and Fourier Neural Operators (FNO) [8] learn mappings between infinite-dimensional function spaces, typically via *regression* on pre-computed high-fidelity snapshots. However, this supervised paradigm faces two critical limitations: the prohibitive cost of generating large training datasets

---

[*]Corresponding author.

*Email addresses:* `yuan.qiu@gatech.edu` (Yuan Qiu), `dahmen@math.sc.edu` (Wolfgang Dahmen), `pchen402@gatech.edu` (Peng Chen)

and the reliance on standard $L^2$-type losses. Crucially, these regression norms are often not *PDE-compliant*, meaning a small training error does not necessarily guarantee accuracy in the physically relevant norms.

To reduce reliance on high-fidelity data, *residual minimization* (as used in Physics-Informed Neural Networks or PINNs [9, 10]) offers an alternative by embedding the PDE directly into the loss function. However, standard residual losses often lack *variational correctness* [11]: they fail to provide uniform upper and lower bounds on the true error in a proper PDE norm. This typically occurs because bulk residuals are measured in norms that are too strong (e.g., $L^2$ for second-order operators), while boundary conditions are enforced via penalty terms in norms that are too weak or inconsistent [12]. Consequently, minimizing such losses does not rigorously control the solution error. While recent research has proposed *a posteriori* error estimation strategies—including learning error certificates [13], developing computable bounds [14], and utilizing functional-type norms [15]—achieving variational correctness *by design* requires a stable variational formulation. Specifically, one must identify a pair of *PDE-compliant norms* for the trial and test spaces such that the residual in the *dual test norm* is uniformly proportional to the error in the trial norm. This condition holds precisely when the variational formulation is stable in the sense of the Babuška-Nečas Theorem [11].

A practical obstruction, well-known in adaptive finite elements, is that evaluating the dual norm is generally difficult; we refer to [11] for various approaches to overcome this. In this work, we focus on a strategy that seeks stable variational formulations where the test space is *self-dual* (i.e., an $L^2$-space), allowing the residual to be computed simply as an $L^2$-norm. While standard higher-order PDE formulations typically do not admit such stability, reformulating them as a *first-order system* often does. This approach, known as First-Order System Least-Squares (FOSLS) [16, 17], has been successfully applied to space-time formulations of parabolic and wave equations [18, 19], as well as stationary diffusion and linear elasticity [20]. Deep learning methods utilizing FOSLS have already demonstrated significant promise for solving PDEs [21, 22, 23, 24, 20]. However, these studies have primarily focused on solving individual PDE instances, rather than addressing the more challenging task of learning the solution operator over a distribution of parameters.

In this work, we bridge this gap by developing a variationally correct operator learning framework for parametric PDEs. To demonstrate the methodology, we restrict our attention here to linear parametric PDEs—specifically, stationary diffusion and linear elasticity. Our approach diverges from standard methods in two key ways. First, we handle mixed Dirichlet-Neumann boundary conditions through *variational lifts* rather than penalty terms, preserving rigorous norm equivalence. Second, to inherently satisfy the *conformity requirements* of our loss function, we propose a *Reduced Basis Neural Operator* (RBNO). While standard neural operators effectively learn mappings between generic function spaces, enforcing specific regularity constraints (such as flux continuity) remains non-trivial. The RBNO bypasses this difficulty by predicting coefficients for a basis that is conforming by construction, thereby ensuring variational stability by design while mitigating the complexity of high-dimensional outputs.

To practically construct this basis, we exploit the rapid decay of Kolmogorov $n$-widths for elliptic problems [25, 26]. We project the solution manifold onto a low-dimensional linear space generated by Proper Orthogonal Decomposition (POD) of high-fidelity snapshots. Since these snapshots are computed using a conforming discretization, their linear combinations naturally preserve the required function space regularity. The neural operator then learns the mapping from parameters to the *coefficients* of this reduced basis. This hybrid architecture functions analogously to an encoder-decoder with a fixed linear decoder [27, 28], yet distinctively, it is trained by minimizing the rigorous FOSLS residual rather than a regression loss. We note that while data-driven reduced basis neural networks have been explored elsewhere [29, 30, 31, 32, 33, 34, 35, 36], our approach uniquely leverages this structure to enable variationally correct training.

**Contributions.** The specific contributions of this paper are as follows:

1. *Rigorous construction of variationally correct FOSLS objectives.* We derive $L^2$-based FOSLS residual losses for stationary diffusion and linear elasticity. Crucially, we show how mixed Dirichlet/Neumann boundary data can be incorporated via *variational lifts*, provably preserving the equivalence between the loss and the solution error in PDE-induced norms without relying on ad hoc boundary penalties.

2. *Finite element realization and discrete error control.* We formulate conforming finite element discretizations that make the quadratic residual loss efficiently computable. We establish a theoretical link between this discrete loss and the discretization errors, proving convergence rates that depend on

the mesh size, polynomial order, and PDE regularity. This ensures the training objective remains a reliable proxy for the true error even after discretization.

3. *Reduced Basis Neural Operator (RBNO) with theoretical guarantees.* We propose the RBNO architecture, which predicts coefficients for a pre-computed, conforming POD basis. By construction, this ensures the network's output satisfies the function space conformity required by the FOSLS loss. We provide a convergence analysis that bounds the total error by the sum of finite element discretization bias, reduced basis projection error, neural network approximation error, statistical estimation errors arising from finite sampling, and optimization error, rigorously justifying the hybrid approach.

4. *Numerical validation.* We demonstrate the method on diffusion benchmarks and a linear elasticity problem. We numerically validate the convergence analysis with respect to finite element discretization, reduced basis projection, statistical estimation and optimization errors in RBNO training. The results confirm that the proposed residual loss serves as a tight, computable *a posteriori* error estimator and that the RBNO achieves superior accuracy in PDE-compliant norms compared to two baselines.

The remainder of the paper is organized as follows. Section 2 reviews the abstract framework of operator learning and variational correctness. Section 3 derives the variationally correct FOSLS formulations for diffusion and elasticity, including the variational lifting strategy for boundary conditions. Section 4 details the finite element realization of the loss and the associated discrete error control. Section 5 introduces the RBNO architecture and establishes the theoretical convergence analysis. Section 6 presents numerical validation of the error estimates and performance comparisons with standard neural operators. Finally, Section 7 discusses limitations and future directions.

## 2 Conceptual background and orientation

We briefly outline the abstract problem formulation, upon which subsequent developments will be based. We consider a *parametric family* of linear partial differential equations (PDEs) in residual form

$$\mathcal{R}(u; \mathfrak{p}) = \mathcal{B}_{\mathfrak{p}} u - f = 0 \tag{1}$$

where $\mathfrak{p}$ stands for a scalar/vector-valued parameter—or more generally for a parameter field—that may range over a given parameter range or space $\mathfrak{P}$ with measure $\mu$. Given $\mathfrak{p} \in \mathfrak{P}$, we are interested in assessing the solution $u = u(\mathfrak{p})$ of (1) for the parameter-instance $\mathfrak{p}$. In what follows, we assume that for each $\mathfrak{p} \in \mathfrak{P}$, (1) is well-posed (in a sense detailed below). Hence, the *solution operator* or *parameter-to-solution map* $\mathfrak{p} \mapsto u(\mathfrak{p})$ is well-defined and its range $\mathcal{M} = \mathcal{M}(\mathfrak{P})$—the set of all solution-states, obtained by traversing $\mathfrak{P}$—is often referred to as *solution manifold*. It will be important to specify for each $\mathfrak{p} \in \mathfrak{P}$ a suitable (infinite dimensional) trial space $\mathbb{U}$ that is to accommodate $\mathcal{M}$ and could, in principle, depend on $\mathfrak{p}$.

We postpone a discussion on this issue for the moment and recall that a common approach to learning $\mathfrak{p} \mapsto u(\mathfrak{p})$ is completely data-driven and employs regression in $\mathbb{U}$ (or for convenience in $L_2$). This requires computing first a large number of high-fidelity solutions, e.g., for randomly chosen parameter samples. A large number of additional test samples is needed to assess the generalization error, as the inherent uncertainty in optimization success renders a priori expressivity results insufficient. In summary, given the at best achievable Monte Carlo rates, depending on the accuracy requirements, the associated computational cost may be prohibitive.

Therefore, we focus in this paper on contriving *residual-type* loss functions. In this setting, the unknown solution is encoded by known problem data (such as source terms), allowing one to directly optimize over the trainable weights that describe the hypothesis class. The number of training samples scales with the number of residual evaluations, not with the number of high-fidelity solution snapshots. The price is that we can no longer assess directly the (generalization) error of the optimization outcome in the chosen model-compliant norm, *unless* we can show that the size of the loss is *uniformly proportional* to that error. We refer to such a loss as *variationally correct*, see [11].

As shown in [11], variational correctness is intimately related to a *stable variational formulation* for each *fiber problem* (1). To be specific, suppose one has found for each $\mathfrak{p} \in \mathfrak{P}$ a pair of Hilbert spaces $\mathbb{U}_\mathfrak{p}, \mathbb{V}_\mathfrak{p}$ such that the family of bilinear forms $b_\mathfrak{p}(\cdot, \cdot)$, $\mathfrak{p} \in \mathfrak{P}$, defined by $b_\mathfrak{p}(w, v) = (\mathcal{B}_\mathfrak{p} w)(v)$, $w \in \mathbb{U}_\mathfrak{p}, v \in \mathbb{V}_\mathfrak{p}$. Then well-posedness of (1) is equivalent to saying that $b_\mathfrak{p}(\cdot, \cdot)$ satisfies for each $\mathfrak{p}$ a continuity condition, an inf-sup condition and a surjectivity condition according to the Babuška-Nečas-Theorem, see e.g. [37, Chapter 1] or [38, 39]. This in turn means that $\mathcal{B}_\mathfrak{p}$, defined weakly as above, is an *isomorphism* from $\mathbb{U}_\mathfrak{p}$ onto $\mathbb{V}'_\mathfrak{p}$. Noting that for any $w \in \mathbb{U}_\mathfrak{p}$, $\mathcal{B}_\mathfrak{p}(u(\mathfrak{p}) - w) = f - \mathcal{B}_\mathfrak{p} w$, this is equivalent to saying that there exist constants $0 < c_\mathfrak{p} \leq C_\mathfrak{p} < \infty$ such that

$$c_\mathfrak{p} \|u(\mathfrak{p}) - w\|_{\mathbb{U}_\mathfrak{p}} \leq \|f - \mathcal{B}_\mathfrak{p} w\|_{\mathbb{V}'_\mathfrak{p}} = \|\mathcal{R}(w; \mathfrak{p})\|_{\mathbb{V}'} \leq C_\mathfrak{p} \|u(\mathfrak{p}) - w\|_{\mathbb{U}_\mathfrak{p}}, \quad \forall w \in \mathbb{U}_\mathfrak{p}. \tag{2}$$

Thus, the residual, measured in the dual test-norm, is a lower and upper bound for the error in the model-compliant norm $\|\cdot\|_{\mathbb{U}_\mathfrak{p}}$, albeit with proportionality constants that may depend on $\mathfrak{p}$, see e.g. [40]. For elliptic (coercive) problems, the choice $\mathbb{V}_\mathfrak{p} = \mathbb{U}_\mathfrak{p}$ is natural, but for other problems, finding a pair $\mathbb{U}_\mathfrak{p}, \mathbb{V}_\mathfrak{p}$ that gives rise to a stable variational formulation (and hence to (29)) is part of the homework.

First, the above notation indicates that the dependence of a stable pair of trial and test space may be *essential*, meaning when $\mathfrak{p}$ varies, then $\mathbb{U}_\mathfrak{p}$ or $\mathbb{V}_\mathfrak{p}$ or both may vary even as sets, see [11] for related results. For the problems studied here, this is not the case, i.e., the spaces agree as sets, and the norms are equivalent, so that we henceforth write $\mathbb{U}_\mathfrak{p} = \mathbb{U}$ and $\mathbb{V}_\mathfrak{p} = \mathbb{V}$. In particular, it makes sense to view the solution manifold as a subset of $\mathbb{U}$. We refer to [11] for scenarios where $\mathbb{U}_\mathfrak{p}$ and/or $\mathbb{V}_\mathfrak{p}$ depend on $\mathfrak{p}$ in an essential way, i.e., even as sets.

Note that learning the map $\mathfrak{p} \mapsto u(\mathfrak{p})$ (respecting the correct metrics) is tantamount to approximating the function $u = u(x, \mathfrak{p})$, as a function of spatial and parametric variables in a model-compliant norm. It is then natural to interpret this function as the solution of a *single* (lifted) variational problem, whose residual enters the training loss. Following [11, 40], consider the Bochner spaces

$$\mathbb{X} := L_2(\mathfrak{P}; \mathbb{U}), \quad \mathbb{Y} := L_2(\mathfrak{P}; \mathbb{V})$$

and the bilinear form

$$b(w, v) := \int_\mathfrak{P} b_\mathfrak{p}(w(\mathfrak{p}), v(\mathfrak{p})) d\mu(\mathfrak{p}), \quad \mathbb{X} \times \mathbb{Y} \to \mathbb{R},$$

where $\mu$ is a (fixed) probability measure on $\mathfrak{P}$, reflecting the probabilistic nature of the parameter dependence. It has been shown in [11] that the "lifted" problem: given $f \in \mathbb{Y}' = L_2(\mathfrak{P}; \mathbb{V}')$, find $u \in \mathbb{X}$ such that

$$b(u, v) = f(v), \quad v \in \mathbb{Y}, \tag{3}$$

is well-posed if and only if the family of fiber problems (1) is *uniformly* well-posed over $\mathfrak{P}$, which then means that for $\mathcal{B} : \mathbb{X} \to \mathbb{Y}'$, defined by $(\mathcal{B}w)(v) := b(w, v)$, $w \in \mathbb{X}, v \in \mathbb{Y}$,

$$c_0 \|u - w\|_\mathbb{X} \leq \|f - \mathcal{B}w\|_{\mathbb{Y}'} =: \|\mathcal{R}(w)\|_{\mathbb{Y}'} \leq C_0 \|u - w\|_\mathbb{X}, \quad w \in \mathbb{X}, \tag{4}$$

where $c_0 = \inf_{\mathfrak{p} \in \mathfrak{P}} c_\mathfrak{p}$, $C_0 = \sup_{\mathfrak{p} \in \mathfrak{P}} C_\mathfrak{p}$, for $c_\mathfrak{p}, C_\mathfrak{p}$ from (29).

Therefore, $\|\mathcal{R}(w)\|_{\mathbb{Y}'}$ can be viewed as an *ideal* residual loss. However, first, integration over $\mathfrak{P}$ is not practically feasible, and second, in view of the definition $\|\mathcal{R}(w)\|_{\mathbb{Y}'} := \sup_{v \in \mathbb{Y}} \frac{b(w,v) - f(v)}{\|v\|_\mathbb{Y}}$, the (exact) evaluation of a non-trivial ($\mathbb{Y} \neq \mathbb{Y}'$) dual norm is not practical.

Regarding the first issue, note that (4) describes proportionalities between *expectations* because, by definition

$$\|u - w\|_\mathbb{X}^2 = \mathbb{E}_{\mathfrak{p} \sim \mu}\big[\|w(\mathfrak{p}) - u(\mathfrak{p})\|_\mathbb{U}^2\big], \quad \|\mathcal{R}(w)\|_{\mathbb{Y}'}^2 = \mathbb{E}_{\mathfrak{p} \sim \mu}\big[\|\mathcal{R}(w; \mathfrak{p})\|_{\mathbb{V}'}^2\big]. \tag{5}$$

A common response to the first issue is to approximate the expectation by its empirical counterpart

$$\mathcal{L}(w; \widehat{\mathfrak{P}}) := \frac{1}{\#\widehat{\mathfrak{P}}} \sum_{\mathfrak{p} \in \widehat{\mathfrak{P}}} \mathcal{L}(w; \mathfrak{p}) =: \|\mathcal{R}(w)\|_{\ell_2(\widehat{\mathfrak{P}}; \mathbb{V}')}^2, \quad \text{where } \mathcal{L}(w; \mathfrak{p}) := \|\mathcal{R}(w; \mathfrak{p})\|_{\mathbb{V}'}^2. \tag{6}$$

4

Here $\widehat{\mathfrak{P}}$ denotes a collection of finite random samples of $\mathfrak{p}$ sampled from $\mu$ and $\#\widehat{\mathfrak{P}}$ denotes the sample size. Of course, (4) implies its discrete counterpart

$$c_0\|u - w\|_{\ell_2(\widehat{\mathfrak{P}};\mathbb{U})} \leq \mathcal{L}(w;\widehat{\mathfrak{P}}) \leq C_0\|u - w\|_{\ell_2(\widehat{\mathfrak{P}};\mathbb{U})}, \quad w \in \mathbb{X}. \tag{7}$$

Specifically, when working with a hypothesis class $\mathcal{H} = \mathcal{H}(\Theta) = \{w(\cdot, \mathfrak{p}; \theta) : \theta \in \Theta\}$ determined by a given budget of trainable weights $\theta \in \Theta$, the regression ansatz

$$\theta^* \in \operatorname*{argmin}_{\theta \in \Theta} \frac{1}{\#\widehat{\mathfrak{P}}} \sum_{\mathfrak{p} \in \widehat{\mathfrak{P}}} \|u(\mathfrak{p}) - w(\mathfrak{p}; \theta)\|_{\mathbb{U}}^2 \tag{8}$$

is replaced by the "equivalent" task of finding

$$\theta^* \in \operatorname*{argmin}_{\theta \in \Theta} \mathcal{L}(w(\theta); \widehat{\mathfrak{P}}), \tag{9}$$

that spares us from computing the snapshots $u(\mathfrak{p})$, which could be done approximately by minimizing the corresponding *empirical risk* with finite training samples. The deviation between $\mathcal{L}(w; \widehat{\mathfrak{P}})$ and $\mathbb{E}_{\mathfrak{p}\sim\mu}\big[\|\mathcal{R}(w; \mathfrak{p})\|_{\mathbb{V}'}^2\big]$, when $\#\widehat{\mathfrak{P}} \to \infty$, the generalization error, so to speak, is a matter of statistical estimation.

In summary, $\|f - \mathcal{B}w\|_{\mathbb{V}'} = \|\mathcal{R}(w)\|_{\mathbb{V}'}$, or its discrete counterpart (6), is "ideal" residual loss; however, the second issue remains how to evaluate (7) for training purposes. A central message from [11], which we reinforce here, is that one should reformulate (1) (if it is not already in this form) as a *first-order system* of PDEs, since this often provides greater flexibility in constructing stable variational formulations. In particular, it is often possible to find a formulation in which the test space $\mathbb{V}$ is an $L_2$-type space (a direct product of $L_2$-spaces) and hence self-dual $\mathbb{V}' = \mathbb{V}$. This is, for instance, the case for space-time variational formulations of parabolic problems or the wave equation, [18, 19, 41], as well as for parametric diffusion problems and static elasticity models, which are our focus here; see [16, 17, 11, 40, 20].

When $\mathbb{V} = \mathbb{V}'$ is an $L_2$-space, the loss (6) is a least squares $L_2$-residual and can therefore (up to quadrature errors) be evaluated. Nevertheless, the ensuing optimization task (9), obtained when $w(\cdot, \cdot, \theta)$ belongs to a hypothesis class comprising deep neural networks with input variables $(x, \mathfrak{p}) \in \Omega \times \mathfrak{P}$, is, in general, hard to handle. This optimization task can be substantially alleviated when using a *hybrid* format of low-rank type, where we separate dependence on spatial and parametric variables. This has been done, for instance, in [11, 40], where approximation systems are employed that contain elements, which are finite elements as functions of spatial variables with $\mathfrak{p}$-dependent expansion coefficients that can be represented, for instance, by DNNs. In the present paper, we also employ a hybrid format that differs in terms of the spatially dependent modes.

The underlying rationale is that for elliptic models and their close relatives, the map $\mathfrak{p} \to u(\mathfrak{p})$ is (under mild assumptions on the parametric coefficient fields) even holomorphic (see [25, 20]) and the so-called *Kolmogorov $n$-widths* decay robustly in the parametric dimension. Recall that for a compact subset $\mathcal{K}$ in a Banach space $\mathbb{X}$, the Kolmogorov $n$-widths (as a measure of thickness of $\mathcal{K}$ in $\mathbb{X}$) are given by

$$d_n(\mathcal{K})_{\mathbb{X}} := \inf_{\dim \mathbb{V}_n = n} \sup_{v \in \mathcal{K}} \inf_{v_n \in \mathbb{V}_n} \|v - v_n\|_{\mathbb{X}}, \quad n \in \mathbb{N}. \tag{10}$$

For $\mathcal{K} = \mathcal{M}$ and $\mathbb{X} = L_\mu^2(\mathfrak{P}; \mathbb{U})$ the $d_n(\mathcal{K})_{\mathbb{X}}$ measure how well the solution manifold can be approximated from a single linear space in a worst-case sense, i.e., in $L_\infty(\mathfrak{P}; \mathbb{U})$. Thinking of $\mathfrak{p}$ as a random variable distributed according to some measure $\mu$ on $\mathfrak{P}$, the Bochner space $L_\mu^2(\mathfrak{P}; \mathbb{U})$ is more akin to a machine learning approach, and (5) suggests measuring accuracy in a mean-squared sense. The corresponding optimality benchmark then reads

$$\delta_n(\mathcal{M}, \mu)_{\mathbb{U}}^2 := \inf_{\dim \mathbb{U}_n = n} \int_{\mathfrak{P}} \min_{w \in \mathbb{U}_n} \|u(\mathfrak{p}) - w\|_{\mathbb{U}}^2 d\mu(\mathfrak{p}). \tag{11}$$

It is well-known that for this metric the *best linear approximation spaces* result from the Hilbert-Schmidt decomposition of the operator $M_u : L_\mu^2(\mathfrak{P}) \to \mathbb{U}$ defined by

$$M_u w := \int_{\mathfrak{P}} u(\mathfrak{p}) w(\mathfrak{p}) d\mu(\mathfrak{p}),$$

has a Hilbert-Schmidt decomposition

$$M_u := \sum_{k=1}^{\infty} s_k \langle \cdot, \phi_k \rangle_{L_\mu^2(\mathfrak{P})} u_k,$$

where $\mathbf{s} = (s_k)_{k\in\mathbb{N}} \in \ell_2(\mathbb{N})$ is non-negative and non-increasing while $(u_k)_{k\in\mathbb{N}}$ and $(\phi_k)_{k\in\mathbb{N}}$ are orthonormal systems in $\mathbb{U}$ and $L_\mu^2(\mathfrak{P})$, respectively. In particular, the spaces $\mathbb{U}_n := \mathrm{span}\,\{u_k : k \leq n\}$ realize the benchmark (11) with

$$\delta_n(\mathcal{M}, \mu)_{\mathbb{U}}^2 = \sum_{k>n} s_k^2, \quad n \in \mathbb{N}. \tag{12}$$

While it is impossible to determine $\mathbb{U}_n$ exactly, one can approximate the $u_k$ with the aid of the well-established strategy of *Proper Orthogonal Decomposition* (POD). Roughly, this is done by choosing a sufficiently large "truth space" $\mathbb{U}_h \subset \mathbb{U}$, typically a finite element space with mesh size $h$, compute for a sufficiently large number of random samples $\mathfrak{p}^i$, $i = 1, \ldots, N_{\mathrm{POD}}$, corresponding (approximate) solutions $u^i = \tilde{u}(\mathfrak{p}^i)$ of $\mathcal{R}(u; \mathfrak{p}^i)$ in $\mathbb{U}_h$, and compute a Singular Value Decomposition of $u^i$, $i = 1, \ldots, N_{\mathrm{POD}}$, in $\mathbb{U}$. When $\mu$ is a probability measure it readily follows that $\delta_n(\mathcal{M}, \mu)_{\mathbb{U}} \leq d_n(\mathcal{M})_{\mathbb{U}}$ so that a rapid decay of the Kolmogorov $n$-widths implies a rapid decay of the $\delta_n(\mathcal{M}, \mu)_{\mathbb{U}}$. In such a situation, we expect that for a sufficiently large training set $\mathfrak{P}_{N_{\mathrm{POD}}} = \{\mathfrak{p}^i : i = 1, \ldots, N_{\mathrm{POD}}\}$, the decay of the approximate singular values $\tilde{s}_i$ closely reflects the decay of the exact singular values $s_i$ and that the tail $\sum_{k=n+1}^{N_{\mathrm{POD}}} \tilde{s}_k^2$ becomes smaller than a given target accuracy $\varepsilon$ already for some $n_\varepsilon \ll \dim \mathbb{U}_h$ of moderate size. In this case

$$\mathbb{U}_{n_\varepsilon} := \mathrm{span}\,\{u_k : k = 1, \ldots, n_\varepsilon\}, \tag{13}$$

with the POD bases $u_k$, $k = 1, \ldots, n_\varepsilon$, serves as a *reduced space*. More details will be provided later below for the respective examples.

One should keep in mind that a rapid decay of (11) requires a significantly smaller number of high-fidelity approximate solutions of (1) than a purely data-driven approach, based on (8).

Given $\mathbb{U}_{n_\varepsilon}$ we will then seek approximations to the solutions $u$ of (1) in the "hybrid" form

$$u(x, \mathfrak{p}; \theta) := \sum_{k=1}^{n_\varepsilon} \phi_k(\mathfrak{p}; \theta) u_k(x). \tag{14}$$

Aside from facilitating optimization, this format offers advantages in handling boundary conditions. As mentioned earlier, one major objective is to incorporate non-trivial mixed Dirichlet-Neumann conditions in the training objective in a variationally correct way. First, due to the nature of the POD basis functions $u_k$, it is easy to incorporate essential Dirichlet conditions. As detailed below for both models, the stationary diffusion equation and linear elasticity, the first-order formulations allow us, in particular, to incorporate Neumann conditions by an $L_2$ source in the least-squares residual, which is determined via solving a single auxiliary second-order elliptic problem.

# 3 First-order system least-squares loss formulation

## 3.1 Stationary diffusion equation with mixed boundary conditions

### 3.1.1 Second order formulation

Let $\Omega \subset \mathbb{R}^d$ be a bounded domain. The strong form of a stationary diffusion equation with heterogeneous diffusivity $\mathfrak{p}(x) \in \mathbb{R}^{d \times d}$ and mixed boundary conditions (Dirichlet and Neumann) is given by

$$\begin{cases} -\mathrm{div}\,(\mathfrak{p}\nabla u) = f, & \text{in } \Omega, \\ u = u_0, & \text{on } \Gamma_D, \\ \mathfrak{p}\nabla u \cdot n = g, & \text{on } \Gamma_N, \end{cases} \tag{15}$$

where $u(x) \in \mathbb{R}$ is the unknown scalar field and $f$ is a given source term with properties specified below. The Dirichlet boundary is denoted by $\Gamma_D \subset \partial\Omega$ with prescribed boundary data $u_0$. The Neumann boundary is given by $\Gamma_N = \partial\Omega \setminus \Gamma_D$, where $n = n(x)$ is (for almost all $x \in \partial\Omega$) the outward unit normal at $x \in \partial\Omega$. Finally, $g$ stands for the prescribed normal flux.

To explain in which sense a first-order formulation is equivalent, we need to refer in both cases to an appropriate weak formulation. To that end, we recapitulate for the convenience of the reader, a few classical facts, and define as usual

$$\mathbb{U} := H^1_{0,\Gamma_D}(\Omega) = \mathrm{clos}_{H^1(\Omega)} \{\phi \in C^\infty(\Omega) : \mathrm{supp}\, \phi \cap \Gamma_D = \emptyset\},$$

where we adopt common terminology when abbreviating $H^1_0(\Omega) := H^1_{0,\partial\Omega}(\Omega)$ when the Dirichlet boundary $\Gamma_D$ is all of $\partial\Omega$. Moreover, let $\mathbb{U}' := (H^1_{0,\Gamma_D}(\Omega))'$ denote its normed dual, i.e., the space of all bounded linear functionals on $H^1_{0,\Gamma_D}(\Omega)$. The classical weak formulation of (15) then reads: given $f \in \mathbb{U}'$ find $u \in H^1(\Omega)$ satisfying

$$u|_{\Gamma_D} = u_0 \quad \text{and} \quad b_{\mathfrak{p}}(u,v) := \int_\Omega \mathfrak{p}\nabla u \cdot \nabla v \, dx = f(v) + \langle g, v\rangle_{\Gamma_N}, \quad v \in \mathbb{U}, \tag{16}$$

where the boundary constraint is understood in the sense of traces. Moreover, $\langle \cdot, \cdot\rangle_{\Gamma_N}$ denotes the duality pairing on $H^{-1/2}(\Gamma_N) \times H^{1/2}_{00}(\Gamma_N)$, where $H^{-1/2}(\Gamma_N) := (H^{1/2}_{00}(\Gamma_N))'$. Here, $H^{1/2}_{00}(\Gamma_N)$ is the subspace of $H^{1/2}(\Gamma_N)$ that consists of those elements whose extension by zero to the rest of $\partial\Omega$ belongs to $H^{1/2}(\partial\Omega)$. Note that $v \in \mathbb{U}$ implies $v|_{\Gamma_N} \in H^{1/2}_{00}(\Gamma_N)$, so that $\langle g, \cdot\rangle_{\Gamma_N} : \mathbb{U} \to \mathbb{R}$ indeed defines a bounded linear functional on $\mathbb{U}$, hence belongs to $\mathbb{U}'$. The operator induced by this weak formulation is affine.

A proper weak first-order formulation requires a different representation of boundary conditions, which will follow from a different but equivalent alternate weak formulation of (15). This will be obtained by solving two auxiliary parameter-independent elliptic problems representing harmonic extensions for both types of boundary data and taking a Riesz representation of the right-hand side functional.

Regarding the Dirichlet conditions, choose a fixed $w \in H^1(\Omega)$ such that $w|_{\Gamma_D} = u_0$ and $\nabla w \cdot n|_{\Gamma_N} = 0$ (in the sense of traces), by solving the parameter-independent problem: find $w \in H^1(\Omega)$ such that

$$(\nabla w, \nabla v)_\Omega = 0, \quad v \in H^1(\Omega), \ w|_{\Gamma_D} = u_0. \tag{17}$$

Then consider the equivalent variational problem: find $u^\circ \in \mathbb{U}$ such that

$$b_{\mathfrak{p}}(u^\circ, v) = -b_{\mathfrak{p}}(w, v) + f(v) + \langle g, v\rangle_{\Gamma_N}, \quad v \in \mathbb{U}, \tag{18}$$

so that $u := u^\circ + w$ obviously satisfies (16). Note that $-b_{\mathfrak{p}}(w, \cdot) + f + \langle g, \cdot\rangle_{\Gamma_N}$ defines a bounded linear functional on $\mathbb{U}$ - as needed to render (16) well-posed - provided that $f \in \mathbb{U}'$.

As such, $f$ could, in principle, incorporate an additional flux-condition on $\Gamma_N$, violating the explicit specification in terms of $g$. To avoid this we confine $f$ to be *flux-free* which in the present context means that

$$f(v) = (f_2, v)_\Omega + (\mathrm{div}\, f_1)(v) = (f_2, v)_\Omega - (f_1, \nabla v)_\Omega, \quad v \in \mathbb{U}. \tag{19}$$

where $f_2 \in L_2(\Omega)$ and $f_1 \in L_2(\Omega; \mathbb{R}^d)$. We refer to a detailed discussion in Appendix A that explains in which sense the implied condition $\langle f_1 \cdot n, v\rangle_{\Gamma_N} = 0$ for $v \in \mathbb{U}$, is to be understood.

Regarding the Neumann condition, consider the second similar auxiliary problem: find $q \in \mathbb{U}$ such that

$$(\nabla q, \nabla v)_\Omega = \langle g, v\rangle_{\Gamma_N}, \ v \in \mathbb{U}, \tag{20}$$

which has a unique solution $q \in \mathbb{U}$ satisfying the distributional relation

$$-(\Delta q)(v) = 0, \qquad \langle \nabla q \cdot n, v\rangle_{\Gamma_N} = \langle g, v\rangle_{\Gamma_N}, \quad v \in \mathbb{U}. \tag{21}$$

Setting $z := \nabla q \in L^2(\Omega; \mathbb{R}^d)$, we have $\langle g, v\rangle_{\Gamma_N} = (z, \nabla v) = -(\mathrm{div}\, z)(v) + \langle z \cdot n, v\rangle_{\Gamma_N} = \langle z \cdot n, v\rangle_{\Gamma_N}$.

In summary, we obtain an equivalent weak formulation to (16): find $u^\circ \in \mathbb{U}$ such that

$$(\mathfrak{p}\nabla u^\circ + \mathfrak{p}\nabla w - z + f_1, \nabla v)_\Omega = (f_2, v)_\Omega, \quad v \in \mathbb{U}, \tag{22}$$

where $u = u^\circ + w$ is the solution (16) and the normal trace term is realized through the field $z$.

7

### 3.1.2 First-order system least-squares formulation

Introducing the auxiliary flux variable

$$\sigma^\circ := \mathfrak{p}\nabla u^\circ + \mathfrak{p}\nabla w - z + f_1, \tag{23}$$

(22) takes the simple form

$$(\sigma^\circ, \nabla v)_\Omega = -(\operatorname{div} \sigma^\circ)(v) + \langle \sigma^\circ \cdot n, v\rangle_{\Gamma_N} = (f_2, v)_\Omega, \quad v \in \mathbb{U}. \tag{24}$$

Testing first with $v \in H_0^1(\Omega)$, shows that $-(\operatorname{div} \sigma^\circ)(v) = (f_2, v)_\Omega$, $v \in H_0^1(\Omega)$. Since $H_0^1(\Omega)$ is dense in $L_2(\Omega)$, we conclude that $-\operatorname{div} \sigma^\circ = f_2$ in $L_2(\Omega)$, hence $\sigma^\circ \in H(\operatorname{div};\Omega)$. Testing subsequently with all $v \in \mathbb{U}$, yields $\langle \sigma^\circ \cdot n, v\rangle_{\Gamma_N} = 0$ for all $v \in \mathbb{U}$, which says that $\sigma^\circ \cdot n|_{\Gamma_N} = 0$ in $H^{-1/2}(\Gamma_N)$.

This shows that we have to seek $\sigma^\circ$ in $\Sigma$ and $(\operatorname{div} \sigma^\circ + f_2, v)_\Omega$ is well-defined for $v \in L_2(\Omega)$, where

$$\Sigma := \{\eta \in H(\operatorname{div};\Omega) : \eta \cdot n|_{\Gamma_N} = 0\}. \tag{25}$$

Hence, we arrive at the well-known FOSLS formulation from [16, 17]: find $(\sigma^\circ, u^\circ) \in \mathbb{H} := \Sigma \times \mathbb{U}$ such that

$$(\sigma^\circ - \mathfrak{p}\nabla u^\circ, \tau)_\Omega = (F, \tau)_\Omega, \quad -(\operatorname{div} \sigma^\circ, v)_\Omega = (f_2, v)_\Omega, \quad (\tau, v) \in \mathbb{L}_2 := L_2(\Omega; \mathbb{R}^{d+1}), \tag{26}$$

where $F := \mathfrak{p}\nabla w - z + f_1 \in L_2(\Omega; \mathbb{R}^d)$. Defining the operator

$$(\mathcal{B}_\mathfrak{p}[\sigma^\circ, u^\circ])([\tau, v]) := (\sigma^\circ - \mathfrak{p}\nabla u^\circ, \tau)_\Omega - (\operatorname{div} \sigma^\circ, v)_\Omega, \quad [\tau, v] \in \mathbb{L}_2 = L_2(\Omega; \mathbb{R}^{d+1}), \tag{27}$$

it has been shown in [16, 17] that $\mathcal{B}_\mathfrak{p}$ is an *isomorphism* from $\mathbb{H}$ onto $\mathbb{L}_2 := L_2(\Omega; \mathbb{R}^{d+1})$, i.e.,

$$\|\mathcal{B}_\mathfrak{p}[\sigma', u']\|_{\mathbb{L}_2}^2 \eqsim \|[\sigma', u']\|_{\mathbb{H}}^2 := \|\sigma'\|_\Sigma^2 + \|u'\|_{\mathbb{U}}^2, \quad \forall [\sigma', u'] \in \mathbb{H}, \tag{28}$$

where the proportionality constants depend on lower and upper bounds on $\mathfrak{p} \in \mathfrak{P}$ in $\Omega$. For the convenience of the reader, we provide a proof adapted to the parametric case in [Appendix B.1](#).

For any approximation $[\tilde\sigma^\circ, \tilde u^\circ] \in \mathbb{H}$ of the exact solution $[\sigma^\circ, u^\circ] \in \mathbb{H}$ at any parameter $\mathfrak{p} \in \mathfrak{P}$, applying this to the error $[\sigma', u'] = [\tilde\sigma^\circ - \sigma^\circ, \tilde u^\circ - u^\circ]$ yields the desired error-residual relation

$$\|[\tilde\sigma^\circ, \tilde u^\circ] - [\sigma^\circ, u^\circ]\|_{\mathbb{H}}^2 \eqsim \mathcal{L}([\tilde\sigma^\circ, \tilde u^\circ]; \mathfrak{p}), \quad [\tilde\sigma^\circ, \tilde u^\circ] \in \mathbb{H}, \mathfrak{p} \in \mathfrak{P}, \tag{29}$$

where the residual fiber-loss $\mathcal{L}([\tilde\sigma^\circ, \tilde u^\circ]; \mathfrak{p})$ of the approximate solution $[\tilde\sigma^\circ, \tilde u^\circ]$ at parameter $\mathfrak{p}$ is given by

$$\mathcal{L}([\tilde\sigma^\circ, \tilde u^\circ]; \mathfrak{p}) := \|\tilde\sigma^\circ - (\mathfrak{p}\nabla \tilde u^\circ + \mathfrak{p}\nabla w - z + f_1)\|_{L_2(\Omega; \mathbb{R}^d)}^2 + \|\operatorname{div} \tilde\sigma^\circ + f_2\|_{L_2(\Omega)}^2. \tag{30}$$

Note that the loss function involves only $L_2$-norms and no explicit boundary term in a broken or dual norm. It is therefore computable and can be used to form the mean squared loss $\mathcal{L}([\tilde\sigma^\circ, \tilde u^\circ]; \widehat{\mathfrak{P}})$ in (6), on which our surrogate models will be trained:

$$\theta^* \in \operatorname*{argmin}_{\theta \in \Theta} \frac{1}{\#\widehat{\mathfrak{P}}} \sum_{\mathfrak{p} \in \widehat{\mathfrak{P}}} \|\sigma^\circ(\mathfrak{p}; \theta) - (\mathfrak{p}\nabla u^\circ(\mathfrak{p}; \theta) + \mathfrak{p}\nabla w - z + f_1)\|_{L^2(\Omega; \mathbb{R}^d)}^2 + \|\operatorname{div} \sigma^\circ(\mathfrak{p}; \theta) + f_2\|_{L^2(\Omega)}^2, \tag{31}$$

To put the above developments into perspective, while the training loss does not contain any explicit boundary term, the homogeneous boundary conditions on $\sigma^\circ$ and $u^\circ$ need to be built into the trial system. Since we will be using a hybrid representation system, the functions $\sigma^\circ(\cdot, \mathfrak{p}; \theta)$ and $u^\circ(\cdot, \mathfrak{p}; \theta)$, as functions of the spatial variables $x$, are piecewise polynomial, such boundary conditions are easy to realize. Second, the loss is variationally correct, i.e., its size bounds the error in the trial norm from below and above.

Alternatively, recombining $u^\circ$ and $w$, one can omit the source term $w$ and absorb the data term $\mathfrak{p}\nabla w$ back in $u = u^\circ + w$. In this case, one has to enforce the inhomogeneous boundary condition $u|_{\Gamma_D} = u_0$ in the trial space so that the POD reduced space needs to be affine. Likewise, $\sigma := \sigma^\circ + z$ would satisfy the correct Neumann conditions so that one might consider a corresponding affine subspace of $H(\operatorname{div};\Omega)$ as well.

Instead of incorporating any boundary conditions into the trial system, it is common practice in machine learning to enforce boundary conditions by penalizing boundary residuals. Again, to preserve variational correctness (the loss is a sharp error bound), one could consider the residual with respect to an "extended" operator

$$\hat{\mathcal{B}}_{\mathfrak{p}} : H(\operatorname{div};\Omega) \times H^1(\Omega) \to \mathbb{L}_2 \times H^{-1/2}(\Gamma_N) \times H_{00}^{1/2}(\Gamma_D), \quad \hat{\mathcal{B}}_{\mathfrak{p}}[\sigma', u'] := \begin{pmatrix} \mathcal{B}_{\mathfrak{p}}[\sigma', u'] \\ \sigma' \cdot n|_{\Gamma_N} \\ u'|_{\Gamma_D} \end{pmatrix}. \tag{32}$$

If one can show that this is also an isomorphism, a corresponding variationally correct fiber-residual relation would now read for each $\mathfrak{p} \in \mathfrak{P}$

$$\|[\tilde{\sigma}, \tilde{u}] - [\sigma, u]\|_{H(\operatorname{div};\Omega) \times H^1(\Omega)}^2 \eqsim \|\tilde{\sigma} - \mathfrak{p}\nabla \tilde{u} - f_1\|_{L_2(\Omega;\mathbb{R}^d)}^2 + \|\operatorname{div}\tilde{\sigma} + f_2\|_{L_2(\Omega)}^2$$
$$+ \|\tilde{u} - u_0\|_{H^{1/2}(\Gamma_D)}^2 + \|\tilde{\sigma} \cdot n - g\|_{H^{-1/2}(\Gamma_N)}^2,$$
$$[\tilde{\sigma}, \tilde{u}] \in H(\operatorname{div};\Omega) \times H^1(\Omega). \tag{33}$$

Thus, the residual now involves a broken Sobolev and dual Sobolev trace-norm whose evaluation is far from straightforward (although possible for the Dirichlet conditions using the intrinsic definitions as a double integral over $\Gamma_D$). Simply incorporating the Dirichlet boundary conditions by adding the term $\|\tilde{u} - u_0\|_{L^2(\Gamma_D)}^2$ instead of $\|\tilde{u} - u_0\|_{H^{1/2}(\Gamma_D)}^2$ is **not** variationally correct, i.e., does **not** provide a rigorous error bound. Replacing $\|\tilde{\sigma} \cdot n - g\|_{H^{-1/2}(\Gamma_N)}^2$ by $\|\tilde{\sigma} \cdot n - g\|_{L_2(\Gamma_N)}^2$ would at least yield an upper bound for the error, but may still, depending on the data, ask too much regularity.

In summary, the above findings will be used as follows. Given a hypothesis class $\mathcal{H}(\theta)$, comprised of functions $[\sigma(\cdot;\theta), u(\cdot;\theta)]$ that depend on trainable weights $\theta \in \Theta$, we can now follow the lines outlined at the end of Section 2. Most importantly, in the present situation, we have $\mathbb{V} = L_2(\Omega;\mathbb{R}^{d+1})$ so that the losses in (6) and (31), respectively, are computable.

## 3.2 Linear elasticity equation

We turn now to the classical model for linear elasticity. We consider a convex polygonal domain $\Omega \subset \mathbb{R}^d$ (typically $d = 1, 2, 3$ for practical scenarios). While plain letters denote scalar-valued functions, we distinguish vectors and (rank-2) tensors by single or double underscores such as

$$\underline{f} = (f_1, f_2, \ldots, f_d) \text{ and } \underline{\underline{f}} = \begin{pmatrix} f_{11} & f_{12} & \cdots & f_{1d} \\ f_{21} & f_{22} & \cdots & f_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ f_{d1} & f_{d2} & \cdots & f_{dd} \end{pmatrix}.$$

$$\operatorname{grad} f = \left( \frac{\partial f}{\partial x_1}, \quad \frac{\partial f}{\partial x_2}, \quad \cdots, \quad \frac{\partial f}{\partial x_d} \right), \quad \operatorname{div}\underline{f} = \sum_{i=1}^d \frac{\partial f_i}{\partial x_i}, \quad \underline{\underline{\operatorname{grad}}}\,\underline{f} = \begin{pmatrix} \underline{\operatorname{grad}} f_1 \\ \underline{\operatorname{grad}} f_2 \\ \vdots \\ \underline{\operatorname{grad}} f_d \end{pmatrix}, \quad \underline{\operatorname{div}}\,\underline{\underline{f}} = \begin{pmatrix} \operatorname{div}\underline{f_1} \\ \operatorname{div}\underline{f_2} \\ \vdots \\ \operatorname{div}\underline{f_d} \end{pmatrix}.$$

The vector-vector and tensor-tensor inner products "$\cdot$" respectively "$:$" are then given by

$$\underline{f} \cdot \underline{g} = \sum_{i=1}^d f_i g_i, \quad \underline{\underline{f}} : \underline{\underline{g}} = \sum_{i=1}^d \sum_{j=1}^d f_{ij} g_{ij},$$

and the tensor-vector inner product reads

$$\underline{\underline{f}} \cdot \underline{g} = \left( \underline{f_1} \cdot \underline{g}, \underline{f_2} \cdot \underline{g}, \ldots, \underline{f_d} \cdot \underline{g} \right)^T.$$

9

In these terms, the strong formulation of a linear elasticity equation with mixed boundary conditions (Dirichlet and Neumann) can be written as

$$\begin{cases} -\underline{\mathrm{div}}\,\underline{\underline{\sigma}}(\underline{u}) = \underline{f}, & \text{in } \Omega, \\ \underline{u} = \underline{u}_0, & \text{on } \Gamma_D, \\ \underline{\underline{\sigma}}(\underline{u}) \cdot \underline{n} = \underline{t}, & \text{on } \Gamma_N, \end{cases} \tag{34}$$

where the Cauchy stress tensor $\underline{\underline{\sigma}}$ is a function of the displacement vector field $\underline{u}$

$$\underline{\underline{\sigma}}(\underline{u}) = 2\mu\,\underline{\underline{\varepsilon}}(\underline{u}) + \lambda\,\mathrm{tr}(\underline{\underline{\varepsilon}}(\underline{u}))\underline{\underline{I}}_d. \tag{35}$$

Here $\lambda, \mu \in \mathbb{R}$ are (possibly spatially dependent) parameters associated with the material, $\underline{\underline{I}}_d$ is the $d \times d$ identity matrix, and $\underline{\underline{\varepsilon}}(\underline{u}) = \frac{1}{2}(\underline{\mathrm{grad}}\,\underline{u} + (\underline{\mathrm{grad}}\,\underline{u})^T)$ is the strain tensor. Note that in these terms, we have

$$\int_\Omega (\underline{\mathrm{div}}\,\underline{\underline{\sigma}}) \cdot \underline{v}\,dx = -\int_\Omega \underline{\underline{\sigma}} : \underline{\mathrm{grad}}\,\underline{v}\,dx + \int_{\partial\Omega} (\underline{\underline{\sigma}} \cdot \underline{n}) \cdot \underline{v}\,d\gamma. \tag{36}$$

As in the previous section, to properly handle boundary conditions, we start with the standard weak formulation of (34): find $\underline{u} \in H^1(\Omega; \mathbb{R}^d)$ such that

$$b(\underline{u}, \underline{v}) := \int_\Omega \underline{\underline{\sigma}}(\underline{u}) : \underline{\mathrm{grad}}\,\underline{v}\,dx = \int_{\Gamma_N} \underline{t} \cdot \underline{v}\,d\gamma + \int_\Omega \underline{f} \cdot \underline{v}\,dx, \quad v \in \mathbb{U} = H^1_{0,\Gamma_D}(\Omega; \mathbb{R}^d), \tag{37}$$

where we assume here for convenience that $\underline{f} \in L^2(\Omega; \mathbb{R}^d)$. Otherwise, we would need to consider a decomposition of flux-free $\underline{f} \in (H^1_{0,\Gamma_D}(\Omega; \mathbb{R}^d))'$, by applying the scalar-valued decompositions (19) componentwise,

$$\underline{f} = \underline{f}_2 + \underline{\mathrm{div}}\,\underline{\underline{f}}_1, \tag{38}$$

where $\underline{f}_2 \in L_2(\Omega; \mathbb{R}^d)$ and $\underline{\underline{f}}_1 \in L_2(\Omega; \mathbb{R}^{d \times d})$.

As in the previous section, we wish to represent the boundary conditions as source terms. Regarding the displacement condition, let $\underline{w} \in H^1(\Omega; \mathbb{R}^d)$ be any fixed function satisfying $\underline{w}|_{\Gamma_D} = \underline{u}_0|_{\Gamma_D}$ and $\underline{\mathrm{grad}}\,\underline{w} \cdot \underline{n} = \underline{0}$ on $\Gamma_N$ (componentwise in the sense of traces). For instance, one can solve

$$a(\underline{w}, \underline{v}) := \int_\Omega \underline{\mathrm{grad}}\,\underline{w} : \underline{\mathrm{grad}}\,\underline{v}\,dx = 0, \quad \underline{w}|_{\Gamma_D} = \underline{u}_0|_{\Gamma_D}, \tag{39}$$

approximately in some finite element space. Then (37) is equivalent to finding $\underline{u}^\circ \in \mathbb{U}$ such that

$$b(\underline{u}^\circ, \underline{v}) = -b(\underline{w}, \underline{v}) + \int_{\Gamma_N} \underline{t} \cdot \underline{v}\,d\gamma + \int_\Omega \underline{f} \cdot \underline{v}\,dx, \quad v \in \mathbb{U}. \tag{40}$$

As in the previous section, $\int_{\Gamma_N} \underline{t} \cdot \underline{v}\,d\gamma$ is to be understood as a dual pairing, denoted by $\langle \underline{t}, \underline{v} \rangle_{\Gamma_N}$. Likewise, we abbreviate for convenience $(\underline{f}, \underline{v})_\Omega := \int_\Omega \underline{f} \cdot \underline{v}\,dx$. We can view $-b(\underline{w}, \cdot)$ as a functional in $\mathbb{U}'$ acting on $v \in \mathbb{U}$, namely, for $\underline{w}$ as above we have (see (36))

$$-b(\underline{w}, \underline{v}) = \underline{\mathrm{div}}(\underline{\underline{\sigma}}(\underline{w}))(\underline{v}), \tag{41}$$

and hence is the divergence of an $L_2(\Omega; \mathbb{R}^{d \times d})$-tensor field. To represent the normal-trace integral in $\mathbb{U}'$, we proceed as before and consider the auxiliary problem: find $\underline{q} \in \mathbb{U} = H^1_{0,\Gamma_D}(\Omega; \mathbb{R}^d)$ such that

$$a(\underline{q}, \underline{v}) = \int_{\Gamma_N} \underline{t} \cdot \underline{v}\,d\gamma. \tag{42}$$

In strong form, this solves $-\underline{\mathrm{div}}\,\underline{\mathrm{grad}}\,(\underline{q}) = 0$ with $\underline{\mathrm{grad}}\,\underline{q} \cdot \underline{n} = \underline{t}$ on $\Gamma_N$ and $\underline{q} = 0$ on $\Gamma_D$.

Hence, (recalling that $\underline{f} \in L_2(\Omega; \mathbb{R}^d)$), in complete analogy to the previous section, a weak formulation of (34) can be stated with $\underline{z} := \underline{\operatorname{grad}}\, q$ as

$$b(\underline{u}^\circ + \underline{w}, \underline{v}) = (\underline{f}, \underline{v})_\Omega + a(q, \underline{v}), \quad \underline{u}^\circ|_{\Gamma_D} = 0, \tag{43}$$

which is a linear elasticity system with homogeneous displacement boundary conditions on boundary $\Gamma_D$. In subsequent experiments, we consider the parametric linear elasticity equation with the following parametrization: $(\lambda, \mu) : \Omega \to \mathbb{R}^2$ are the Lambda parameters depending on the Young's modulus $E : \Omega \to \mathbb{R}$ and Poisson ratio $\nu \in \mathbb{R}$ as

$$\mu(x) = \frac{E(x)}{2(1+\nu)}, \quad \lambda(x) = \frac{\nu E(x)}{(1+\nu)(1-2\nu)}.$$

We assume the Young's modulus is a random field parameter given by $E(x) = \exp\{\mathfrak{p}(x)\} + 1$, where $\mathfrak{p}$ is a random field with Gaussian measure $\mathcal{N}(\bar{\mathfrak{p}}, \mathcal{C})$ with mean $\bar{\mathfrak{p}}$ and covariance operator $\mathcal{C} := (\delta I - \gamma \Delta)^{-\alpha}$, with $\delta, \gamma, \alpha > 0$ collectively determining the correction, variance, and smoothness. We restrict our analysis to random field parameters satisfying $\mathfrak{p} \leq \hat{\mathfrak{p}}$ for a sufficiently large constant $\hat{\mathfrak{p}}$. This ensures that the Young's modulus $E$ is uniformly bounded, with $1 \leq E \leq \exp(\hat{\mathfrak{p}}) + 1$. We denote the parametric linear mapping $\mathcal{C}_\mathfrak{p}$, a rank-4 stiffness tensor, that maps the strain $\underline{\underline{\varepsilon}}$ to the stress $\underline{\underline{\sigma}}$ as

$$\underline{\underline{\sigma}} = \mathcal{C}_\mathfrak{p} \underline{\underline{\varepsilon}} := 2\mu \underline{\underline{\varepsilon}} + \lambda \operatorname{tr}(\underline{\underline{\varepsilon}}) I_d.$$

Under our assumptions on the range of parameters, its inverse $\mathcal{C}_\mathfrak{p}^{-1}$ is given by

$$\mathcal{C}_\mathfrak{p}^{-1} \underline{\underline{\sigma}} = \frac{1}{2\mu} \underline{\underline{\sigma}} - \frac{\lambda}{2\mu(\lambda d + 2\mu)} \operatorname{tr}(\underline{\underline{\sigma}}) I_d. \tag{44}$$

The following first-order formulation of Equation (43) is obtained by treating the stress tensor $\underline{\underline{\sigma}}$ as an independent unknown variable and (in its preconditioned form) reads (recall that $\underline{z} = \underline{\operatorname{grad}}\, q$)

$$\begin{cases} -\underline{\operatorname{div}}\, \underline{\underline{\sigma}}^\circ = \underline{f}, & \text{in } \Omega, \\ \mathcal{C}_\mathfrak{p}^{-1/2} \underline{\underline{\sigma}}^\circ = \mathcal{C}_\mathfrak{p}^{1/2} \underline{\underline{\varepsilon}}(\underline{u}^\circ + \underline{w}) - \mathcal{C}_\mathfrak{p}^{-1/2} \underline{\underline{z}}, & \text{in } \Omega, \\ \underline{u}^\circ = 0, & \text{on } \Gamma_D, \\ \underline{\underline{\sigma}}^\circ \cdot \underline{n} = \underline{0}, & \text{on } \Gamma_N. \end{cases} \tag{45}$$

In this form we can apply the results in [16, 17, 20] to arrive at the weak (fiber) formulation: find $\underline{\underline{\sigma}}^\circ \in \Sigma := H_{0,\Gamma_N}(\operatorname{div}; \Omega; \mathbb{R}^{d \times d})$ and $\underline{u}^\circ \in \mathbb{U} := H_{0,\Gamma_D}^1(\Omega; \mathbb{R}^d)$ such that

$$b_\mathfrak{p}([\underline{\underline{\sigma}}^\circ, \underline{u}^\circ], [\underline{\underline{\tau}}, \underline{v}]) := (-\underline{\operatorname{div}}\, \underline{\underline{\sigma}}^\circ, \underline{v})_\Omega + (\mathcal{C}_\mathfrak{p}^{-1/2}(\underline{\underline{\sigma}}^\circ - \mathcal{C}_\mathfrak{p} \underline{\underline{\varepsilon}}(\underline{u}^\circ)), \underline{\underline{\tau}})_\Omega$$
$$= (\underline{f}, \underline{v})_\Omega + (\mathcal{C}_\mathfrak{p}^{1/2} \underline{\underline{\varepsilon}}(\underline{w}) - \mathcal{C}_\mathfrak{p}^{-1/2} \underline{\underline{z}}, \underline{\underline{\tau}})_\Omega, \quad \underline{\underline{\tau}} \in L^2(\Omega; \mathbb{R}^{d \times d}), \underline{v} \in L^2(\Omega; \mathbb{R}^d). \tag{46}$$

We define as in the previous section the operator $\mathcal{B}_\mathfrak{p} : \mathbb{H} := \Sigma \times \mathbb{U} \to \mathbb{L}_2 := L_2(\Omega; \mathbb{R}^{d \times d}) \times L_2(\Omega; \mathbb{R}^d)$ by $(\mathcal{B}_\mathfrak{p}[\underline{\underline{\sigma}}^\circ, \underline{u}^\circ])([\underline{\underline{\tau}}, \underline{v}]) = b_\mathfrak{p}([\underline{\underline{\sigma}}^\circ, \underline{u}^\circ], [\underline{\underline{\tau}}, \underline{v}])$ for all $([\underline{\underline{\sigma}}^\circ, \underline{u}^\circ], [\underline{\underline{\tau}}, \underline{v}]) \in \mathbb{H} \times \mathbb{L}_2$, i.e.,

$$\mathcal{B}_\mathfrak{p}[\underline{\underline{\sigma}}^\circ, \underline{u}^\circ] = \begin{pmatrix} \mathcal{C}_\mathfrak{p}^{-1/2}(\underline{\underline{\sigma}}^\circ - \mathcal{C}_\mathfrak{p} \underline{\underline{\varepsilon}}(\underline{u}^\circ)) \\ -\underline{\operatorname{div}}\, \underline{\underline{\sigma}}^\circ \end{pmatrix}.$$

It has been shown in [16, 17, 20] that $\mathcal{B}_\mathfrak{p}$ is for each $\mathfrak{p} \in \mathfrak{P}$ an isomorphism. Corresponding inf-sup- or continuity constants depend in general on $\mathfrak{p}$ so that any uniform stability depends on the range permitted in $\mathfrak{P}$. Similarly to the error-residual relation (29) for the case of the stationary diffusion problem, we obtain,

11

for any approximation $[\underline{\tilde{\sigma}}^\circ, \tilde{u}^\circ]$ of the exact solution $[\underline{\sigma}^\circ, \underline{u}^\circ]$ at any given parameter $\mathfrak{p}$, the fiber-error-residual relation

$$\|[\underline{\tilde{\sigma}}^\circ, \tilde{u}^\circ] - [\underline{\sigma}^\circ, \underline{u}^\circ]\|_{\mathbb{H}}^2 \eqsim \mathcal{L}([\underline{\tilde{\sigma}}^\circ, \tilde{u}^\circ]; \mathfrak{p}) \quad \forall [\underline{\tilde{\sigma}}^\circ, \tilde{u}^\circ] \in \mathbb{H}, \ \mathfrak{p} \in \mathfrak{P}, \tag{47}$$

where the residual $\mathcal{L}([\underline{\tilde{\sigma}}^\circ, \tilde{u}^\circ]; \mathfrak{p})$ of the approximation $[\underline{\tilde{\sigma}}^\circ, \tilde{u}^\circ]$ at parameter $\mathfrak{p}$ is given by

$$\mathcal{L}([\underline{\tilde{\sigma}}^\circ, \tilde{u}^\circ]; \mathfrak{p}) := \|\mathcal{C}_{\mathfrak{p}}^{-1/2}\underline{\tilde{\sigma}}^\circ - \mathcal{C}_{\mathfrak{p}}^{1/2}\underline{\varepsilon}(\tilde{u}^\circ + \underline{w}) + \mathcal{C}_{\mathfrak{p}}^{-1/2}\underline{z}\|_{L^2(\Omega;\mathbb{R}^{d\times d})}^2 + \|\underline{\mathrm{div}}\,\underline{\tilde{\sigma}}^\circ + \underline{f}\|_{L^2(\Omega;\mathbb{R}^d)}^2. \tag{48}$$

This is for each $\mathfrak{p} \in \mathfrak{P}$ an easily computable loss function that enters the corresponding mean-squared loss in (6) which is used to train the surrogate models.

Training yields then $[\underline{\sigma}^\circ(\theta^*), \underline{u}^\circ(\theta^*)] \in \mathbb{X} = L_\mu^2(\mathfrak{P}; \mathbb{H})$ from which we recover an approximation to the solution of (45) by

$$\tilde{u}(\mathfrak{p}) = \underline{u}^\circ(\mathfrak{p}; \theta^*)) + \underline{w}, \quad \underline{\tilde{\sigma}}(\mathfrak{p}) = \underline{\sigma}^\circ(\mathfrak{p}; \theta^*) + \underline{z}. \tag{49}$$

We immediately infer from [11] that the exact solution $[\sigma^\circ, u^\circ] \in \mathbb{X}$ of (45) satisfies

$$\|\underline{\sigma}^\circ(\theta^*), \underline{u}^\circ(\theta^*)] - [\underline{\sigma}^\circ, \underline{u}^\circ]\|_{\ell_2(\widehat{\mathfrak{P}};\mathbb{H})}^2 \eqsim \mathcal{L}([\underline{\sigma}^\circ(\theta^*), \underline{u}^\circ(\theta^*)]; \widehat{\mathfrak{P}}). \tag{50}$$

# 4 Conforming finite element approximations

To evaluate the residual losses (30) (diffusion) and (48) (elasticity), we employ conforming Galerkin FOSLS discretizations with mixed finite-element (FE) pairs $\Sigma_h \times \mathbb{U}_h$ and discrete boundary lifts. We formulate the discrete problems with lifted boundary conditions and, under uniform ellipticity, establish (i) variational error–residual equivalence and (ii) convergence rates for the FE losses evaluated at the FE solutions.

## 4.1 Conforming finite element approximation of the diffusion problem

For the diffusion problem and each parameter $\mathfrak{p} \in \mathfrak{P}$, we compute the conforming FE approximation $[\tilde{\sigma}^\circ, \tilde{u}^\circ] = [\sigma_h^\circ, u_h^\circ] \in \mathbb{H}_h$ of the lifted fiber solution as

$$\sigma_h^\circ(x, \mathfrak{p}) = \sum_{n=1}^{N_h^\sigma} \sigma_n^\circ(\mathfrak{p})\varphi_n^\sigma(x) \text{ and } u_h^\circ(x, \mathfrak{p}) = \sum_{n=1}^{N_h^u} u_n^\circ(\mathfrak{p})\varphi_n^u(x), \tag{51}$$

where $\boldsymbol{\sigma}^\circ(\mathfrak{p}) = (\sigma_1^\circ(\mathfrak{p}), \ldots, \sigma_{N_h^\sigma}^\circ(\mathfrak{p}))^\top$ and $\boldsymbol{u}^\circ(\mathfrak{p}) = (u_1^\circ(\mathfrak{p}), \ldots, u_{N_h^u}^\circ(\mathfrak{p}))^\top$ are the parametric coefficient vectors, $(\varphi_n^\sigma)_{n=1}^{N_h^\sigma}$ and $(\varphi_n^u)_{n=1}^{N_h^u}$ are the basis functions of the FE space $\mathbb{H}_h = \Sigma_h \times \mathbb{U}_h \subset H(\mathrm{div}\,; \Omega) \times H^1(\Omega)$, defined over a mesh with element size $h$. Here $N_h^\sigma$ and $N_h^u$ are the number of degrees of freedom (DoFs) of the FE spaces $\Sigma_h$ and $\mathbb{U}_h$, respectively. Specifically, we employ the conforming FE pairs

$$\Sigma_h \times \mathbb{U}_h = \mathrm{RT}_k^\circ \times \mathrm{CG}_m^\circ, \quad k \geq 0, \ m \geq 1,$$

where $\mathrm{RT}_k$ is the *Raviart–Thomas element* [42] of order $k$ (which has polynomial degree $\leq k+1$), $H(\mathrm{div}\,; \Omega)$-conforming vector fields with continuous normal fluxes (see, e.g., DefElement for details), and $\mathrm{CG}_m$ is the *continuous Galerkin element* (or Lagrange element) of degree $m$, with piecewise polynomials of degree $\leq m$ on each cell and $C^0$-continuous across cell interfaces, thus $H^1(\Omega)$-conforming scalar fields. Here $\mathrm{RT}_k^\circ$ and $\mathrm{CG}_m^\circ$ denote the subspaces with the essential boundary conditions incorporated, i.e., $\tau \cdot n = 0$ on $\Gamma_N$ for any $\tau \in \mathrm{RT}_k^\circ$ and $v = 0$ on $\Gamma_D$ for any $v \in \mathrm{CG}_m^\circ$, practically by setting the corresponding DoFs to zero.

For each $\mathfrak{p}$ we compute the FE approximation of the solution $[\sigma^\circ(\mathfrak{p}), u^\circ(\mathfrak{p})]$ by solving the Galerkin system

$$(\mathcal{B}_{\mathfrak{p}}[\sigma_h^\circ, u_h^\circ], \mathcal{B}_{\mathfrak{p}}[\tau_h, v_h])_\Omega = ([\mathfrak{p}\nabla w_h - z_h + f_1, f_2], \mathcal{B}_{\mathfrak{p}}[\tau_h, v_h])_\Omega, \quad \forall [\tau_h, v_h] \in \mathbb{H}_h, \tag{52}$$

12

which reads in explicit terms as: find $[\sigma_h^\circ, u_h^\circ] \in \Sigma_h \times \mathbb{U}_h$ such that

$$(\sigma_h^\circ - \mathfrak{p}\nabla u_h^\circ, \tau_h - \mathfrak{p}\nabla v_h)_\Omega + (\operatorname{div}\sigma_h^\circ, \operatorname{div}\tau_h)_\Omega$$
$$= (\mathfrak{p}\nabla w_h - z_h + f_1, \tau_h - \mathfrak{p}\nabla v_h)_\Omega - (f_2, \operatorname{div}\tau_h)_\Omega, \quad \forall[\tau_h, v_h] \in \mathbb{H}_h. \tag{53}$$

Here the auxiliary variables $w_h \in \mathrm{CG}_m^\circ$ and $q_h \in \mathrm{CG}_m^\circ$ are conforming Galerkin approximations of the lifts $w$ and $q$ from (17) and (20), respectively, on the same mesh and with the same polynomial degree $m$ as $u_h^\circ$; we set $z_h := \nabla q_h$. The FE loss is assembled with these discrete lifts $(w_h, z_h)$. Note that the solution of the Galerkin problem (53) satisfies the minimal residual property

$$[\sigma_h^\circ(\mathfrak{p}), u_h^\circ(\mathfrak{p})] = \operatorname*{argmin}_{[\tau_h, v_h] \in \mathbb{H}_h} \mathcal{L}([\tau_h, v_h]; \mathfrak{p}). \tag{54}$$

By construction, $[\sigma_h^\circ, u_h^\circ]$ minimizes the fiber loss over $\Sigma_h \times \mathbb{U}_h$, i.e., it is the orthogonal projection in the FOSLS inner product induced by $(\tau - \mathfrak{p}\nabla v, \eta - \mathfrak{p}\nabla s) + (\operatorname{div}\tau, \operatorname{div}\eta)$.

Under standard regularity assumptions for the exact fiber solution $[\sigma^\circ, u^\circ]$ and uniform ellipticity of the parameter $\mathfrak{p}$, the FE loss evaluated at the FE solution $[\sigma_h^\circ, u_h^\circ]$ satisfies the following properties; a proof is given in Appendix C.1.

**Theorem 1.** *Assume $\Omega$ is Lipschitz and $\mathfrak{p} \in L^\infty(\Omega)$ is uniformly bounded with $0 < \alpha \leq \mathfrak{p}(x) \leq \beta < \infty$ independent of $\mathfrak{p} \in \mathfrak{P}$. Let $\Sigma_h = \mathrm{RT}_k^\circ$ and $\mathbb{U}_h = \mathrm{CG}_m^\circ$ on a shape-regular mesh of size $h$ with $k \geq 0$, $m \geq 1$. Let $[\sigma^\circ, u^\circ]$ be the exact solution and $[\sigma_h^\circ, u_h^\circ]$ the Galerkin solution in $\Sigma_h \times \mathbb{U}_h$. Then, for each fixed $\mathfrak{p} \in \mathfrak{P}$, one has the error-residual equivalence*

$$\mathcal{L}([\sigma_h^\circ, u_h^\circ]; \mathfrak{p}) \approx \|\sigma^\circ - \sigma_h^\circ\|_{H(\operatorname{div};\Omega)}^2 + \|u^\circ - u_h^\circ\|_{H^1(\Omega)}^2, \tag{55}$$

*with equivalence constants depending only on $\alpha, \beta$ and the domain.*

*Moreover, if $\sigma^\circ \in H^{s_\sigma}(\Omega; \mathbb{R}^d)$, $\operatorname{div}\sigma^\circ \in H^{s_{\mathrm{div}}}(\Omega)$ with $s_\sigma, s_{\mathrm{div}} \geq 0$, and $u^\circ \in H^{s_u}(\Omega)$ with $s_u \geq 1$, then*

$$\mathcal{L}([\sigma_h^\circ, u_h^\circ]; \mathfrak{p}) \lesssim h^{2\min(k+1, s_\sigma)}\|\sigma^\circ\|_{H^{s_\sigma}}^2 + h^{2\min(k+1, s_{\mathrm{div}})}\|\operatorname{div}\sigma^\circ\|_{H^{s_{\mathrm{div}}}}^2 + h^{2\min(m, s_u - 1)}\|u^\circ\|_{H^{s_u}}^2. \tag{56}$$

*In particular, if $\sigma^\circ \in H^{k+1}$, $\operatorname{div}\sigma^\circ \in H^{k+1}$, and $u^\circ \in H^{m+1}$, then $\mathcal{L}([\sigma_h^\circ, u_h^\circ]; \mathfrak{p}) \lesssim h^{2(k+1)} + h^{2m}$, and the balanced choice $m = k + 1$ yields the optimal scaling $\mathcal{L}([\sigma_h^\circ, u_h^\circ]; \mathfrak{p}) \lesssim h^{2(k+1)}$.*

## 4.2 Conforming finite element approximation of the elasticity problem

For the elasticity problem and each parameter $\mathfrak{p} \in \mathfrak{P}$, we approximate the lifted fiber solution $[\underline{\underline{\sigma}}^\circ(\mathfrak{p}), \underline{u}^\circ(\mathfrak{p})]$ by a conforming Galerkin FOSLS discretization in the mixed FE space

$$\mathbb{H}_h = \Sigma_h \times \mathbb{U}_h \subset H_{0,\Gamma_N}(\operatorname{div};\Omega; \mathbb{R}^{d\times d}) \times H_{0,\Gamma_D}^1(\Omega; \mathbb{R}^d),$$

where we take

$$\Sigma_h \times \mathbb{U}_h = (\mathrm{RT}_k^\circ)^d \times (\mathrm{CG}_m^\circ)^d, \quad k \geq 0, \ m \geq 1.$$

Note the essential boundary conditions are built into the spaces: $\underline{\underline{\tau}}_h \cdot \underline{n} = \underline{0}$ on $\Gamma_N$ for $\underline{\underline{\tau}}_h \in \Sigma_h$ and $\underline{v}_h = \underline{0}$ on $\Gamma_D$ for $\underline{v}_h \in \mathbb{U}_h$. We parametrize the FE approximations $[\underline{\underline{\sigma}}_h^\circ(\mathfrak{p}), \underline{u}_h^\circ(\mathfrak{p})]$ in the same form as in (51) with the tensor- and vector-valued basis functions of $\Sigma_h$ and $\mathbb{U}_h$, respectively, with corresponding coefficient vectors $\boldsymbol{\sigma}^\circ(\mathfrak{p})$ and $\boldsymbol{u}^\circ(\mathfrak{p})$. Discrete boundary lifts are computed analogously to the diffusion case. Specifically, we compute $\underline{w}_h \in (\mathrm{CG}_m)^d$ as the vector-valued harmonic lift of the Dirichlet data by solving the auxiliary problem (39), and $\underline{q}_h \in (\mathrm{CG}_m^\circ)^d$ as the solution of the auxiliary problem (42); we then set the tensor lift $\underline{\underline{z}}_h := \underline{\underline{\operatorname{grad}}}\,\underline{q}_h$. The FE loss and right-hand side are assembled with these discrete lifts $(\underline{w}_h, \underline{\underline{z}}_h)$.

For each $\mathfrak{p} \in \mathfrak{P}$, we compute the FE solution $[\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ] \in \Sigma_h \times \mathbb{U}_h$ by solving the Galerkin system

$$(\mathcal{B}_\mathfrak{p}[\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ], \mathcal{B}_\mathfrak{p}[\underline{\underline{\tau}}_h, \underline{v}_h])_\Omega = ([\mathcal{C}_\mathfrak{p}^{1/2}\underline{\underline{\varepsilon}}(\underline{w}_h) - \mathcal{C}_\mathfrak{p}^{-1/2}\underline{\underline{z}}_h, \underline{f}], \mathcal{B}_\mathfrak{p}[\underline{\underline{\tau}}_h, \underline{v}_h])_\Omega, \quad \forall[\underline{\underline{\tau}}_h, \underline{v}_h] \in \Sigma_h \times \mathbb{U}_h, \tag{57}$$

which in expanded form reads: find $[\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ] \in \Sigma_h \times \mathbb{U}_h$ such that

$$(\mathcal{C}_{\mathfrak{p}}^{-1/2}(\underline{\underline{\sigma}}_h^\circ - \mathcal{C}_{\mathfrak{p}}\underline{\underline{\varepsilon}}(\underline{u}_h^\circ)), \mathcal{C}_{\mathfrak{p}}^{-1/2}(\underline{\underline{\tau}}_h - \mathcal{C}_{\mathfrak{p}}\underline{\underline{\varepsilon}}(\underline{v}_h)))_\Omega + (\mathrm{div}\,\underline{\underline{\sigma}}_h^\circ,\,\mathrm{div}\,\underline{\underline{\tau}}_h)_\Omega$$
$$= (\mathcal{C}_{\mathfrak{p}}^{1/2}\underline{\underline{\varepsilon}}(\underline{w}_h) - \mathcal{C}_{\mathfrak{p}}^{-1/2}\underline{\underline{z}}_h, \mathcal{C}_{\mathfrak{p}}^{-1/2}(\underline{\underline{\tau}}_h - \mathcal{C}_{\mathfrak{p}}\underline{\underline{\varepsilon}}(\underline{v}_h)))_\Omega - (\underline{f},\,\mathrm{div}\,\underline{\underline{\tau}}_h)_\Omega, \quad \forall[\underline{\underline{\tau}}_h, \underline{v}_h] \in \Sigma_h \times \mathbb{U}_h. \tag{58}$$

Equivalently, $[\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]$ minimizes the FE fiber loss over $\Sigma_h \times \mathbb{U}_h$, i.e.,

$$[\underline{\underline{\sigma}}_h^\circ(\mathfrak{p}), \underline{u}_h^\circ(\mathfrak{p})] = \operatorname*{argmin}_{[\underline{\underline{\tau}}_h, \underline{v}_h] \in \mathbb{H}_h} \mathcal{L}([\underline{\underline{\tau}}_h, \underline{v}_h]; \mathfrak{p}). \tag{59}$$

Note that the Cauchy stress tensor $\underline{\underline{\sigma}}$ in (35) is symmetric, while the FE approximation $\underline{\underline{\sigma}}_h = \underline{\underline{\sigma}}_h^\circ + \underline{\underline{z}}_h$, with $\underline{\underline{\sigma}}_h^\circ \in (\mathrm{RT}_k^\circ)^d$ and $\underline{\underline{z}}_h = \mathrm{grad}\,\underline{q}_h$ with $\underline{q}_h \in (\mathrm{CG}_m^\circ)^d$, is not necessarily symmetric. To address this, we can impose weak symmetry by adding a penalty term $\|\underline{\underline{\sigma}}_h - \underline{\underline{\sigma}}_h^\top\|_{L^2(\Omega)}^2$ to the FE loss, leading to corresponding modifications in the Galerkin system (57). We omit this for simplicity.

Under standard regularity assumptions for the exact fiber solution $[\underline{\underline{\sigma}}^\circ, \underline{u}^\circ]$ and uniform bounds for the elasticity tensor $\mathcal{C}_{\mathfrak{p}}$, the FE loss evaluated at the FE solution satisfies similar properties as in Theorem 1. We present the theorem and proof in Appendix C.2 for completeness.

# 5 Reduced basis neural operators

## 5.1 Reduced basis approximations

Given the high-dimensional nature of the FE coefficient vectors $\boldsymbol{\sigma}_h^\circ(\mathfrak{p}) \in \mathbb{R}^{N_h^\sigma}$ and $\boldsymbol{u}_h^\circ(\mathfrak{p}) \in \mathbb{R}^{N_h^u}$ in (51), directly parameterizing these full DoF representations with a neural network leads to very large models and expensive residual evaluations, especially for large-scale FE discretizations. Instead, we seek a low-dimensional representation of the parametric solution manifold that is conformal with the underlying variational structure and allows for efficient evaluation of the loss.

The following facts and notational conventions will be frequently used in what follows. For the parametric coefficient vector $\boldsymbol{s}_h(\mathfrak{p}) = (\boldsymbol{\sigma}_h^\circ(\mathfrak{p}); \boldsymbol{u}_h^\circ(\mathfrak{p})) \in \mathbb{R}^{N_h^s}$ of the FE functions $s_h(\mathfrak{p}) = [\sigma_h^\circ(\mathfrak{p}), u_h^\circ(\mathfrak{p})]$ in the diffusion case, we first equip $\mathbb{R}^{N_h^s}$ with the discrete $H(\mathrm{div}) \times H^1$-norm

$$\|\boldsymbol{s}_h\|_{X_h}^2 := \|[\sigma_h^\circ, u_h^\circ]\|_{\mathbb{H}_h}^2 = (\sigma_h^\circ, \sigma_h^\circ)_{H(\mathrm{div})} + (u_h^\circ, u_h^\circ)_{H^1}, \tag{60}$$

where each term is assembled using the same FE pair as in the high-fidelity discretization. The corresponding Gram matrix $X_h$ defines the discrete norm $\|\boldsymbol{s}_h\|_{X_h}^2 = \boldsymbol{s}_h^\top X_h \boldsymbol{s}_h$. We then perform Proper Orthogonal Decomposition (POD) in this $X_h$-inner product. Specifically, given $N_s$ snapshot solutions for a training sample set $\mathfrak{P}_{N_s} = \{\mathfrak{p}_1, \ldots, \mathfrak{p}_{N_s}\} \subset \mathfrak{P}$, stored in the matrix $\mathbf{S} = \mathbf{S}_{\mathfrak{P}_{N_s}} = (\boldsymbol{s}_h(\mathfrak{p}_1), \ldots, \boldsymbol{s}_h(\mathfrak{p}_{N_s})) \in \mathbb{R}^{N_h^s \times N_s}$, we solve the eigenvalue problem

$$\mathbf{C}\boldsymbol{v}_k = \lambda_k \boldsymbol{v}_k, \quad k = 1, \ldots, N_s, \tag{61}$$

where $\mathbf{C} := \mathbf{S}^\top X_h \mathbf{S}/N_s \in \mathbb{R}^{N_s \times N_s}$. The eigenpairs $(\lambda_k, \boldsymbol{v}_k)$, with $\lambda_1 \geq \lambda_2 \geq \cdots$, yield the POD basis vectors $\pi_k = \mathbf{S}\boldsymbol{v}_k/\sqrt{\lambda_k} \in \mathbb{R}^{N_h^s}$, $k = 1, \ldots, N_s$, which are orthonormal with respect to the $X_h$-inner product; see [26, Chapter 6].

The functions $\phi_k \in \mathbb{H}_h$, $k = 1, \ldots, r$, whose FE coefficients are given by the POD basis vectors $\pi_k$, thus form an $\mathbb{H}$-orthonormal basis of an $r$-dimensional subspace $\mathbb{H}_r = \mathrm{span}\{\phi_1, \ldots, \phi_r\}$, which we refer to as the POD-RB subspace of $\mathbb{H}_h$. Hence, the matrix $\Pi_r = (\pi_1, \ldots, \pi_r) \in \mathbb{R}^{N_h^s \times r}$ with columns $\pi_k$, defines for any $\boldsymbol{s}_h \in \mathbb{R}^{N_h^s}$ via $\boldsymbol{s}_r = \Pi_r^\top X_h \boldsymbol{s}_h \in \mathbb{R}^r$ a projection onto an $r$-dimensional subspace of $\mathbb{R}^{N_h^s}$. Specifically,

$$s_r = \Phi_r \boldsymbol{s}_r \in \mathbb{H}_r, \text{ with } \Phi_r = (\phi_1, \ldots, \phi_r) \text{ and } \boldsymbol{s}_r = \Pi_r^\top X_h \boldsymbol{s}_h \in \mathbb{R}^r,$$

is the $\mathbb{H}$-orthogonal projection of the FE function $s_h$ with coefficient vector $\boldsymbol{s}_h$. Conversely, given a coefficient vector $\boldsymbol{s}_r \in \mathbb{R}^r$ of a function $s_r$ in $\mathbb{H}_r$, the coefficient vector of $s_h$ is given by

$$\boldsymbol{s}_h = \Pi_r \boldsymbol{s}_r \in \mathbb{R}^{N_h^s}. \tag{62}$$

14

In practice, we take the projection dimension $r$ such that [26]

$$||s_h - s_r||^2_{L^2(\mathfrak{P};\mathbb{H}_h)} = ||\boldsymbol{s}_h - \Pi_r \boldsymbol{s}_r||^2_{L^2(\mathfrak{P};X_h)} \approx \sum_{k=r+1}^{N_s} \lambda_k \leq \tau \sum_{k=1}^{N_s} \lambda_k, \tag{63}$$

for a desired tolerance $\tau \ll 1$. Equivalently, $s_r$ approximates the "truth-space" solution $s_h$ in $\mathbb{H}_h$ with relative error $\sqrt{\tau}$. When the solution manifold has rapidly decaying Kolmogorov $r$-widths, we have $r \ll N_h^s$.

We will learn reduced basis neural operator (RBNO) *surrogates* $s_r(\mathfrak{p};\theta) = [\sigma_r^\circ, u_r^\circ](\mathfrak{p};\theta) \in \mathbb{H}_r$ to the parameter-to-solution map $\mathfrak{p} \mapsto s(\mathfrak{p}) = [\sigma^\circ(\mathfrak{p}), u^\circ(\mathfrak{p})] \in \mathbb{H}$ in terms of elements from the *hypothesis class*

$$\mathcal{H}_r(\Theta) := \left\{ s_r(\mathfrak{p};\theta) = \Phi_r \boldsymbol{s}_r(\mathfrak{p};\theta), \; \boldsymbol{s}_r(\mathfrak{p};\theta) \in \mathbb{R}^r, \mathfrak{p} \in \mathfrak{P}, \theta \in \Theta \right\} \subset L^2_\mu(\mathfrak{P}, \mathbb{H}_r), \tag{64}$$

where the expansion coefficient vector $\boldsymbol{s}_r(\mathfrak{p};\theta)$ is represented by neural networks with input variables $\mathfrak{p} \in \mathfrak{P}$ and trainable parameters $\theta \in \Theta$ in a suitable parameter space $\Theta$.

On account of FOSLS stability in (B.1) and (B.7), the loss $\mathcal{L}$ quantifies the error

$$\|s(\mathfrak{p}) - s_r\|^2_{\mathbb{H}} \eqsim \mathcal{L}(s_r; \mathfrak{p}) = \|\mathcal{B}_\mathfrak{p} s_r - S\|^2_{\mathbb{L}_2}$$

uniformly in $\mathfrak{p} \in \mathfrak{P}$, where $S$ is the source term, e.g., $S = [\mathfrak{p}\nabla w - z + f_1, f_2]$ for the diffusion problem.

In order to see how errors of this type depend on the two constituents $\mathbb{H}_h, \mathbb{H}_r$, we derive first corresponding decompositions of the loss. To that end, we make standard use of the fact that loss minimizers over any subspace of $\mathbb{H}$ are solutions of the normal equation posed on this subspace (with source data projected accordingly) and hence enjoy Galerkin orthogonality.

**Lemma 1.** *Assume that $\tilde{\mathbb{H}}$ is a closed subspace of $\mathbb{H}$ and that $\tilde{s}(\mathfrak{p}) = [\tilde{\sigma}^\circ(\mathfrak{p}), \tilde{u}^\circ(\mathfrak{p})] \in \tilde{\mathbb{H}}$ solves*

$$(\mathcal{B}_\mathfrak{p}\tilde{s}(\mathfrak{p}), \mathcal{B}_\mathfrak{p}\tilde{s}')_\Omega = (S, \mathcal{B}_\mathfrak{p}\tilde{s}')_\Omega, \quad \forall \tilde{s}' \in \tilde{\mathbb{H}}. \tag{65}$$

*Then we have*

$$\mathcal{L}(\tilde{s}'; \mathfrak{p}) = \|\mathcal{B}_\mathfrak{p}(\tilde{s}' - \tilde{s}(\mathfrak{p}))\|^2_{\mathbb{L}_2} + \|\mathcal{B}_\mathfrak{p}\tilde{s}(\mathfrak{p}) - S\|^2_{\mathbb{L}_2}, \quad \forall \tilde{s}' \in \tilde{\mathbb{H}}. \tag{66}$$

*In other words,*

$$\tilde{s}(\mathfrak{p}) = \operatorname*{argmin}_{\tilde{s}' \in \tilde{\mathbb{H}}} \mathcal{L}(\tilde{s}'; \mathfrak{p}) = \operatorname*{argmin}_{\tilde{s}' \in \tilde{\mathbb{H}}} \|\mathcal{B}_\mathfrak{p}(\tilde{s}' - \tilde{s}(\mathfrak{p}))\|^2_{\mathbb{L}_2}. \tag{67}$$

*Proof.* By definition, we have for the exact FOSLS solution $s(\mathfrak{p}) = [\sigma^\circ(\mathfrak{p}), u^\circ(\mathfrak{p})] \in \mathbb{H}$

$$\mathcal{L}(\tilde{s}'; \mathfrak{p}) = \|\mathcal{B}_\mathfrak{p}\tilde{s}' - S\|^2_{\mathbb{L}_2} = \|\mathcal{B}_\mathfrak{p}(\tilde{s}' - \tilde{s}(\mathfrak{p})) + \mathcal{B}_\mathfrak{p}(\tilde{s}(\mathfrak{p}) - s(\mathfrak{p}))\|^2_{\mathbb{L}_2}$$
$$= \|\mathcal{B}_\mathfrak{p}(\tilde{s}' - \tilde{s}(\mathfrak{p}))\|^2_{\mathbb{L}_2} + \|\mathcal{B}_\mathfrak{p}(\tilde{s}(\mathfrak{p}) - s(\mathfrak{p}))\|^2_{\mathbb{L}_2} + 2(\mathcal{B}_\mathfrak{p}(\tilde{s}' - \tilde{s}(\mathfrak{p})), \mathcal{B}_\mathfrak{p}(\tilde{s}(\mathfrak{p}) - s(\mathfrak{p})))_\Omega.$$

Since by (65), $(\mathcal{B}_\mathfrak{p}(\tilde{s}(\mathfrak{p}) - s(\mathfrak{p})), \mathcal{B}_\mathfrak{p}\tilde{s}')_\Omega = 0$ for all $\tilde{s}' \in \tilde{\mathbb{H}}$, and since $\tilde{s}' - \tilde{s}(\mathfrak{p}) \in \tilde{\mathbb{H}}$, we conclude that $(\mathcal{B}_\mathfrak{p}(\tilde{s}' - \tilde{s}(\mathfrak{p})), \mathcal{B}_\mathfrak{p}(\tilde{s}(\mathfrak{p}) - s(\mathfrak{p})))_\Omega = 0$ for all $\tilde{s}' \in \tilde{\mathbb{H}}$, whence the assertion follows. $\square$

We are now in a position to interrelate best approximation errors for later purposes.

**Theorem 2.** *As before, let $\mathbb{H}_h = \Sigma_h \times \mathbb{U}_h \subset H(\operatorname{div};\Omega) \times H^1(\Omega)$ be the $H(\operatorname{div}) \times H^1$-conforming FE space and, for each fixed parameter $\mathfrak{p} \in \mathfrak{P}$, let $s_h(\mathfrak{p}) = [\sigma_h^\circ(\mathfrak{p}), u_h^\circ(\mathfrak{p})] \in \mathbb{H}_h$ be the FE FOSLS solution and*

$$s_r(\mathfrak{p}) = [\sigma_r^\circ(\mathfrak{p}), u_r^\circ(\mathfrak{p})] = \operatorname*{argmin}_{[\tau_r, v_r] \in \mathbb{H}_r} \mathcal{L}([\tau_r, v_r]; \mathfrak{p}) \tag{68}$$

*be the RB FOSLS solution in the POD subspace $\mathbb{H}_r \subset \mathbb{H}_h$. Then the following hold:*

1. *For any $s_r \in \mathbb{H}_r$, there holds*

$$\mathcal{L}(s_r; \mathfrak{p}) \eqsim \| s_r(\mathfrak{p}) - s_r \|^2_{\mathbb{H}} + \mathcal{L}(s_r(\mathfrak{p}); \mathfrak{p}). \tag{69}$$

*In addition, the reduced FOSLS minimizer $s_r(\mathfrak{p})$ in (68) satisfies the quasi-optimality bound*

$$\| s_r(\mathfrak{p}) - s(\mathfrak{p}) \|_{\mathbb{H}} \eqsim \min_{t_r \in \mathbb{H}_r} \| t_r - s(\mathfrak{p}) \|_{\mathbb{H}}. \tag{70}$$

15

2. *For any $s_r \in \mathbb{H}_r$, there holds*

$$\mathcal{L}(s_r; \mathfrak{p}) \eqsim \| s_h(\mathfrak{p}) - s_r \|_{\mathbb{H}}^2 + \mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p}). \tag{71}$$

*In addition, the reduced FOSLS minimizer $s_r(\mathfrak{p})$ in (68) satisfies the quasi-optimality bound*

$$\| s_r(\mathfrak{p}) - s_h(\mathfrak{p}) \|_{\mathbb{H}} \eqsim \min_{t_r \in \mathbb{H}_r} \| t_r - s_h(\mathfrak{p}) \|_{\mathbb{H}}. \tag{72}$$

3. *The function $s_r \in \mathbb{X}_r := L_\mu^2(\mathfrak{P}; \mathbb{H}_r)$ that for each $\mathfrak{p} \in \mathfrak{P}$ solves (65) for $\tilde{\mathbb{H}} = \mathbb{H}_r$, is the unique minimizer*

$$s_r = \operatorname*{argmin}_{t \in \mathbb{X}_r} \mathbb{E}_{\mathfrak{p} \sim \mu} \big[ \mathcal{L}(t(\mathfrak{p}); \mathfrak{p}) \big] =: \operatorname*{argmin}_{t \in \mathbb{X}_r} \mathcal{L}(t; \mathfrak{P}) \tag{73}$$

*and*

$$\mathbb{E}_{\mathfrak{p} \sim \mu} \big[ \mathcal{L}(s_r(\mathfrak{p}); \mathfrak{p}) \big] \eqsim \mathbb{E}_{\mathfrak{p} \sim \mu} \big[ \| s_r(\mathfrak{p}) - s(\mathfrak{p}) \|_{\mathbb{H}}^2 \big] \eqsim \min_{t_r \in \mathbb{X}_r} \mathbb{E}_{\mathfrak{p} \sim \mu} \big[ \| t_r(\mathfrak{p}) - s(\mathfrak{p}) \|_{\mathbb{H}}^2 \big], \tag{74}$$

*with a constant that depends only on $\mathfrak{P}$.*

*Proof.* 1. Taking $\tilde{\mathbb{H}} := \mathbb{H}_r$, (69) follows from (66) in Lemma 1 with $\tilde{s}(\mathfrak{p}) = s_r(\mathfrak{p})$ and $\tilde{s}' = s_r$, combined with FOSLS stability (B.1) and (B.7) for the diffusion and elasticity model, respectively. We obtain (70) by the optimality of $s_r(\mathfrak{p})$ and FOSLS stability, i.e., for any $t_r \in \mathbb{H}_r$,

$$\| s_r(\mathfrak{p}) - s(\mathfrak{p}) \|_{\mathbb{H}} \eqsim \| \mathcal{B}_\mathfrak{p} s_r(\mathfrak{p}) - S \|_{\mathbb{L}_2}^2 = \mathcal{L}(s_r(\mathfrak{p}); \mathfrak{p}) \leq \mathcal{L}(t_r; \mathfrak{p}) = \| \mathcal{B}_\mathfrak{p} t_r - S \|_{\mathbb{L}_2}^2 \eqsim \| t_r - s(\mathfrak{p}) \|_{\mathbb{H}}. \tag{75}$$

2. The same reasoning for $\tilde{\mathbb{H}} = \mathbb{H}_h$, $\tilde{s}(\mathfrak{p}) = [\sigma_h^\circ(\mathfrak{p}), u_h^\circ(\mathfrak{p})]$, confirms (71), noting that $s_r \in \mathbb{H}_r$ are also elements in $\mathbb{H}_h$. Applying minimization of $s_r \in \mathbb{H}_r$ over (71), we have

$$\min_{s_r \in \mathbb{H}_r} \mathcal{L}(s_r; \mathfrak{p}) \simeq \min_{s_r \in \mathbb{H}_r} \| s_h(\mathfrak{p}) - s_r \|_{\mathbb{H}}^2 + \mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p}), \tag{76}$$

or equivalently,

$$\mathcal{L}(s_r(\mathfrak{p}) \; \mathfrak{p}) \simeq \min_{t_r \in \mathbb{H}_r} \| t_r - s_h(\mathfrak{p}) \|_{\mathbb{H}}^2 + \mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p}). \tag{77}$$

Then (72) holds when we combing (77) with $\mathcal{L}(s_r(\mathfrak{p}); \mathfrak{p}) \eqsim \| s_h(\mathfrak{p}) - s_r(\mathfrak{p}) \|_{\mathbb{H}}^2 + \mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p})$ which is obtained by plugging $s_r = s_r(\mathfrak{p})$ in (71).

3. This is an immediate consequence of 1. and uniform FOSLS stability (B.1) and (B.7). $\qquad \square$

**Remark 1.** *Theorem 2 says that $s_r(\mathfrak{p})$ minimizes for each $\mathfrak{p}$ the loss over $\mathbb{H}_r$ and (near-)best approximates $s(\mathfrak{p})$. Statement 2 says that $s_r(\mathfrak{p})$ (near-)best approximates the FE-FOSLS solution from $\mathbb{H}_r$. Overall this states that the error incurred by $s_r(\mathfrak{p})$ with respect to the exact solution is uniformly proportional to the sum of the best approximation errors of $s(\mathfrak{p})$ from the FE space $\mathbb{H}_h$ and the error of the approximation of the FE-FOSLS solution by the RB-FOSLS solution.*

**Corollary 1.** *Adhering to the above assumptions, let*

$$s_r(\theta^*) \in \operatorname*{argmin}_{t_r \in \mathcal{H}_r(\Theta)} \mathbb{E}_{\mathfrak{p} \sim \mu} \big[ \mathcal{L}(t_r(\mathfrak{p}); \mathfrak{p}) \big]. \tag{78}$$

*Then, with a proportionality constant, depending only on uniform FOSLS-stability (B.1) and (B.7), there holds*

$$\mathbb{E}_{\mathfrak{p} \sim \mu} \big[ \| s(\mathfrak{p}) - s_r(\mathfrak{p}; \theta^*) \|_{\mathbb{H}}^2 \big] \lesssim h^{2\eta} + \sum_{j > r} \lambda_j^2 + \mathbb{E}_{\mathfrak{p} \sim \mu} \big[ \| s_r(\mathfrak{p}) - s_r(\mathfrak{p}; \theta^*) \|_{\mathbb{H}}^2 \big], \tag{79}$$

*where $\eta$ depends on the polynomial orders of the FE approximation and the regularity of $s(\mathfrak{p}) \in \mathbb{H}$, $\mathfrak{p} \in \mathfrak{P}$, as in Theorem 1. Note that the last summand reflects the best approximation from the hypothesis class to the best approximation of the exact FOSLS solutions $s(\mathfrak{p})$ from $\mathbb{H}_r$.*

*Proof.* We infer from Theorem 2, (69), that, with $s_r(\theta^*) \in \mathcal{H}_r(\Theta)$, and FOSLS stability (B.1) and (B.7),

$$
\begin{aligned}
&\mathbb{E}_{\mathfrak{p} \sim \mu}\big[\|s(\mathfrak{p}) - s_r(\mathfrak{p}; \theta^*)\|_{\mathbb{H}}^2\big] \\
&\eqsim \mathbb{E}_{\mathfrak{p} \sim \mu}\Big[\|s_r(\mathfrak{p}; \theta^*) - s_r(\mathfrak{p})\|_{\mathbb{H}}^2\Big] + \mathbb{E}_{\mathfrak{p} \sim \mu}\Big[\mathcal{L}(s_r(\mathfrak{p}); \mathfrak{p})\Big] \\
&\eqsim \mathbb{E}_{\mathfrak{p} \sim \mu}\Big[\|s_r(\mathfrak{p}; \theta^*) - s_r(\mathfrak{p})\|_{\mathbb{H}}^2\Big] + \mathbb{E}_{\mathfrak{p} \sim \mu}\Big[\|s_h(\mathfrak{p}) - s_r(\mathfrak{p})\|_{\mathbb{H}}^2\Big] + \mathbb{E}_{\mathfrak{p} \sim \mu}\Big[\mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p})\Big],
\end{aligned}
\tag{80}
$$

where we have invoked Theorem 2, Statement 2 in the last step. Bounding the first term in the last inequality by (63), using again that, uniformly in $\mathfrak{p} \in \mathfrak{P}$, on account of FOSLS-stability and the best approximation property of the FE-FOSLS solutions, we conclude that

$$
\mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p}) \eqsim \inf_{t_h(\mathfrak{p}) \in \mathbb{H}_h} \|t_h(\mathfrak{p}) - s(\mathfrak{p})\|_{\mathbb{H}}^2.
$$

Finally, invoking the error bounds in Appendix C.1 and Appendix C.2, finishes the proof. $\square$

## 5.2 Loss evaluation

We learn an approximation to the optimal approximation $s_r(\mathfrak{p}) \in \mathbb{H}_r$ from the hypothesis class $\mathcal{H}_r(\Theta) \subset \mathbb{X}_r$ (see (64)) by minimizing an empirical loss over $\Theta$. Since this will require frequent evaluations of fiber losses $\mathcal{L}(\cdot; \mathfrak{p})$ at randomly drawn samples $\mathfrak{p} \in \mathfrak{P}$, the cost of such evaluations is an issue. In fact, substituting the FE approximation (51) in the residual fiber loss function (30), we obtain

$$
\mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p}) = \boldsymbol{w}(\mathfrak{p})^\top W_\mathfrak{p} \, \boldsymbol{w}(\mathfrak{p}) = \boldsymbol{s}_h(\mathfrak{p})^\top W_\mathfrak{p}^\circ \boldsymbol{s}_h(\mathfrak{p}) + 2\boldsymbol{s}_h(\mathfrak{p})^\top \boldsymbol{\alpha}_\mathfrak{p} + \beta_\mathfrak{p},
\tag{81}
$$

where $\boldsymbol{w}(\mathfrak{p}) = (\boldsymbol{s}_h(\mathfrak{p}); 1) \in \mathbb{R}^{N_h^s+1}$ is a concatenation of the parametric coefficient vectors $\boldsymbol{s}_h(\mathfrak{p}) \in \mathbb{R}^{N_h^s}$ and scalar 1. Here $W_\mathfrak{p} \in \mathbb{R}^{(N_h^s+1)\times(N_h^s+1)}$ is a sparse symmetric matrix assembled from the FE approximation at each parameter $\mathfrak{p} \in \mathfrak{P}$, which has the block structure

$$
W_\mathfrak{p} = \begin{pmatrix} W_\mathfrak{p}^\circ & \boldsymbol{\alpha}_\mathfrak{p} \\ \boldsymbol{\alpha}_\mathfrak{p}^\top & \beta_\mathfrak{p} \end{pmatrix},
$$

with $W_\mathfrak{p}^\circ \in \mathbb{R}^{N_h^s \times N_h^s}$, $\boldsymbol{\alpha}_\mathfrak{p} \in \mathbb{R}^{N_h^s}$, and $\beta_\mathfrak{p} \in \mathbb{R}$ assembled from the bilinear, linear, and constant forms in the loss function (30). Note that evaluation of the loss function (81) has complexity $O(N_h^s)$ when employing sparse matrix vector products, which is expensive for large-scale FE discretizations.

To reduce the computational complexity, we introduce a RB computation of the loss function (81), when evaluated at elements in the RB space $\mathbb{H}_r$. Let $s_r$ denote the RB approximation of the FE solution $s_h$, with corresponding approximation of the coefficient vectors by the projection in (62). Then the RB loss function can be evaluated as

$$
\mathcal{L}(s_r; \mathfrak{p}) = \boldsymbol{s}_r(\mathfrak{p})^\top W_\mathfrak{p}^r \boldsymbol{s}_r(\mathfrak{p}) + 2\boldsymbol{s}_r(\mathfrak{p})^\top \boldsymbol{\alpha}_\mathfrak{p}^r + \beta_\mathfrak{p},
\tag{82}
$$

where reduced weights $W_\mathfrak{p}^r := \Pi_r^\top W_\mathfrak{p}^\circ \Pi_r \in \mathbb{R}^{r \times r}$ and vectors $\boldsymbol{\alpha}_\mathfrak{p}^r := \Pi_r^\top \boldsymbol{\alpha}_\mathfrak{p} \in \mathbb{R}^r$ can be precomputed and stored for the training samples of the parameter, allowing a fast evaluation of the loss function with complexity $O(r^2)$. The RB FOSLS minimizer $s_r(\mathfrak{p}) \in \mathbb{H}_r$ from (68) can then be computed with the RB coefficient vector $\boldsymbol{s}_r(\mathfrak{p}) \in \mathbb{R}^r$ obtained by solving the reduced normal equation

$$
W_\mathfrak{p}^r \, \boldsymbol{s}_r(\mathfrak{p}) = -\boldsymbol{\alpha}_\mathfrak{p}^r.
\tag{83}
$$

For the elasticity problem, we follow the same computation of the RB approximation as in the diffusion case. Specifically, we adopt the following discrete $H(\mathrm{div}; \Omega; \mathbb{R}^{d \times d}) \times H^1(\Omega; \mathbb{R}^d)$-norm as in [20]

$$
\|\boldsymbol{s}_h\|_{X_h}^2 := (\underline{\mathrm{div}}\, \underline{\underline{\sigma}}_h^\circ, \underline{\mathrm{div}}\, \underline{\underline{\sigma}}_h^\circ)_{L^2} + (\underline{\underline{\sigma}}_h^\circ, \mathcal{C}_{\bar{\mathfrak{p}}}^{-1}\underline{\underline{\sigma}}_h^\circ)_{L_2} + (\underline{\underline{\varepsilon}}(\underline{u}_h^\circ), \mathcal{C}_{\bar{\mathfrak{p}}}\underline{\underline{\varepsilon}}(\underline{u}_h^\circ))_{L^2},
\tag{84}
$$

where $\mathcal{C}_{\bar{\mathfrak{p}}}$ and $\mathcal{C}_{\bar{\mathfrak{p}}}^{-1}$ are the stiffness tensor and its inverse at the mean of the parameter $\bar{\mathfrak{p}} \in \mathfrak{P}$. The same statements as in Proposition 2 hold for the RB approximation of the elasticity problem with the loss function (48) and the discrete norm (84).

We conclude this section with one further prerequisite that will be needed in the next section.

**Proposition 1.** *Assume $\mathfrak{p} \in \mathfrak{P} \subset L_\infty(\Omega)$ is uniformly bounded as in Theorem 1. Then for all $\mathfrak{p} \in \mathfrak{P}$, the reduced weights $W_{\mathfrak{p}}^r$, vectors $\boldsymbol{\alpha}_{\mathfrak{p}}^r$, and scalars $\beta_{\mathfrak{p}}$, defined in (82), satisfy the uniform bounds*

$$\gamma_- I \preceq W_{\mathfrak{p}}^r \preceq \gamma_+ I, \qquad \|\boldsymbol{\alpha}_{\mathfrak{p}}^r\|_2 \le \sqrt{\gamma_+}\, S_{\max}, \qquad |\beta_{\mathfrak{p}}| \le S_{\max}^2.$$

*Here $c \le \gamma_- \le \gamma_+ \le C$ with the stability constants $0 < c < C < \infty$ in Lemma 3, and $S_{\max}$ is the uniform bound for the input source $\|S(\mathfrak{p})\|_{\mathbb{L}^2} \le S_{\max}$, $\mathfrak{p} \in \mathfrak{P}$. For instance, for diffusion one has $S = [\mathfrak{p}\nabla w - z + f_1, f_2]$.*

*Proof.* Since $W_{\mathfrak{p}}^r = \Pi_r^\top W_{\mathfrak{p}}^\circ \Pi_r$, for any $\boldsymbol{s}_r \in \mathbb{R}^r$ we have $\boldsymbol{s}_r^\top W_{\mathfrak{p}}^r \boldsymbol{s}_r = (\Pi_r \boldsymbol{s}_r)^\top W_{\mathfrak{p}}^\circ (\Pi_r \boldsymbol{s}_r)$. Uniform FOSLS stability (see Lemma 3) implies there exist $0 < \gamma_- \le \gamma_+ < \infty$ such that the stability estimates (B.1) hold in FE spaces $\Sigma_h \times \mathbb{U}_h \subset \Sigma \times \mathbb{U}$ with constants $c \le \gamma_-$ and $\gamma_+ \le C$, i.e., for all FE coefficient vectors $\boldsymbol{s}_h \in \mathbb{R}^{N_h^s}$,

$$\gamma_-\, \boldsymbol{s}_h^\top X_h \boldsymbol{s}_h \ \le \ \boldsymbol{s}_h^\top W_{\mathfrak{p}}^\circ \boldsymbol{s}_h \ \le \ \gamma_+\, \boldsymbol{s}_h^\top X_h \boldsymbol{s}_h, \qquad \mathfrak{p} \in \mathfrak{P}.$$

Therefore, by setting $\boldsymbol{s}_h = \Pi_r \boldsymbol{s}_r$, there holds

$$\gamma_-\, \boldsymbol{s}_r^\top (\Pi_r^\top X_h \Pi_r) \boldsymbol{s}_r \ \le \ \boldsymbol{s}_r^\top W_{\mathfrak{p}}^r \boldsymbol{s}_r \ \le \ \gamma_+\, \boldsymbol{s}_r^\top (\Pi_r^\top X_h \Pi_r) \boldsymbol{s}_r.$$

By the $X_h$-orthonormality of $\Pi_r$, $\Pi_r^\top X_h \Pi_r = I_r$, hence $\gamma_- \|\boldsymbol{s}_r\|_2^2 \le \boldsymbol{s}_r^\top W_{\mathfrak{p}}^r \boldsymbol{s}_r \le \gamma_+ \|\boldsymbol{s}_r\|_2^2$, proving the spectral bounds. Next, from the quadratic expansion $\mathcal{L}(\boldsymbol{s}_h; \mathfrak{p}) = \boldsymbol{s}_h^\top W_{\mathfrak{p}}^\circ \boldsymbol{s}_h + 2\, \boldsymbol{s}_h^\top \boldsymbol{\alpha}_{\mathfrak{p}} + \beta_{\mathfrak{p}}$, the linear term obeys, for any FE coefficient vector $\boldsymbol{s}_h$ with corresponding FE functions $s_h$,

$$|\boldsymbol{s}_h^\top \boldsymbol{\alpha}_{\mathfrak{p}}| \ = \ (\mathcal{B}_{\mathfrak{p}} s_h, S(\mathfrak{p})) \ \le \ \|\mathcal{B}_{\mathfrak{p}} s_h\|_{\mathbb{L}^2}\, \|S(\mathfrak{p})\|_{\mathbb{L}^2} \ \le \ \sqrt{\gamma_+}\, \|S(\mathfrak{p})\|_{\mathbb{L}^2}\, (\boldsymbol{s}_h^\top X_h \boldsymbol{s}_h)^{1/2},$$

where we have used Cauchy–Schwarz in the first inequality and the uniform stability bound in the second. Taking $\boldsymbol{s}_h = \Pi_r \boldsymbol{s}_r$ with $\|\boldsymbol{s}_r\|_2 = 1$ and using $\boldsymbol{s}_r^\top (\Pi_r^\top X_h \Pi_r) \boldsymbol{s}_r = 1$, yields $\|\boldsymbol{\alpha}_{\mathfrak{p}}^r\|_2 = \sup_{\|\boldsymbol{s}_r\|_2 = 1} (\Pi_r \boldsymbol{s}_r)^\top \boldsymbol{\alpha}_{\mathfrak{p}} \le \sqrt{\gamma_+}\, \|S(\mathfrak{p})\|_{\mathbb{L}^2} \le \sqrt{\gamma_+}\, S_{\max}$. Finally, $\beta_{\mathfrak{p}} = \|S(\mathfrak{p})\|_{\mathbb{L}^2}^2 \le S_{\max}^2$. Uniformity follows from the hypotheses. $\qquad \square$

## 5.3 Learning the reduced basis neural operator

Recall that the elements of the hypothesis class $\mathcal{H}_r(\Theta)$, defined in (64), belong for each $\mathfrak{p} \in \mathfrak{P}$, as functions of $x \in \Omega$, to the $r$–dimensional RB space constructed in Section 5.1. As indicated in Section 2, learning a parametric map will be based on minimizing an empirical risk of the type

$$\min_{\theta \in \Theta} \hat{R}_{\widehat{\mathfrak{P}}}(\theta), \quad \text{where} \quad \hat{R}_{\widehat{\mathfrak{P}}}(\theta) := \hat{R}_{\widehat{\mathfrak{P}}}(s_r(\cdot; \theta)) := \frac{1}{\#\widehat{\mathfrak{P}}} \sum_{\mathfrak{p} \in \widehat{\mathfrak{P}}} \mathcal{L}(s_r(\mathfrak{p}; \theta); \mathfrak{p}). \tag{85}$$

In what follows

$$\widehat{\mathfrak{P}} = \{\mathfrak{p}_1, \ldots, \mathfrak{p}_N\} \subset \mathfrak{P},$$

stands for a collection of finite i.i.d. random samples from $\mathfrak{P}$. In computations, the evaluation of the RB loss $\mathcal{L}(s_r(\mathfrak{p}; \theta); \mathfrak{p})$ is always based on the right hand side expression (82) for $\boldsymbol{s}_r(\mathfrak{p}) = \boldsymbol{s}_r(\mathfrak{p}; \theta)$. The coefficients in $\boldsymbol{s}_r(\mathfrak{p}; \theta)$ need to approximate for every $\mathfrak{p} \in \mathfrak{P}$ the coefficient vector $\boldsymbol{s}_r(\mathfrak{p})$ of function $s_r(\mathfrak{p})$. Recall that those are obtained when minimizing the "ideal continuous" loss which we also abbreviate, for convenience, as follows: for any $\theta \in \Theta$ we write $s_r(\cdot; \theta) = [\sigma_r^\circ(\cdot; \theta), u_r^\circ(\cdot; \theta)] \in \mathcal{H}_r(\Theta)$ and set

$$R(s_r(\cdot; \theta)) =: R(\theta) := \mathbb{E}_{\mathfrak{p} \sim \mu}\Big[\mathcal{L}(s_r(\mathfrak{p}; \theta); \mathfrak{p})\Big].$$

We refer to functions $s_r(\cdot; \theta) = [\sigma_r^\circ(\cdot; \theta), u_r^\circ(\cdot; \theta)] \in \mathcal{H}_r(\Theta)$ as RBNO approximations to the solution of the RB-FOSLS problem.

In order to estimate the accuracy of empirical loss minimizers with respect to the Bochner norm $\|w\|_{\mathbb{X}}^2 = \int_{\mathfrak{P}} \|w(\mathfrak{p})\|_{\mathbb{H}}^2 d\mu(\mathfrak{p}) = \mathbb{E}_{\mathfrak{p} \sim \mu}\big[\|w(\mathfrak{p})\|_{\mathbb{H}}^2\big]$, we combine the previous error bounds with (rather standard) concepts from Learning Theory. Here we are content with a relatively simple version that may not be the best possible. Nevertheless, their applicability imposes some conditions on the hypothesis class $\mathcal{H}_r(\Theta)$ that need to be ensured.

*Boundedness:* Since $s_r(\mathfrak{p})$ minimizes the (ideal) loss over $\mathbb{H}_r$ and hence solves the normal equations in $\mathbb{H}_r$, uniform ellipticity of the normal equations implies boundedness of the solution by the data, which in our case means, on account of the orthogonality of the reduced POD basis $\Phi_r$, that

$$\|\boldsymbol{s}_r(\mathfrak{p})\|_{\ell_2} = \|[\sigma_r^\circ(\mathfrak{p}), u_r^\circ(\mathfrak{p})]\|_{\mathbb{H}} \le \sqrt{\gamma_+} S_{\max}, \quad \mathfrak{p} \in \mathfrak{P},$$

provided that $\mathfrak{P}$ remains bounded in a suitable domain. Therefore, restraining the neural network approximations to the $\mathfrak{p}$-dependent coefficients to remain bounded as well will not impede their expressivity.

**Assumption 1.** *There exists a uniform bound $B > 0$ such that the neural network output $\|\boldsymbol{s}_r(\mathfrak{p}; \theta)\|_{\ell_2} \le B$ for all $(\mathfrak{p}, \theta) \in \mathfrak{P} \times \Theta$ (enforced, e.g., by range clipping or weight penalties).*

The second important property is the uniform bound and Lipschitz continuity of the loss $\mathcal{L}(t_r; \mathfrak{p})$ with respect to the first argument $t_r \in \mathcal{H}_r(\Theta)$.

**Lemma 2.** *Assume that Assumption 1 is valid. Then, there exist constants $M, L < \infty$ such that*

$$\mathcal{L}(t_r; \mathfrak{p}) \le M, \quad (t_r, \mathfrak{p}) \in \mathcal{H}_r(\Theta) \times \mathfrak{P}, \tag{86}$$

*and*

$$\left|\mathcal{L}(t_r; \mathfrak{p}) - \mathcal{L}(t_r'; \mathfrak{p})\right| \le L\|t_r - t_r'\|_{\mathbb{H}}, \quad t_r, t_r' \in \mathcal{H}_r(\Theta), \mathfrak{p} \in \mathfrak{P}. \tag{87}$$

*Proof.* By definition,

$$\mathcal{L}(t_r; \mathfrak{p})^{\frac{1}{2}} = \|\mathcal{B}_{\mathfrak{p}} t_r - S\|_{\mathbb{L}_2} \le \|\mathcal{B}_{\mathfrak{p}} t_r\|_{\mathbb{L}_2} + S_{\max} \le \sqrt{\gamma_+}\|t_r\|_{\mathbb{H}} + S_{\max}.$$

Since $t_r \in \mathcal{H}_r(\Theta)$ is of the form $t_r(x; \mathfrak{p}) = \Phi_r \boldsymbol{t}_r(\mathfrak{p}; \theta)$, orthogonality of the reduced bases $\Phi_r$, gives $\|t_r\|_{\mathbb{H}} = \|\boldsymbol{t}_r(\mathfrak{p}; \theta)\|_{\ell_2} \le B$, $(\mathfrak{p}, \theta) \in \mathfrak{P} \times \Theta$, by Assumption 1, confirming (86) with

$$M \le \left(\sqrt{\gamma_+} B + S_{\max}\right)^2. \tag{88}$$

Similarly,

$$\begin{aligned}
\left|\mathcal{L}(t_r; \mathfrak{p}) - \mathcal{L}(t_r'; \mathfrak{p})\right| &= \left|\|\mathcal{B}_{\mathfrak{p}} t_r - S\|_{\mathbb{L}_2}^2 - \|\mathcal{B}_{\mathfrak{p}} t_r' - S\|_{\mathbb{L}_2}^2\right| \\
&= \left|\|\mathcal{B}_{\mathfrak{p}} t_r - S\|_{\mathbb{L}_2} + \|\mathcal{B}_{\mathfrak{p}} t_r' - S\|_{\mathbb{L}_2}\right| \left|\|\mathcal{B}_{\mathfrak{p}} t_r - S\|_{\mathbb{L}_2} - \|\mathcal{B}_{\mathfrak{p}} t_r' - S\|_{\mathbb{L}_2}\right| \\
&\le 2\left\{\sqrt{\gamma_+} B + S_{\max}\right\}\|\mathcal{B}_{\mathfrak{p}} t_r - \mathcal{B}_{\mathfrak{p}} t_r'\|_{\mathbb{L}_2} \\
&\le 4\left\{\sqrt{\gamma_+} B + S_{\max}\right\}\sqrt{\gamma_+}\|t_r - t_r'\|_{\mathbb{H}},
\end{aligned}$$

which finishes the proof with $L = 4\left\{\sqrt{\gamma_+} B + S_{\max}\right\}\sqrt{\gamma_+}$. $\qquad\square$

As a final prerequisite, we recall the notion of *pseudo-dimension* of a function class $\mathcal{F}$. The pseudo-dimension is the VC-dimension of the set of their epi-graphs. For precise definitions, see, for instance, [43, Chapter 3].

**Remark 2.** *When $\mathcal{H}$ is a linear space, one has $\mathrm{Pdim}\,\mathcal{H} = \dim \mathcal{H} + 1$. When $\mathcal{H}$ is a fully connected neural network of depth $D$, it has been shown in [44] that $\mathrm{Pdim}\,\mathcal{H} \lesssim \#\Theta\, D \log(\#\Theta)$.*

To highlight the roles of the involved constituents in our particular setting, we present the following uniform convergence result in high probability.

**Theorem 3.** *We assume the validity of Assumption 1. For any $\mathfrak{p} \in \mathfrak{P}$, let $s(\mathfrak{p}) = [\sigma^{\circ}(\mathfrak{p}), u^{\circ}(\mathfrak{p})] \in \mathbb{H}$ be the solution of the FOSLS problem (26) and let $s_r(\mathfrak{p}, \hat{\theta}) = [\sigma_r^{\circ}, u_r^{\circ}](\mathfrak{p}; \hat{\theta}) \in \mathbb{H}_r$ denote the RBNO approximation, obtained by minimizing the empirical risk (85) over $\Theta$. More precisely, $s_r(\mathfrak{p}, \hat{\theta})$ satisfies*

$$\hat{R}_{\widehat{\mathfrak{P}}}(\hat{\theta}) \leq \inf_{\theta \in \Theta} \hat{R}_{\widehat{\mathfrak{P}}}(\theta) + \varepsilon_{\mathrm{opt}}, \tag{89}$$

*where $\varepsilon_{\mathrm{opt}} \geq 0$ accounts for a possible optimization error. Then, for any $\delta \in (0,1)$, with probability at least $1 - \delta$, there holds*

$$\mathbb{E}_{\mathfrak{p} \sim \mu}\Big[ \|s(\mathfrak{p}) - s_r(\mathfrak{p}; \hat{\theta})\|_{\mathbb{H}}^2 \Big] \leq c^{-1} \inf_{\theta \in \Theta} \mathbb{E}_{\mathfrak{p} \sim \mu}\big[\mathcal{L}(s_r(\mathfrak{p}; \theta); \mathfrak{p})\big] + M\sqrt{\frac{\log \frac{1}{\delta}}{2N}} + CL\sqrt{\frac{2P \log \frac{eN}{P}}{N}} + \varepsilon_{\mathrm{opt}}, \tag{90}$$

*where $c$ is the constant from (B.1) or (B.7), $C$ is an absolute constant, $P$ is the pseudo-dimension of the hypothesis class $\mathcal{H}_r(\Theta)$, and $M, L$ are from Lemma 2. Moreover, there exist constants $C_1$ and $C_2$, depending on FE discretization and uniform FOSLS stability, respectively, such that the approximation error can be decomposed as follows*

$$\inf_{\theta \in \Theta} \mathbb{E}_{\mathfrak{p} \sim \mu}\big[\mathcal{L}(s_r(\mathfrak{p}; \theta); \mathfrak{p})\big] \leq C_1\Big\{h^{2\eta} + \sum_{k > r} \lambda_k\Big\} + C_2 \inf_{\theta \in \Theta} \mathbb{E}_{\mathfrak{p} \sim \mu}\Big[\|s_r(\mathfrak{p}) - s_r(\mathfrak{p}; \theta)\|_{\mathbb{H}}^2\Big]. \tag{91}$$

*Here the exponent $\eta$ depends on the regularity of the exact FOSLS solution $s(\mathfrak{p})$, see Appendix C.*

Before turning to the proof, a few comments are in order. (91) says that the first term on the right hand side of (90) could be reduced to an approximation error within the subspace $\mathbb{H}_r$ by choosing a sufficiently fine FE discretization and a correspondingly accurate reduced basis. The second group of terms in (90), representing the estimation error, is only meaningful for $N > P$, hence in an *under-parametrized* regime, which is in accordance with the fact that resulting surrogates are supposed to act as reduced models. For large $N$, the estimation error behaves like $O(N^{-1/2})$, which is a slow rate. Faster rates would also require variance information, which is not required here.

Bounds of the above type are, in essence, standard. For the convenience of the reader we devote the remainder of this section to highlighting the main ingredients of the proof.

*Proof.* Note first that in the present case, the minimal risk

$$R^* = \min_{s \in L_{\mu}^2(\mathfrak{P}; \mathbb{H})} \mathbb{E}_{\mathfrak{p} \sim \mu}\big[\mathcal{L}(s; \mathfrak{p})\big] = 0.$$

Thus, by uniform FOSLS stability

$$\mathbb{E}_{\mu}\Big[\|s(\mathfrak{p}) - s_r(\mathfrak{p}; \hat{\theta})\|_{\mathbb{H}}^2\Big] \eqsim \mathbb{E}_{\mathfrak{p} \sim \mathfrak{P}}\big[\mathcal{L}(s_r(\mathfrak{p}; \hat{\theta}); \mathfrak{p})\big] = R(\hat{\theta}) = R(\hat{\theta}) - R^*,$$

so that, in the above terms, our task is to bound $R(\hat{\theta}) = R(s_r(\cdot; \hat{\theta}))$. We follow the first standard lines and define

$$s_r^* \in \operatorname*{argmin}_{s_r \in \mathcal{H}_r(\Theta)} R(s_r),$$

and decompose the excess risk

$$R(\hat{\theta}) = R(\hat{\theta}) - R(s_r^*) + R(s_r^*). \tag{92}$$

By FOSLS-stability,

$$R(s_r^*) \eqsim \min_{s_r \in \mathcal{H}_r(\Theta)} \mathbb{E}_{\mathfrak{p} \sim \mu}\big[\|s(\mathfrak{p}) - s_r(\mathfrak{p})\|_{\mathbb{H}}^2\big]$$

20

represents the *best-approximation error* from the hypothesis class $\mathcal{H}_r(\Theta)$, which is a deterministic quantity. Instead, the first summand, $R(\hat{\theta}) - R(s_r^*)$, is usually referred to as *estimation error*, a stochastic quantity, which we further have to analyze. We employ a further decomposition to obtain

$$R(\hat{\theta}) - R(s_r^*) = \{R(\hat{\theta}) - \hat{R}_{\widehat{\mathfrak{P}}}(\hat{\theta})\} + \{\hat{R}_{\widehat{\mathfrak{P}}}(\hat{\theta}) - \hat{R}_{\widehat{\mathfrak{P}}}(s_r^*)\} + \{\hat{R}_{\widehat{\mathfrak{P}}}(s_r^*) - R(s_r^*)\}$$
$$\leq 2 \sup_{s_r \in \mathcal{H}_r(\Theta)} \{R(s_r) - \hat{R}_{\widehat{\mathfrak{P}}}(s_r)\} + \varepsilon_{\text{opt}},$$

where we have used (89). This gives

$$R(\hat{\theta}) \leq 2 \sup_{s_r \in \mathcal{H}_r(\Theta)} \{R(s_r) - \hat{R}_{\widehat{\mathfrak{P}}}(s_r)\} + R(s_r^*) + \varepsilon_{\text{opt}}, \tag{93}$$

To further estimate the only stochastic quantity on the right hand side of (93), note that the estimation error

$$E(\mathfrak{p}_1, \ldots, \mathfrak{p}_N) := \sup_{s_r \in \mathcal{H}_r(\Theta)} \{R(s_r) - \hat{R}_{\widehat{\mathfrak{P}}}(s_r)\}, \tag{94}$$

satisfies, on account of Lemma 2, (86),

$$\left| E(\mathfrak{p}_1, \ldots, \mathfrak{p}_{i-1}, \mathfrak{p}_i, \mathfrak{p}_{i+1}, \ldots, \mathfrak{p}_N) - E(\mathfrak{p}_1, \ldots, \mathfrak{p}_{i-1}, \mathfrak{p}'_i, \mathfrak{p}_{i+1}, \ldots, \mathfrak{p}_N) \right| \leq \frac{M}{N}, \quad \mathfrak{p}_i, \mathfrak{p}'_i \in \mathfrak{P}.$$

Applying Mc'Diarmid's inequality (see [45, Propositio 1.3]), yields

$$\text{Prob}\left\{ \left| E(\mathfrak{p}_1, \ldots, \mathfrak{p}_N) - \mathbb{E}_{\widehat{\mathfrak{P}} \sim \mu^N}[E(\mathfrak{p}_1, \ldots, \mathfrak{p}_N)] \right| \geq t \right\} \leq 2 \exp\left\{ -2t^2/(N(M/N)^2) \right\}$$
$$= 2 \exp\left\{ -\frac{2Nt^2}{M^2} \right\}. \tag{95}$$

Hence, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, one has

$$\sup_{s_r \in \mathcal{H}_r(\Theta)} \{R(s_r) - \hat{R}_{\widehat{\mathfrak{P}}}(s_r)\} \leq \mathbb{E}_{\widehat{\mathfrak{P}} \sim \mu^N}\left[ \sup_{s_r \in \mathcal{H}_r(\Theta)} \{R(s_r) - \hat{R}_{\widehat{\mathfrak{P}}}(s_r)\} \right] + M\sqrt{\frac{\log \frac{1}{\delta}}{2N}}. \tag{96}$$

It remains to bound expectations of the above type on the right of (96). To that end, one can resort to the so-called (empirical) *Rademacher-Complexity*, defined as follows (see e.g. [43, Definition 3.1]). For $\widehat{\mathfrak{P}}$ as above let $\mathcal{F}$ denote a class of functions mapping $\mathfrak{P}$ into $[0, M]$. The empirical Rademacher Complexity of $\mathcal{F}$ is defined by

$$\hat{\mathfrak{R}}_{\widehat{\mathfrak{P}}}(\mathcal{F}) := \mathbb{E}_\epsilon\left[ \sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N \epsilon_i f(\mathfrak{p}_i) \right],$$

where $\epsilon \in \{-1, 1\}^N$ are independent random variables taking values $\pm 1$ with equal probability (Rademacher variables). The expectation of $\hat{\mathfrak{R}}_{\widehat{\mathfrak{P}}}(\mathcal{F})$

$$\mathfrak{R}_N(\mathcal{F}) := \mathbb{E}_{\widehat{\mathfrak{P}} \sim \mu^N}\left[ \hat{\mathfrak{R}}_{\widehat{\mathfrak{P}}}(\mathcal{F}) \right]$$

is referred to as *Rademacher Complexity*. Specifically consider functions of the form

$$f(\mathfrak{p}; \theta) := \mathcal{L}(s_r(\mathfrak{p}; \theta); \mathfrak{p}),$$

and let the class $\mathcal{F}$ be comprised of all such functions obtained when $s_r$ ranges over $\mathcal{H}_r(\Theta)$, and $\mathfrak{p} \in \mathfrak{P}$, $\theta \in \Theta$. In view of Assumption 1 and Lemma 2, each $f \in \mathcal{F}$ maps $\mathfrak{P}$ into $[0, M]$. Notice next that, in these terms (94) can be rewritten as

$$\sup_{s_r \in \mathcal{H}_r} \{R(s_r) - \hat{R}_{\widehat{\mathfrak{P}}}(s_r)\} = \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{\mathfrak{p} \sim \mu}[f(\mathfrak{p})] - \frac{1}{N} \sum_{\mathfrak{p} \in \widehat{\mathfrak{P}}} f(\mathfrak{p}) \right\}.$$

21

It is well-known (see for instance, [46] or [45, Proposition 4.3]) that the expectation of this quantity can be bounded in terms of the Rademacher Complexity as

$$\mathbb{E}_{\widehat{\mathfrak{P}}\sim\mu^N}\left[\sup_{f\in\mathcal{F}}\left|\mathbb{E}_{\mathfrak{p}\sim\mu}[f(\mathfrak{p})]-\frac{1}{N}\sum_{\mathfrak{p}\in\widehat{\mathfrak{P}}}f(\mathfrak{p})\right|\right]\leq 4\mathfrak{R}_N(\mathcal{F}). \tag{97}$$

Moreover, since Lemma 2, (87), applies to the class $\mathcal{F}$, so that the "contraction principle" yields (see [45, Proposition 4.3.])

$$\mathfrak{R}_N(\mathcal{F})\leq L\mathfrak{R}_N(\mathcal{H}_r(\Theta)), \tag{98}$$

where $L$ is the Lipschitz constant. Combining this with (96) and (97), provides that with probability at least $1-\delta$ we have

$$\sup_{s_r\in\mathcal{H}_r}\left\{R(s_r)-\hat{R}_{\widehat{\mathfrak{P}}}(s_r)\right\}\leq 4L\mathfrak{R}_N(\mathcal{H}_r(\Theta))+M\sqrt{\frac{\log\frac{1}{\delta}}{2N}}, \tag{99}$$

which now shows a dependence on the complexity of the hypothesis class $\mathcal{H}_r(\Theta)$. This can be alternatively described in terms of the pseudo-dimension of $\mathcal{H}_r(\Theta)$. Relating Rademacher Complexity to covering numbers and using Dudley's integral, it can eventually be shown (see [46, 47]) that

$$\mathfrak{R}_N(\mathcal{H}_r(\Theta))\leq C\sqrt{\frac{2P\log\frac{eN}{P}}{N}}, \tag{100}$$

where $P=\operatorname{Pdim}\mathcal{H}_r(\Theta)$ is the pseudo-dimension of $\mathcal{H}_r(\Theta)$ and $C$ is an absolute constant. Substituting this bound into (99), combining the result with (93), we conclude that with probability at least $1-\delta$

$$R(\hat{\theta})\leq R(s_r^*)+4LC\sqrt{\frac{2P\log\frac{eN}{P}}{N}}+M\sqrt{\frac{\log\frac{1}{\delta}}{2N}}. \tag{101}$$

The assertion follows now from bounding the approximation error $R(s_r^*)$ with the aid of Corollary 1. □

# 6 Numerical experiments

All computations were run on a Linux x86_64 workstation equipped with dual AMD EPYC 9334 CPUs and an NVIDIA L40S GPU. For all finite element solves and weight-matrix assembly we used Python 3.13 with the libraries DOLFINx 0.9.0 [48] and scifem 0.7.0, together with hIPPYlib 3.1.0 [49] (built on legacy FEniCS 2019.1.0 [42]) to generate samples of the Gaussian random parameter fields. Neural networks were implemented and trained in PyTorch 2.6.0+cu124 [50].

## 6.1 Problem setup

We consider two stationary diffusion problems—a heat conduction model with a piecewise-constant conductivity field and a Darcy flow model with a lognormal random permeability—and one linear elasticity problem describing a clamped beam under traction. Together, these examples cover both smooth and rough random parameter fields, as well as scalar-valued (diffusion) and vector-valued (elasticity) PDEs, and include mixed Dirichlet–Neumann boundary conditions typical of the three applications.

### 6.1.1 Heat conduction

We consider a steady-state heat conduction problem with a piecewise-constant conductivity field and a nonzero external heat source, motivated by electronics thermal management. The conductivity field $\mathfrak{p}(x)$

is assumed to be constant over a $4 \times 4$ arrangement of mini-squares uniformly spread across the domain $\Omega = (0,1) \times (0,1)$, that is,

$$\mathfrak{p}(x) = \mathbf{1}_{\Omega \setminus (\cup \bar{\Omega}_i)}(x) + \sum_{i=1}^{16} 10^{\mu_i} \mathbf{1}_{\bar{\Omega}_i}(x),$$

where $\mathbf{1}_{\bar{\Omega}_i}(x)$ denotes the indicator function of subset $\bar{\Omega}_i \subset \Omega$. Each $\bar{\Omega}_i$ is defined as

$$\bar{\Omega}_i = \left[\frac{m}{8} - \frac{1}{16}, \frac{m}{8} + \frac{1}{16}\right] \times \left[\frac{n}{8} - \frac{1}{16}, \frac{n}{8} + \frac{1}{16}\right], \qquad m, n \in \{1, 3, 5, 7\},$$

so that the 16 pairs $(m, n)$ define 16 mini-squares. Each $\mu_i \in \mathbb{R}$ is independently sampled from the uniform distribution $\mathcal{U}(-1, 1)$, so that the parameter vector $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_{16})$ lies in $\mathbb{R}^{16}$. The boundary data are prescribed as

$$u_0(x) = 0.1(1 - x_1)\sin(4\pi x_2), \quad x \in \Gamma_{\text{left}} \cup \Gamma_{\text{right}} =: \Gamma_D,$$
$$g(x) = 0.1(1 - x_2)\cos(2\pi x_1), \quad x \in \Gamma_{\text{top}} \cup \Gamma_{\text{bottom}} := \Gamma_N,$$

where $\Gamma_{\text{left}}, \Gamma_{\text{right}}, \Gamma_{\text{top}}, \Gamma_{\text{bottom}}$ denote the left, right, top and bottom boundary of the domain $\Omega$. The source term $f$ is defined by its action on any test function $v \in H^1_{0,\Gamma_D}(\Omega)$

$$\langle f, v \rangle = (\mathbf{f}_1, \boldsymbol{\nabla} v)_{L^2} + (f_2, v)_{L^2},$$

where

$$\mathbf{f}_1(x) = (0.5 \times \mathbf{1}_{\cup \bar{\Omega}_i}(x), -0.5 \times \mathbf{1}_{\cup \bar{\Omega}_i}(x))^\top \text{ and } f_2(x) = 1.$$

This representation fits the flux-free decomposition of the source used in the loss formulation. Note that $f \in H^{-1}(\Omega) \subset (H^1_{0,\Gamma_D}(\Omega))'$, while $f \notin L^2(\Omega)$, in accordance with our setup (19). The low regularity of the source term, combined with the piecewise discontinuous conductivity field, results in limited solution regularity, see Figure 1 (top row) for an illustration of one parameter-solution pair at a random parameter sample. In fact, the finite element solution $[u_h^\circ, \sigma_h^\circ]$ (with elements $\text{RT}_1^\circ \times \text{CG}_2^\circ$ on a mesh of size $128 \times 128$, where the edges of the mini-squares are aligned with the mesh edges) exhibits steep gradients around the mini-squares.

### 6.1.2 Darcy flow

In this case, we consider a Darcy flow problem with a lognormal diffusion field to model subsurface flow with a random permeability field. Specifically, we assume $\Omega = (0,1) \times (0,1)$ and

$$\mathfrak{p}(x) = 0.01 + \exp(m(x)), \tag{102}$$

where $m \sim \mathcal{N}(\bar{m}, \mathcal{C})$ with mean $\bar{m} = 0$ and covariance operator $\mathcal{C} := (\delta I - \gamma \Delta)^{-\alpha}$. The parameters $\delta, \gamma, \alpha > 0$ collectively control the correlation, variance, and smoothness of the random field. For demonstration, we set $\delta = 1.5$, $\gamma = 0.15$, and $\alpha = 2$. To sample the Gaussian random field $m$ we use the hIPPYlib implementation [49], which realizes the covariance operator $\mathcal{C}$ via an auxiliary elliptic problem with Robin boundary conditions.

The boundary data for the Darcy solution are specified as

$$u_0(x) = 1 - x_1, \quad x \in \Gamma_{\text{left}} \cup \Gamma_{\text{right}} =: \Gamma_D,$$
$$g(x) = 0, \quad x \in \Gamma_{\text{top}} \cup \Gamma_{\text{bottom}} := \Gamma_N.$$

The source term, representing water extraction at nine wells, is defined as

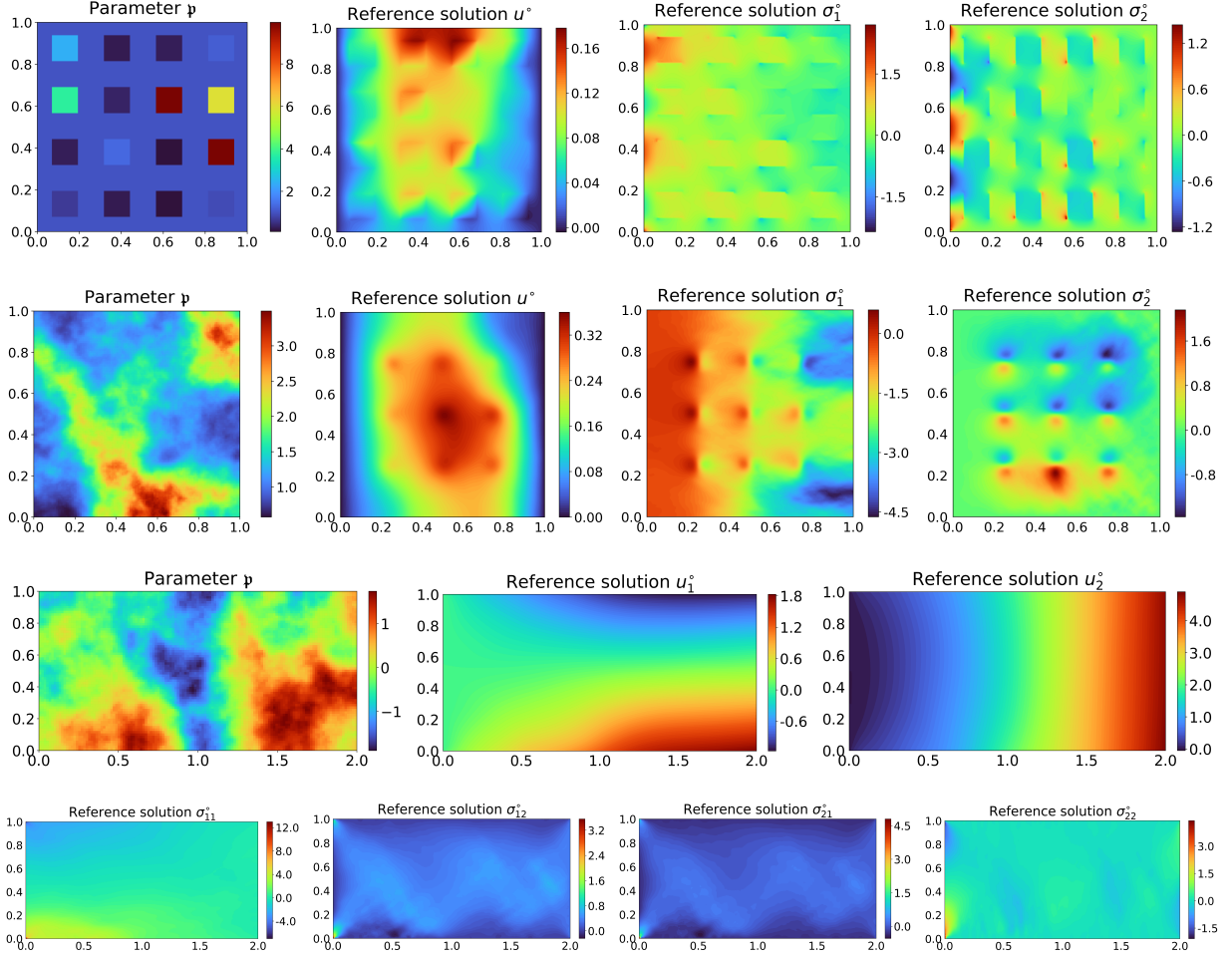$$f(x) = \sum_{i=1}^{9} 100 \exp\left(-\left(\frac{\|x - c_i\|_2}{w}\right)^2\right), \tag{103}$$

23

Figure 1: Visualization of parameter-to-solution map $\mathfrak{p}_h \mapsto [u_h^\circ(\mathfrak{p}_h), \sigma_h^\circ(\mathfrak{p}_h)]$ at a random parameter sample $\mathfrak{p}_h$ (left) for the heat conduction (first row), Darcy flow (second row), and linear elasticity setup (third and fourth rows).

where the centers are given by $c_i = (m/4, n/4)$ for $m, n = 1, 2, 3$ and the width is $w = 1/32$. The second row of Figure 1 displays one parameter-solution pair at a random parameter sample, with the finite element solution $[u_h^\circ, \sigma_h^\circ]$ computed with elements $\mathrm{RT}_1^\circ \times \mathrm{CG}_2^\circ$ on a mesh of size $128 \times 128$, and the finite element parameter sample computed using $\mathrm{CG}_1$ elements on the same mesh.

### 6.1.3 Linear elasticity

We consider a rectangular elastic body clamped along its left edge, with an upward force applied to its right edge. In this experiment, the physical domain is set to $\Omega = (0, 2) \times (0, 1)$. The Gaussian measure for the random parameter field $\mathfrak{p}$ is characterized by $\bar{\mathfrak{p}} = 0, \delta = 1.5, \gamma = 0.15$, and $\alpha = 2$, and it enters the model through spatial variations of the stiffness tensor. The Poisson ratio is set to $\nu = 0.4$. The boundary data are prescribed as

$$\underline{u}_0(x) = (0,0)^\top, \quad x \in \Gamma_{\text{left}},$$

$$\underline{t}(x) = \begin{cases} (0,0)^\top, & x \in \Gamma_{\text{top}} \cup \Gamma_{\text{bottom}}, \\ \left(0.6 \exp\left(-(x_2 - 0.5)^2/4\right), 0.3(1 + x_2/10)\right)^\top, & x \in \Gamma_{\text{right}}, \end{cases}$$

That is, the left edge is clamped, while on the right edge we apply a horizontal traction given by a Gaussian bump centered at $x_2 = 0.5$ and a vertical traction that increases linearly with $x_2$; all other boundary

24

segments are traction-free. Equivalently, $\Gamma_D := \Gamma_{\text{left}}$ and $\Gamma_N := \Gamma_{\text{top}} \cup \Gamma_{\text{bottom}} \cup \Gamma_{\text{right}}$. One parameter-solution pair at a random parameter sample is illustrated in the third and fourth rows of Figure 1, where the finite element solution $[\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]$ is computed with elements $(\text{RT}_1^\circ)^2 \times (\text{CG}_2^\circ)^2$ on a mesh of size $256 \times 128$. The finite element parameter sample is computed using $\text{CG}_1$ elements on the same mesh. We can observe a corner singularity of the stress tensor at the clamped left edge, where the boundary conditions change from Dirichlet to Neumann type, yielding low solution regularity.

### 6.1.4 Variational lift of boundary data

The variational lift of the boundary data by harmonic extension to the domain is shown in Figure 2. These fields are precomputed by solving (17) and (20) for the diffusion problems using $\text{CG}_m$ elements (matching the order $m$ of $u_h^\circ$) and by solving (39) and (42) for the linear elasticity problem using $\text{CG}_m^2$ elements.
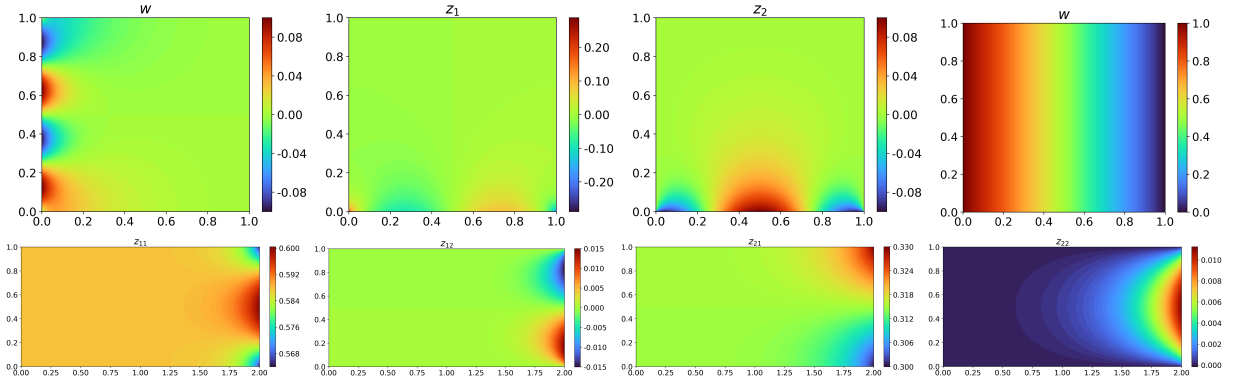


Figure 2: Auxiliary variables $w$ and $z = (z_1, z_2)$ for heat conduction (top left three plots), $w$ for Darcy flow (top right) with $z = (0,0)$, obtained by solving auxiliary problems (17) and (20). Bottom: $\underline{\underline{z}} = (z_{11}, z_{12}; z_{21}, z_{21})$ by solving (42) for linear elasticity with $\underline{w} = (0,0)$. These variables encode the inhomogeneous Dirichlet and Neumann boundary data.

## 6.2 Conforming finite element approximations

We verify that the conforming FOSLS FE discretization yields sufficiently small FE loss and exhibits the expected convergence rates as in Theorem 1 for the diffusion problem and Theorem 4 for the linear elasticity problem. Moreover, we compare the cost of assembling the loss weights to that of solving the underlying PDEs.

To assess the FE loss, we solve the FOSLS normal equations (52) for the diffusion problem using elements $\Sigma_h \times \mathbb{U}_h = \text{RT}_k^\circ \times \text{CG}_{k+1}^\circ$ and (57) for the linear elasticity problem using elements $\Sigma_h \times \mathbb{U}_h = (\text{RT}_k^\circ)^d \times (\text{CG}_{k+1}^\circ)^d$, both with different mesh sizes $h$ and orders $k = 0, 1$, for 100 independent parameter samples, and evaluate the resulting FE loss by inserting the FE solutions into the corresponding residual loss. The mean FE loss over the 100 samples, for various mesh sizes and FE spaces, is recorded in Table 1. For all cases, we also compute the mean FE error (squared error) in the $\mathbb{H} = H(\text{div}) \times H^1$ norm compared to the reference solutions computed using higher order elements ($k = 2$). The close FE losses and FE errors confirm the error-residual equivalence (55) for the diffusion problem and (C.1) for the linear elasticity problem. The computational cost of assembling the sparse weight matrix $W_{\mathfrak{p}}$ (with sparsity pattern shown in Figure D.10) scales linearly (especially for the same FE orders) with the number of degrees of freedom (DoFs) $N_h^s$ and is much smaller than the cost of solving the algebraic system of size $N_h^s \times N_h^s$ arising from the FE discretization of the PDE; see the rightmost columns of Table 1. As expected, refining the mesh and/or increasing the FE order reduces the FE loss and improves the FE approximations, at the expense of larger assembly and solve times.

Figure 3 shows, for a representative random parameter sample, the convergence of FE loss with respect to the mesh size $h$ and the FE order $k$ across the three problem setups. In all cases, the loss decreases

| Problem | Mesh | FE space | # DoFs | FE loss | FE error | $W_\mathfrak{p}$ | Solve |
|---|---|---|---|---|---|---|---|
| Heat Conduction | $64\times64$ | $RT_0\times CG_1$ | 16,641 | $4.86\times10^{-3}$ | $5.83\times10^{-3}$ | 0.015 | 0.10 |
| | | $RT_1\times CG_2$ | 57,857 | $3.52\times10^{-4}$ | $4.55\times10^{-4}$ | 0.051 | 0.38 |
| | $128\times128$ | $RT_0\times CG_1$ | 66,049 | $1.55\times10^{-3}$ | $1.93\times10^{-3}$ | 0.04 | 0.45 |
| | | $RT_1\times CG_2$ | 230,401 | $1.09\times10^{-4}$ | $1.46\times10^{-4}$ | 0.22 | 1.76 |
| Darcy Flow | $128\times128$ | $RT_0\times CG_1$ | 66,049 | $9.54\times10^{-1}$ | $9.55\times10^{-1}$ | 0.27 | 0.52 |
| | | $RT_1\times CG_2$ | 230,401 | $4.00\times10^{-3}$ | $4.01\times10^{-3}$ | 0.49 | 1.88 |
| | $256\times256$ | $RT_0\times CG_1$ | 263,169 | $2.40\times10^{-1}$ | $2.40\times10^{-1}$ | 1.03 | 2.38 |
| | | $RT_1\times CG_2$ | 919,553 | $2.53\times10^{-4}$ | $2.53\times10^{-4}$ | 2.10 | 8.36 |
| Elasticity | $128\times64$ | $RT_0^2\times CG_1^2$ | 66,306 | $6.51\times10^{-3}$ | $2.77\times10^{-2}$ | 0.20 | 0.63 |
| | | $RT_1^2\times CG_2^2$ | 230,914 | $7.37\times10^{-4}$ | $9.29\times10^{-4}$ | 1.41 | 3.22 |
| | $256\times128$ | $RT_0^2\times CG_1^2$ | 263,682 | $2.46\times10^{-3}$ | $5.96\times10^{-3}$ | 0.77 | 3.00 |
| | | $RT_1^2\times CG_2^2$ | 920,578 | $3.00\times10^{-4}$ | $3.32\times10^{-4}$ | 5.59 | 14.41 |

Table 1: Mean FE loss and mean FE error (squared) in $\mathbb{H}$-norm over 100 parameter samples, and wall-clock time (in seconds) to assemble the weight matrix $W_\mathfrak{p}$ in (81) and to solve the PDE using a direct (LU) solver for the diffusion (30) and linear elasticity (48) problems, with different mesh sizes, FE spaces, and the corresponding number of degrees of freedom (DoFs) $N_h^s$.
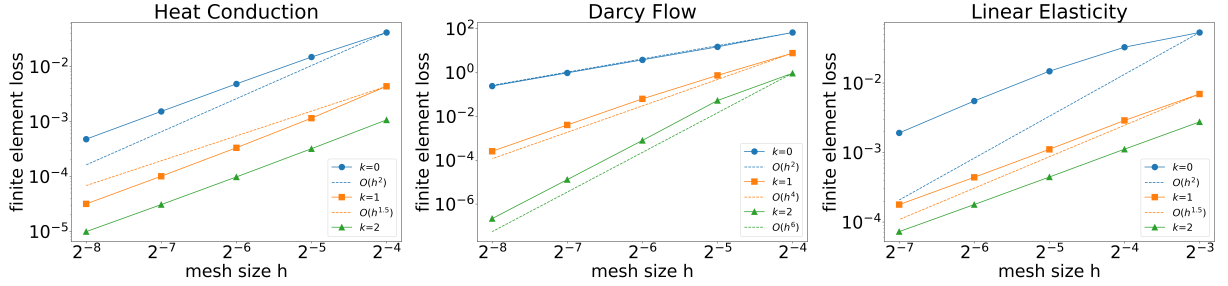


Figure 3: Convergence of the FE loss with respect to mesh size ($h$) and FE order ($k = 0, 1, 2$) for a representative parameter sample. Solid lines show the measured losses and dashed lines the reference rates ($O(h^{2(k+1)})$); Darcy flow follows the predicted asymptotic behavior, while convergence in heat conduction and linear elasticity is limited by reduced solution regularity due to discontinuous coefficients and boundary-induced corner singularities, confirming the analysis in Theorem 1 and 4.

monotonically as the mesh is refined and as the order increases. The measured losses (solid lines) follow closely the expected asymptotic convergence rates $O(h^{2(k+1)})$ (dashed lines) for the Darcy flow, confirming that the FE discretization exhibits the theoretically predicted convergence behavior with respect to both $h$ and $k$ as in (56) of Theorem 1. The observed convergence rates are limited by the lower regularity of the solutions due to the discontinuous conductivity field in heat conduction, and the mixed boundary conditions in the linear elasticity problem that lead to corner singularities.

## 6.3   Reduced basis approximations

We now investigate the equivalence of the RB loss and errors in Theorem 2 and their dependence on the number of RB functions, and how this choice affects computational cost. We first construct the RB spaces as in Section 5.1 from $N_{\mathrm{POD}} = 500$ and 1000 independent parameter samples and the corresponding high-fidelity FE solutions (using $RT_1\times CG_2$ elements on a $128\times128$ mesh) for the heat conduction and Darcy flow problem, respectively. For the linear elasticity problem, we use $N_{\mathrm{POD}} = 1000$ independent parameter samples and the corresponding high-fidelity FE solutions (using $RT_1^2\times CG_2^2$ elements on a $128\times64$ mesh). Then for each of 500 test random parameter samples, we compute the FE weight matrix $W_\mathfrak{p}$ in the FE loss (81), the corresponding RB weight matrix $W_\mathfrak{p}^r$ and vector $\boldsymbol{\alpha}_\mathfrak{p}^r$ in the RB loss (82) by projection, and solve the reduced normal equation (83) for the coefficient vector $\boldsymbol{s}_r(\mathfrak{p})$ of the RB solution $s_r(\mathfrak{p})$. We also compute the FE solutions $s_h(\mathfrak{p})$ and $\bar{s}_h(\mathfrak{p})$, using the same mesh and elements of $RT_1\times CG_2$ and $RT_3\times CG_4$, respectively, where we take the latter as the "ground truth" solution. Finally, we compute the RB loss (82) of the RB solution and the squared error of the RB solution in $\mathbb{H} = H(\mathrm{div}\,)\times H^1$-norm with respect to the corresponding FE solutions.

26

By the decomposition of the RB loss (71) in Theorem 2, we expect the gap between the RB loss and the FE loss to be equivalent (up to problem-dependent constants) to the squared error of the RB solution relative to the FE solution, i.e., $\mathcal{L}(s_r(\mathfrak{p}); \mathfrak{p}) - \mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p}) \approx \|s_h(\mathfrak{p}) - s_r(\mathfrak{p})\|_{\mathbb{H}}^2$. This relationship is demonstrated in Figure 4. We observe this equivalence by comparing an empirical mean squared error approximation to $\mathbb{E}_{\mathfrak{p}\sim\mu}\big[\|\|s_r(\mathfrak{p}) - s_h(\mathfrak{p})\|\|_{\mathbb{H}}^2\big]$ against the empirical mean loss difference, approximating $\mathbb{E}_{\mathfrak{p}\sim\mu}\big[\mathcal{L}(s_r(\mathfrak{p}); \mathfrak{p}) - \mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p})\big]$, for an increasing number of basis functions across the three problems. Both quantities are evaluated via sample average approximation using a shared set of 500 test samples, which was found to be sufficient for low-variance estimation of the mean.



Figure 4: Comparison between the empirical mean squared error $\mathbb{E}_{\mathfrak{p}\sim\mu}\big[\|\|s_r(\mathfrak{p}) - s_h(\mathfrak{p})\|\|_{\mathbb{H}}^2\big]$ and the empirical mean loss difference $\mathbb{E}_{\mathfrak{p}\sim\mu}\big[\mathcal{L}(s_r(\mathfrak{p}); \mathfrak{p}) - \mathcal{L}(s_h(\mathfrak{p}); \mathfrak{p})\big]$. Both quantities are estimated over 500 random samples (using $\mathrm{RT}_1 \times \mathrm{CG}_2$ elements for $s_h$), confirming the RB loss decomposition (71) in Theorem 2.

The comparison of the mean squared error of the RB solution $s_r$ (compared to the FE solution $s_h$), the $X_h$-projection error of the FE solution $s_h \in \mathbb{H}_h$ onto the RB space $\mathbb{H}_r$, and the error estimate by the trailing eigenvalues in (63), is shown in Figure 5, confirming the quasi-optimality of the RB solution in (72) in Theorem 2 and the tight error estimate. The faster decay of the error estimate by the trailing eigenvalues at a large number of RB basis functions is due to the limited number $N_{\mathrm{POD}}$ of samples in computing the eigenvalues (which leads to inaccurate estimation of small eigenvalues).
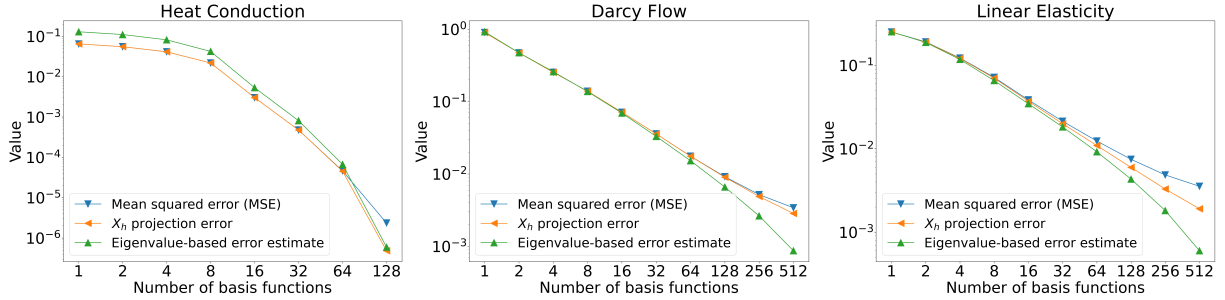


Figure 5: Comparison between the empirical mean squared error $\mathbb{E}_{\mathfrak{p}\sim\mu}\big[\|\|s_r(\mathfrak{p}) - s_h(\mathfrak{p})\|\|_{\mathbb{H}}^2\big]$ of the RB solution $s_r$ (approximated over 500 samples with respect to the FE solution $s_h$ using $\mathrm{RT}_1 \times \mathrm{CG}_2$ elements), the square of the $X_h$-projection error of $s_h$ onto the RB space $\mathbb{H}_r$, and the eigenvalue-based error estimate (63). These results illustrate the quasi-optimality (72) of the RB approximation in Theorem 2 and the tightness of the error estimate.

Moreover, the comparison of the mean RB loss and the mean squared error of the RB solution $s_r$ (compared to the "ground truth" solution $\bar{s}_h$) with respect to the number of RB functions is shown in Figure 6, demonstrating the residual-error equivalence in (74) in Theorem 2. The slower decay of the mean squared error of the RB solution for a large number of RB basis functions is due to the FE discretization errors (using $\mathrm{RT}_1 \times \mathrm{CG}_2$ elements) compared to the "ground truth" (FE discretization using $\mathrm{RT}_3 \times \mathrm{CG}_4$ elements). We observe a noticeably faster decay of both loss and error in heat conduction compared to Darcy flow and the linear elasticity problem, reflecting the smaller Kolmogorov $n$-widths of the manifold of the parameter-to-solution map in heat conduction.
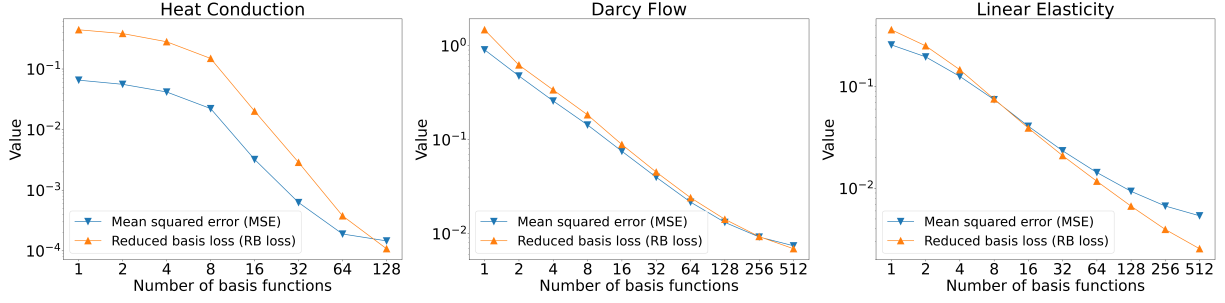
27

Figure 6: Comparison between the empirical mean squared error $\mathbb{E}_{\mathfrak{p}\sim\mu}\left[||s_r(\mathfrak{p}) - \bar{s}_h(\mathfrak{p})||^2_{\mathbb{H}}\right]$ of the RB solution $s_r$ (relative to the "ground truth" FE solution $\bar{s}_h$ with RT$_3\times$CG$_4$ elements) and the empirical mean RB loss $\mathbb{E}_{\mathfrak{p}\sim\mu}\left[\mathcal{L}(s_r(\mathfrak{p});\mathfrak{p})\right]$. Both averages are computed using 500 random samples, confirming the residual-error equivalence (74) in Theorem 2.

## 6.4 Reduced basis neural operator

As defined in hypothesis class (64), the reduced basis neural operator (RBNO) consists of two parts. In the first part, a neural network maps the parameter $\mathfrak{p}$ (represented as a grid-based image of size $\mathbb{R}^{H\times W}$, e.g., $H = W = 129$ for the diffusion problem) to the RB coefficients $\boldsymbol{s}_r(\mathfrak{p};\theta) \in \mathbb{R}^r$ of the approximate solution; this stage is computed in single precision (float32). In the second part, the RBNO solution is reconstructed as a linear combination of the precomputed RB basis functions with the predicted RB coefficients, i.e., $s_r(\mathfrak{p};\theta) = \Phi_r\boldsymbol{s}_r(\mathfrak{p};\theta)$, and this reconstruction is evaluated in double precision (float64). Note that the FE and RB weight matrices are also computed in double precision to ensure accurate loss evaluation.

We approximate the parameter-to-RB coefficient map by a deep neural network

$$\boldsymbol{s}_r(\cdot;\theta) = L_2 \circ \sigma \circ L_1 \circ \text{Flatten} \circ C_L \circ \cdots \circ C_1 : \mathbb{R}^{H\times W} \to \mathbb{R}^r,$$

where $L_1, L_2$ are fully-connected linear layers, and each $C_l$ is a convolutional block followed by a nonlinear activation function $\sigma$, i.e., $C_l = \sigma \circ \text{Conv}_l$. We use the LeakyReLU activation function $\sigma(x) = \max(0, x) + 0.01\min(0, x)$, and initialize all network parameters $\theta$ using the Xavier uniform method [1]. We train the neural network by minimizing the empirical RB loss (85) using the SOAP optimizer [51], which we find to outperform the commonly used AdamW optimizer for faster convergence and higher accuracy. More details on the architecture and training can be found in Appendix D.3.

The decay of the empirical loss on the training set (with 1000 and 3000 samples, respectively) and validation set (with 500 samples) is shown in Figure 7 for the three problems. During training, we monitor the validation loss and retain the network weights that achieve a minimum validation loss (early stopping), rather than those minimizing the training loss, to prevent overfitting and improve generalization to unseen parameter samples. From 1000 to 3000 training samples, we can observe a larger training loss at the end of training, a smaller validation loss at the minimum values, and a closer gap between them.

To demonstrate Theorem 3, we show in Figure 8 how three quantities decay with increasing number of training samples from 16 to 64, 256, 1024, and 4096. Specifically, we report (1) the RBNO loss $\mathbb{E}_{\mathfrak{p}\sim\mu}\left[\mathcal{L}(s_r(\mathfrak{p};\hat{\theta});\mathfrak{p})\right]$ at the RBNO solution $s_r(\mathfrak{p};\hat{\theta})$, (2) the mean square error between the RBNO solution and the RB solution $\mathbb{E}_{\mathfrak{p}\sim\mu}\left[||s_r(\mathfrak{p}) - s_r(\mathfrak{p};\hat{\theta})||^2_{\mathbb{H}}\right]$, and (3) the mean square error between the RBNO solution and the "ground truth" FE solution (with RT$_3\times$CG$_4$ elements) $\mathbb{E}_{\mathfrak{p}\sim\mu}\left[||\bar{s}_h(\mathfrak{p}) - s_r(\mathfrak{p};\hat{\theta})||^2_{\mathbb{H}}\right]$, where the expectation is evaluated with 500 test samples.

In Figure 9, we show the histograms of the ratios between the RBNO error (compared to the "ground truth" FE solution) and the square root of the RBNO residual loss for the three problems with 500 test samples. The ratios cluster near 1, indicating that the residual loss is a good estimator of the error.

## 6.5 Comparison to other neural operators

In this section, we compare RBNO with two popular neural operators, detailed below. To place these comparisons in proper perspective, we note a key distinction in objectives. The central goal of RBNO is to
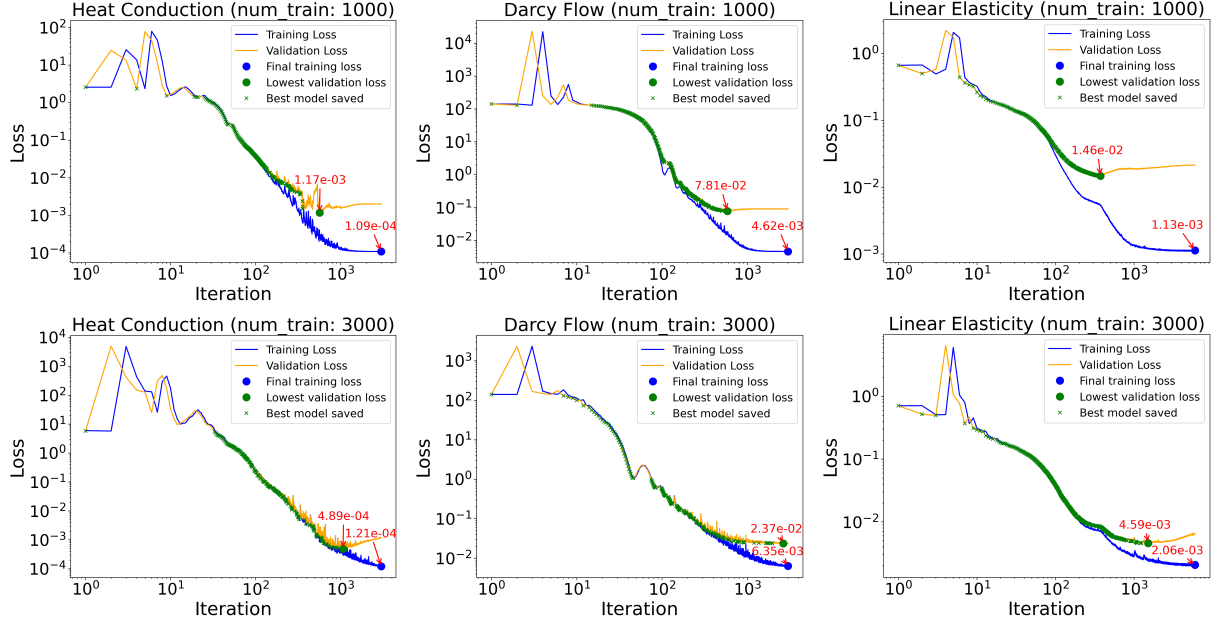
Figure 7: Empirical loss decay over optimization iterations on the training set with 1000 samples (top) and 3000 samples (bottom), and on the validation set with 500 samples, for the three problems. During training, model checkpoints are updated whenever the validation loss decreases, and the final model is selected as the checkpoint with the lowest validation loss.
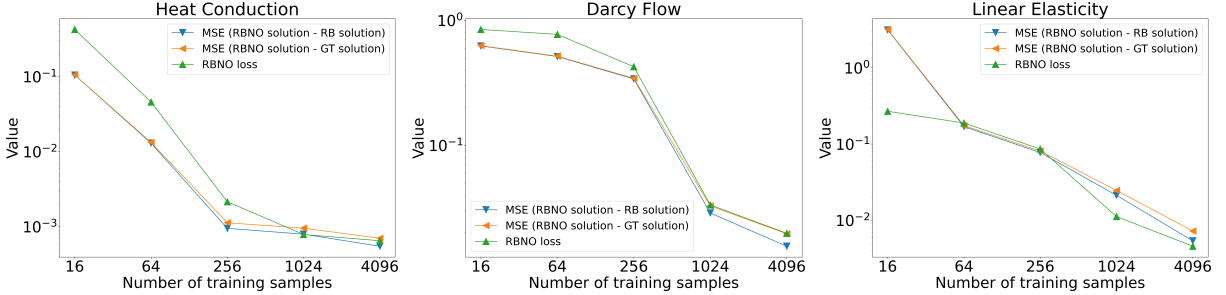


Figure 8: Decay of the empirical RBNO loss $\mathbb{E}_{\mathfrak{p}\sim\mu}\left[\mathcal{L}(s_r(\mathfrak{p};\hat{\theta});\mathfrak{p})\right]$ at the RBNO solution $s_r(\mathfrak{p};\hat{\theta})$, the empirical mean square error between the RBNO solution and the RB solution $\mathbb{E}_{\mathfrak{p}\sim\mu}\left[||s_r(\mathfrak{p})-s_r(\mathfrak{p};\hat{\theta})||_{\mathbb{H}}^2\right]$, and the empirical mean square error between the RBNO solution and the "ground truth" FE solution (with $\mathrm{RT}_3\times\mathrm{CG}_4$ elements) $\mathbb{E}_{\mathfrak{p}\sim\mu}\left[||\bar{s}_h(\mathfrak{p})-s_r(\mathfrak{p};\hat{\theta})||_{\mathbb{H}}^2\right]$ with increasing number of training samples from 16 to 64, 256, 1024, and 4096, where the expectation is evaluated with 500 test samples.

directly output the physically relevant quantities, both $u$ and $\sigma$, e.g., temperature/pressure and thermal/flow flux in stationary diffusion or displacement and stress in linear elasticity. In many scenarios, these are the variables of primary interest, and their errors in PDE-compliant norms can be directly measured by the FOSLS residual loss. In contrast, the neural operators discussed here typically provide pointwise evaluations of the solution $u$ at discrete sites within the computational domain. Consequently, any comparison requires post-processing the neural operator outputs, an additional step that may influence the final solution quality (see detailed investigations in Appendix D.4). To ensure a fair and rigorous comparison, we apply these neural operators to the solution components $u$ and $\sigma$ of the FOSLS formulation—a step rarely taken in standard practice—to promote variational correctness.

The first neural operator is the PCA-based neural network (PCA-Net) [27], which projects both the parameter field and the solution onto low-dimensional PCA bases and then learns a map between their PCA coefficients using a multilayer perceptron (MLP) with a mean squared error loss. The PCA bases
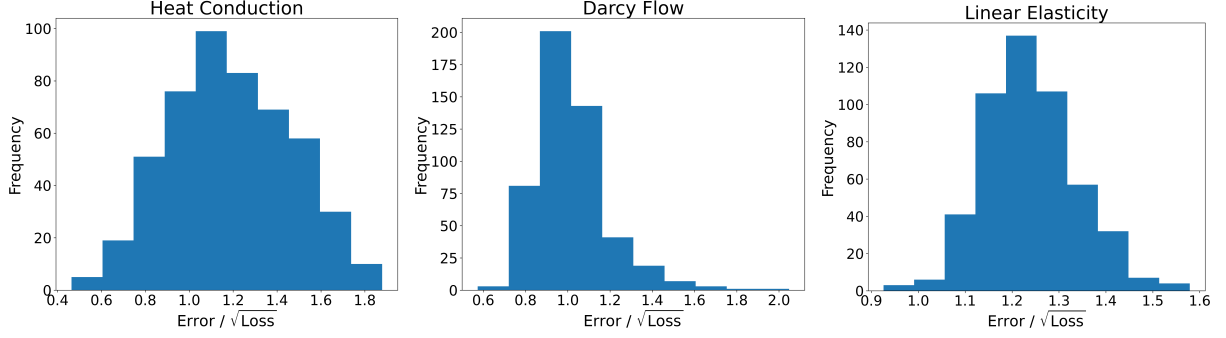
Figure 9: Histograms of the ratio Error/$\sqrt{\text{Loss}}$, between the RBNO error $||\bar{s}_h(\mathfrak{p}) - s_r(\mathfrak{p}; \hat{\theta})||_{\mathbb{H}}$ (with "ground truth" FE solution $\bar{s}_h(\mathfrak{p})$ using $\text{RT}_3 \times \text{CG}_4$ elements) and the square root of the RBNO residual loss $\mathcal{L}(s_r(\mathfrak{p}; \hat{\theta}); \mathfrak{p})$, over 500 test samples for the three problems.

are computed as the eigenvectors of empirical covariance of the pointwise evaluation of the training data $\mathfrak{p} \mapsto [\sigma_h^\circ(\mathfrak{p}), u_h^\circ(\mathfrak{p})]$ as in [27]. The second is the Fourier Neural Operator (FNO) [8], which lifts the input field to a higher-dimensional feature space and applies a sequence of spectral convolution layers: each layer takes a Fourier transform in space, multiplies a fixed number of retained modes by trainable complex weights, and applies an inverse transform followed by a nonlinearity. This architecture learns the operator in frequency space and can be evaluated on arbitrary discretizations. In our experiments, FNO is trained with a relative $L_2$ loss on the pointwise finite element solution $[\sigma_h^\circ, u_h^\circ]$ at the mesh grid points. Details of their architectures and training setups are given in Appendix D.3. To compare them with RBNO in terms of relative errors in $\mathbb{L} = L_2 \times L_2$ and $\mathbb{H} = H(\text{div}) \times H^1$ norms, we first construct piecewise linear ($\text{CG}_1 \times \text{CG}_1$) FE functions from their output values, which are predictions at the grid points, on the same mesh. The errors are then computed with respect to the corresponding $\text{CG}_1 \times \text{CG}_1$ FE representations of the high-fidelity solutions, which are obtained by evaluating the solutions at the grid points and mapping them one-to-one to the $\text{CG}_1$ degrees of freedom. Residual losses are computed using a standard FE assembly. The results are shown in Table 2, with the smallest error or loss for each problem highlighted in bold. We observe that RBNO achieves both the smallest errors in $\mathbb{H}$-norm and the smallest residual loss, while FNO attains the smallest error in $\mathbb{L}$-norm. Notably, RBNO, despite not being trained on pointwise evaluation data, still yields a relatively small error in $\mathbb{L}$-norm. In contrast, both PCA-Net and FNO, which are trained on pointwise evaluation data, produce large relative errors in the $\mathbb{H}$-norm and violate the PDE residual with significantly large values. Constructing $\text{RT}_1 \times \text{CG}_2$ FE representations from the pointwise outputs by projection in the $\mathbb{H}$-norm does not reduce the error or loss. A more detailed investigation of the large errors and losses, as well as the performance of PCA-Net and RBNO trained using the same RB coefficients, is provided in Appendix D.4 and Appendix D.5.

| Method | Relative error in $\mathbb{L}$-norm | Relative error in $\mathbb{H}$-norm | Residual loss |
|---|---|---|---|
| | **Heat Conduction** | | |
| PCA-Net | $1.27 \times 10^{-1}$ $(2.52 \times 10^{-2})$ | $3.32 \times 10^{-1}$ $(5.75 \times 10^{-2})$ | $1.84 \times 10^{1}$ $(3.01\ )$ |
| FNO | $\mathbf{1.82 \times 10^{-2}}(\mathbf{1.66 \times 10^{-3}})$ | $1.46 \times 10^{-1}$ $(1.53 \times 10^{-2})$ | $1.90 \times 10^{1}$ $(4.12\ )$ |
| RBNO | $2.43 \times 10^{-2}$ $(9.35 \times 10^{-3})$ | $\mathbf{1.86 \times 10^{-2}}(\mathbf{6.36 \times 10^{-3}})$ | $\mathbf{5.33 \times 10^{-4}}(\mathbf{1.44 \times 10^{-3}})$ |
| | **Darcy Flow** | | |
| PCA-Net | $1.79 \times 10^{-1}$ $(7.53 \times 10^{-2})$ | $1.13 \times 10^{-1}$ $(2.35 \times 10^{-2})$ | $2.59$ $(2.03\ )$ |
| FNO | $\mathbf{3.30 \times 10^{-2}}(\mathbf{1.18 \times 10^{-2}})$ | $1.60 \times 10^{-1}$ $(3.81 \times 10^{-2})$ | $5.10$ $(2.34\ )$ |
| RBNO | $6.39 \times 10^{-2}$ $(2.26 \times 10^{-2})$ | $\mathbf{1.23 \times 10^{-2}}(\mathbf{3.64 \times 10^{-3}})$ | $\mathbf{2.72 \times 10^{-2}}(\mathbf{2.61 \times 10^{-2}})$ |
| | **Linear Elasticity** | | |
| PCA-Net | $1.34 \times 10^{-1}$ $(1.03 \times 10^{-1})$ | $3.27 \times 10^{-1}$ $(1.14 \times 10^{-1})$ | $2.46 \times 10^{1}$ $(6.11\ )$ |
| FNO | $\mathbf{1.17 \times 10^{-2}}(\mathbf{3.34 \times 10^{-3}})$ | $4.30 \times 10^{-1}$ $(6.11 \times 10^{-2})$ | $2.82 \times 10^{1}$ $(6.45\ )$ |
| RBNO | $2.49 \times 10^{-2}$ $(5.89 \times 10^{-3})$ | $\mathbf{3.91 \times 10^{-2}}(\mathbf{8.68 \times 10^{-3}})$ | $\mathbf{4.87 \times 10^{-3}}(\mathbf{3.56 \times 10^{-3}})$ |

Table 2: Comparison on the mean (and standard deviation, estimated over 500 test samples) of the relative errors in $\mathbb{L}_2 = L_2 \times L_2$-norm and $\mathbb{H} = H(\text{div}) \times H^1$-norm, and residual loss for three neural networks on three problems.

30

# 7  Conclusions, limitations, and future work

This work presented a variationally correct operator learning framework for constructing residual-based loss functions to train accurate surrogates for parameter-to-solution maps of stationary diffusion and linear elasticity models. Our central design principle is *variational correctness*: the residual loss is equivalent (up to constants) to the solution error in the model-induced norm. We achieved this by utilizing FOSLS formulations and incorporating mixed Dirichlet–Neumann boundary conditions via variational lifts, thereby avoiding boundary penalty terms that typically degrade norm equivalence.

Computationally, we bridged the gap between continuous theory and practical implementation by establishing a rigorous link between the discrete loss and the associated discretization errors. We proved that the total error is bounded by a sum of distinct components: finite element discretization bias, reduced basis projection error, neural network approximation error, statistical estimation error, and optimization error. To strictly enforce the function space conformity required by the FOSLS objective while mitigating high-dimensional output effects, we employed the RBNO architecture trained using a computationally efficient RB loss function. Numerical experiments validated this error decomposition and convergence analysis, confirming that the variationally correct residual serves as a tight, computable *a posteriori* error estimator.

While less reliant on data than purely data-driven regression, our method still requires a moderate number of high-fidelity solutions to construct the POD reduced basis. For this strategy to be effective, the reduced dimension must grow slowly as target accuracy increases, a requirement tied to the robust decay of Kolmogorov $n$-widths. While such decay is expected for the dissipative models treated here , our experiments indicate that the required dimension may vary substantially under less benign parameter representations. Furthermore, this favorable decay does not extend to dispersive or transport-dominated problems (e.g., Helmholtz or wave-type regimes), which will require alternative approximation concepts.

Handling boundary conditions remains a delicate issue in scientific machine learning. We addressed this by incorporating mixed boundary conditions via auxiliary lifts, preserving variational correctness without "wrong penalty terms." This comes at the cost of solving two parameter-independent auxiliary boundary value problems, assuming parameter-independent boundary data. For parameter-dependent boundary data, this approach must be adapted; a promising direction is to jointly approximate (or train) the auxiliary lift solutions alongside the main surrogate while maintaining stability.

While we restricted our attention to stationary models, extending this framework to time-dependent PDEs appears viable for parabolic problems, where Kolmogorov $n$-widths decay in a parameter-robust manner and space–time least-squares formulations are available [18, 41]. In contrast, while analogous formulations exist for wave equations [19], their dispersive nature typically precludes the existence of comparable low-dimensional bases, necessitating new hybrid representations and stability-aware learning objectives.

Finally, we aim to further develop RBNO as a surrogate model equipped with rigorous a posteriori error estimates for downstream tasks such as inverse problems [3, 52], data assimilation [53, 54], optimization under uncertainty [55, 56, 57], and optimal experimental design [58, 59, 60]. Since success in these domains often relies on derivative or sensitivity information [34, 61, 57], a critical direction for future work is the development of RBNO surrogates to provide accurate derivative information.

## References

[1] R. C. Smith, Uncertainty quantification: theory, implementation, and applications, SIAM, 2024.

[2] P. Chen, A. Quarteroni, G. Rozza, Reduced basis methods for uncertainty quantification, SIAM/ASA Journal on Uncertainty Quantification 5 (1) (2017) 813–869.

[3] A. M. Stuart, Inverse problems: a Bayesian perspective, Acta numerica 19 (2010) 451–559.

[4] P. Chen, O. Ghattas, Stein variational reduced basis Bayesian inversion, SIAM Journal on Scientific Computing 43 (2) (2021) A1163–A1193.

[5] X. Huan, J. Jagalur, Y. Marzouk, Optimal experimental design: Formulations and computations, Acta Numerica 33 (2024) 715–840.

[6] K. Wu, P. Chen, O. Ghattas, A fast and scalable computational framework for large-scale high-dimensional Bayesian optimal experimental design, SIAM/ASA Journal on Uncertainty Quantification 11 (1) (2023) 235–261.

[7] L. Lu, P. Jin, G. Pang, Z. Zhang, G. E. Karniadakis, Learning nonlinear operators via DeepONet based on the universal approximation theorem of operators, Nature machine intelligence 3 (3) (2021) 218–229.

[8] Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, A. Anandkumar, Fourier neural operator for parametric partial differential equations, arXiv preprint arXiv:2010.08895 (2020).

[9] M. Raissi, P. Perdikaris, G. Karniadakis, Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, Journal of Computational Physics 378 (2019) 686–707.

[10] Z. Dongkun, L. Lu, L. Guo, G. Karniadakis, Quantifying total uncertainty in physics-informed neural networks for solving forward and inverse stochastic problems, Journal of Computational Physics 397 (2019) 108850.

[11] M. Bachmayr, W. Dahmen, M. Oster, Variationally correct neural residual regression for parametric PDEs: On the viability of controlled accuracy, IMA Journal of Numerical Analysis (2025) draf073.

[12] E. Kharazmi, Z. Zhang, G. E. Karniadakis, hp-VPINNs: Variational physics-informed neural networks with domain decomposition, Computer Methods in Applied Mechanics and Engineering 374 (2021) 113547.

[13] V. Fanaskov, A. Rudikov, I. Oseledets, Neural functional a posteriori error estimates, arXiv preprint arXiv:2402.05585 (2024).

[14] L. Ernst, N. Rekatsinas, K. Urban, A posteriori certification for physics-informed neural networks, arXiv preprint arXiv:2502.20336 (2025).

[15] S. Cao, F. Brarda, R. Li, Y. Xi, Spectral-refiner: Fine-tuning of accurate spatiotemporal neural operator for turbulent flows, arXiv preprint arXiv:2405.17211 (2024).

[16] M. D. Gunzburger, P. B. Bochev, Least-Squares Finite Element Methods, Vol. 166 of Applied Mathematical Sciences, Springer, 2009.

[17] Z. Cai, R. Lazarov, T. A. Manteuffel, S. F. McCormick, First-order system least squares for second-order partial differential equations: Part I, SIAM Journal on Numerical Analysis 31 (6) (1994) 1785–1799.

[18] T. Führer, M. Karkulik, Space-time least-squares finite elements for parabolic equations, Comput. Math. Appl. 92 (2021) 27–36.

[19] T. Führer, R. Gonzales, M. Karkulik, Well-posedness of first order acoustic wave equations and space-time finite element approximation, arXiv preprint arXiv:2311.10536 [math.NA] (2023).

[20] J. A. Opschoor, P. C. Petersen, C. Schwab, First order system least squares neural networks, arXiv preprint arXiv:2409.20264 (2024).

[21] Z. Cai, J. Chen, M. Liu, X. Liu, Deep least-squares methods: An unsupervised learning-based numerical method for solving elliptic PDEs, Journal of Computational Physics 420 (2020) 109707.

[22] L. Lyu, Z. Zhang, M. Chen, J. Chen, Mim: A deep mixed residual method for solving high-order partial differential equations, Journal of Computational Physics 452 (2022) 110930.

[23] F. M. Bersetche, J. P. Borthagaray, A deep first-order system least squares method for solving elliptic PDEs, Computers & Mathematics with Applications 129 (2023) 136–150.

[24] X. Li, J. Wu, X. Tai, J. Xu, Y.-G. Wang, Solving a class of multi-scale elliptic PDEs by fourier-based mixed physics informed neural networks, Journal of Computational Physics 508 (2024) 113012.

[25] A. Cohen, R. DeVore, Approximation of high-dimensional parametric PDEs, Acta Numerica 24 (2015) 1–159.

[26] A. Quarteroni, A. Manzoni, F. Negri, Reduced basis methods for partial differential equations: an introduction, Vol. 92, Springer, 2015.

[27] K. Bhattacharya, B. Hosseini, N. B. Kovachki, A. M. Stuart, Model reduction and neural networks for parametric PDEs, The SMAI journal of computational mathematics 7 (2021) 121–157.

[28] L. Lu, X. Meng, S. Cai, Z. Mao, S. Goswami, Z. Zhang, G. E. Karniadakis, A comprehensive and fair comparison of two neural operators (with practical extensions) based on fair data, Computer Methods in Applied Mechanics and Engineering 393 (2022) 114778.

[29] N. Dal Santo, S. Deparis, L. Pegolotti, Data driven approximation of parametrized PDEs by reduced basis and neural networks, Journal of Computational Physics 416 (2020) 109550.

[30] S. Fresca, L. Dede', A. Manzoni, A comprehensive deep learning-based approach to reduced order modeling of nonlinear time-dependent parametrized PDEs, Journal of Scientific Computing 87 (2) (2021) 61.

[31] G. Kutyniok, P. Petersen, M. Raslan, R. Schneider, A theoretical analysis of deep neural networks and parametric PDEs, Constructive Approximation 55 (1) (2022) 73–125.

[32] T. O'Leary-Roseberry, U. Villa, P. Chen, O. Ghattas, Derivative-informed projected neural networks for high-dimensional parametric maps governed by PDEs, Computer Methods in Applied Mechanics and Engineering 388 (2022) 114199.

[33] N. Franco, A. Manzoni, P. Zunino, A deep learning approach to reduced order modelling of parameter dependent partial differential equations, Mathematics of Computation 92 (340) (2023) 483–524.

[34] T. O'Leary-Roseberry, P. Chen, U. Villa, O. Ghattas, Derivative-informed neural operator: an efficient framework for high-dimensional parametric derivative learning, Journal of Computational Physics 496 (2024) 112555.

[35] H. Zheng, Y. Chen, J. Han, Y. Yu, ReBaNO: Reduced basis neural operator mitigating generalization gaps and achieving discretization invariance, arXiv preprint arXiv:2509.09611 (2025).

[36] Y. Wang, G. Lin, Reduced-basis deep operator learning for parametric PDEs with independently varying boundary and source data, arXiv preprint arXiv:2511.18260 (2025).

[37] L. Demkowicz, J. Gopalakrishnan, The discontinuous Petrov-Galerkin method, Acta Numerica 34 (2025) 293–384.

[38] I. Babuška, Error-bounds for finite element method, Numer. Math. 16 (1971) 322–333.

[39] I. Babuška, K. Aziz, G. Fix, R. Kellogg, Survey lectures on the mathematical foundations of the finite element method, in: The Mathematical Foundations of the Finite Element Method with Applications to Partial Differential Equations, K. Aziz, ed., Academic Press, 1972.

[40] P. Cortés Castillo, W. Dahmen, J. Gopalakrishnan, DPG loss functions for learning parameter-to-solution maps by neural networks, arXiv preprint: http://arxiv.org/abs/2506.18773 (2025).

[41] G. Gantner, R. Stevenson, Applications of a space-time FOSLS formulation for parabolic PDEs, IMA J. Numer. Anal. (2023).

[42] A. Logg, K.-A. Mardal, G. Wells, Automated solution of differential equations by the finite element method: The FEniCS book, Vol. 84, Springer Science & Business Media, 2012.

[43] M. Mohri, A. Rostamizadeh, A. Talwalkar, Foundations of machine learning, MIT press, 2018.

[44] P. L. Bartlett, N. Harvey, C. Liaw, A. Mehrabian, Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks, Journal of Machine Learning Research 20 (63) (2019) 1–17.

[45] F. Bach, Learning Theory from First Principles, MIT Press, 2024.

[46] P. L. Bartlett, S. Mendelson, Rademacher and gaussian complexities: Risk bounds and structural results., Journal of Machine Learning Research 3 (2002) 463–482.

[47] M. Anthony, P. L. Bartlett, Neural Network Learning: Theoretical Foundations, Cambridge University Press, 1999.

[48] I. A. Baratta, J. P. Dean, J. S. Dokken, M. Habera, J. S. Hale, C. N. Richardson, M. E. Rognes, M. W. Scroggs, N. Sime, G. N. Wells, DOLFINx: the next generation FEniCS problem solving environment, preprint (2023).

[49] U. Villa, N. Petra, O. Ghattas, hIPPYlib: An extensible software framework for large-scale inverse problems governed by PDEs: Part I: Deterministic inversion and linearized Bayesian inference, ACM Transactions on Mathematical Software (TOMS) 47 (2) (2021) 1–34.

[50] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, Advances in neural information processing systems 32 (2019).

[51] N. Vyas, D. Morwani, R. Zhao, M. Kwun, I. Shapira, D. Brandfonbrener, L. Janson, S. Kakade, SOAP: Improving and stabilizing shampoo using Adam, arXiv preprint arXiv:2409.11321 (2024).

[52] L. Cao, T. O'Leary-Roseberry, O. Ghattas, Derivative-informed neural operator acceleration of geometric MCMC for infinite-dimensional bayesian inverse problems, Journal of Machine Learning Research 26 (78) (2025) 1–68.

[53] P. Si, P. Chen, Latent-EnSF: A Latent Ensemble Score Filter for High-Dimensional Data Assimilation with Sparse Observation Data, in: International Conference on Learning Representations (ICLR), 2025, published as a conference paper at ICLR 2025.

[54] P. Xiao, P. Si, P. Chen, LD-EnSF: Synergizing Latent Dynamics with Ensemble Score Filters for Fast Data Assimilation with Sparse Observations, arXiv preprint arXiv:2411.19305 (2024).

[55] P. Chen, J. O. Royset, Performance bounds for PDE-constrained optimization under uncertainty, SIAM Journal on Optimization 33 (3) (2023) 1828–1854.

[56] D. Luo, T. O'Leary-Roseberry, P. Chen, O. Ghattas, Efficient PDE-constrained optimization under high-dimensional uncertainty using derivative-informed neural operators, SIAM Journal on Scientific Computing 47 (4) (2025) C899–C931.

[57] B. Yao, D. Luo, L. Cao, N. Kovachki, T. O'Leary-Roseberry, O. Ghattas, Derivative-informed Fourier neural operator: Universal approximation and applications to pde-constrained optimization, arXiv preprint arXiv:2512.14086 (2025).

[58] K. Wu, T. O'Leary-Roseberry, P. Chen, O. Ghattas, Large-scale Bayesian optimal experimental design with derivative-informed projected neural network, Journal of Scientific Computing 95 (1) (2023) 30.

[59] J. Go, P. Chen, Accurate, scalable, and efficient Bayesian optimal experimental design with derivative-informed neural operators, Computer Methods in Applied Mechanics and Engineering 438 (2025) 117845.

[60] J. Go, P. Chen, Sequential infinite-dimensional Bayesian optimal experimental design with derivative-informed latent attention neural operator, Journal of Computational Physics 532 (2025) 113976.

[61] Y. Qiu, N. Bridges, P. Chen, Derivative-enhanced deep operator network, Advances in Neural Information Processing Systems 37 (2024) 20945–20981.

[62] D. Boffi, F. Brezzi, M. Fortin, Mixed Finite Element Methods and Applications, Vol. 44 of Springer Series in Computational Mathematics, Springer, 2013.

[63] I. Loshchilov, F. Hutter, Decoupled weight decay regularization, arXiv preprint arXiv:1711.05101 (2017).

# Appendix A   On the Riesz representation of $f \in \mathbb{U}'$

When admitting in (15) a general functional $f$ in $\mathbb{U}'$ we have to ensure that it does not interfere with the specific Neumann boundary condition on $\Gamma_N$. Indeed, in a weak formulation Neumann boundary conditions appear as functionals in $\mathbb{U}'$. We have therefore restricted $f$ in (15) to be of the form $f = f_2 + \operatorname{div} f_1$ where $f_2 \in L_2(\Omega)$ and $f_1 \in L_2(\Omega; \mathbb{R}^d)$ is "flux-free", in the sense that,

$$(\operatorname{div} f_1, v)_\Omega = -(f_1, \nabla v)_\Omega, \quad f(v) = (f_2, v)_\Omega + (\operatorname{div} f_1)(v), \quad \forall v \in \mathbb{U}. \tag{A.1}$$

To explain why such a requirement is reasonable we briefly describe how a (generally non-unique) decomposition (A.1) may come about.

To that end, we'll make use of the following Weyl-type decomposition of a general $F \in L_2 := L_2(\Omega; \mathbb{R}^d)$:

$$F = \nabla p \oplus \zeta \quad \text{for some} \quad p \in \mathbb{U}, \ \zeta \in \Sigma_0, \tag{A.2}$$

where we define

$$\Sigma := \{\eta \in H(\operatorname{div}; \Omega) : \eta \cdot n|_{\Gamma_N} = 0\}, \quad \Sigma_0 := \{\zeta \in \Sigma : \zeta \cdot n|_{\Gamma_N} = 0, \ \operatorname{div} \zeta = 0\}. \tag{A.3}$$

In particular, we then have, by our assumptions on $\zeta$ in (A.2),

$$\begin{aligned}
(F, \nabla v)_\Omega &= (\nabla p, \nabla v)_\Omega + (\zeta, \nabla v)_\Omega = (\nabla p, \nabla v)_\Omega \\
&= -(\operatorname{div} \nabla p)(v) + \langle \nabla p \cdot n, v \rangle_{\Gamma_N} \\
&= -(\operatorname{div} F)(v) + \langle \nabla p \cdot n, v \rangle_{\Gamma_N}, \quad v \in \mathbb{U}.
\end{aligned} \tag{A.4}$$

In other words, we can see the "trace-effect" when taking a weak divergence of an $L_2$-field. To suppress the flux term, consider the subspace

$$\mathbb{U}_{\Gamma_N} := \{v \in \mathbb{U} : \nabla v \cdot n|_{\Gamma_N} = 0\} \subset \mathbb{U}.$$

Thus, whenever, $F \in L_2$ has a decomposition (A.2) with $p \in \mathbb{U}_{\Gamma_N}$, then (A.4) says that

$$(F, \nabla v)_\Omega = -(\operatorname{div} F)(v), \quad v \in \mathbb{U}, \tag{A.5}$$

which explains the notion *flux-free* in (A.1).

To arrive at the desired decomposition of $f$ in (15), consider the *Riesz-lift* of $f$ with respect to the full $H^1$-inner product on $\mathbb{U}$: find $r_f \in \mathbb{U}$ such that

$$(\nabla r_f, \nabla v)_\Omega + (r_f, v)_\Omega = -f(v), \quad v \in \mathbb{U}. \tag{A.6}$$

which has a unique solution $r_f \in \mathbb{U}$. Since

$$f(v) = (\operatorname{div} \nabla r_f)(v) - (r_f, v)_\Omega - \langle \nabla r_f \cdot n, v \rangle_{\Gamma_N}, \quad v \in \mathbb{U}, \tag{A.7}$$

upon setting

$$f_2 := -r_f \in L_2(\Omega), \quad f_1 := \nabla r_f \in L_2(\Omega; \mathbb{R}^d),$$

we can see that the decomposition (A.7) does give rise to the desired property (A.1), provided that $r_f \in \mathbb{U}_{\Gamma_N}$, and how the implied relation $\langle f_1 \cdot n, v \rangle_{\Gamma_N} = 0$, $v \in \mathbb{U}$, is then to be understood.

# Appendix B    Stability and norm equivalence

## Appendix B.1    Stationary diffusion

**Lemma 3.** *Assume that $0 < \alpha \leq \mathfrak{p}(x) \leq \beta < \infty$ holds for fixed $\alpha, \beta$. Then the mapping $\mathcal{B}_{\mathfrak{p}}([\tau, v]) = \binom{\tau - \mathfrak{p} \nabla v}{-\operatorname{div} \tau}$ is an isomorphism from $\mathbb{H} := H_{0, \Gamma_N}(\operatorname{div}; \Omega) \times H^1_{0, \Gamma_D}(\Omega)$ onto $\mathbb{L}_2 := L_2(\Omega; \mathbb{R}^{d+1})$. In particular, the norm equivalence*

$$c\|[\tau, v]\|_{H(\operatorname{div};\Omega) \times H^1(\Omega)} \leq \|\mathcal{B}_{\mathfrak{p}}([\tau, v])\|_{\mathbb{L}_2} \leq C\|[\tau, v]\|_{H(\operatorname{div};\Omega) \times H^1(\Omega)}, \quad [\tau, v] \in \mathbb{H}, \tag{B.1}$$

*holds for constants $c, C$ depending only on $\alpha, \beta$ and the Poincaré constant $C_P$ in $\|v\|_{L_2(\Omega)} \leq C_P \|\nabla v\|_{L_2(\Omega)}$, $v \in \mathbb{U} := H^1_{0, \Gamma_D}(\Omega)$. In fact, denoting by $\kappa(\mathcal{B}_{\mathfrak{p}}) := \|\mathcal{B}_{\mathfrak{p}}\|_{\mathbb{H} \to \mathbb{L}_2} \|\mathcal{B}_{\mathfrak{p}}^{-1}\|_{\mathbb{L}_2 \to \mathbb{H}}$ the condition number of $\mathcal{B}_{\mathfrak{p}}$ we have $\kappa(\mathcal{B}_{\mathfrak{p}}) \leq C\beta/\alpha$ where $C$ is a constant depending only on $C_P$. In these terms*

$$c \geq c_1 \big( \sup_{\mathfrak{p} \in \mathfrak{P}} \kappa(\mathcal{B}_{\mathfrak{p}}) \big)^{-1} \quad while \quad C \leq (2 + \beta^2)^{1/2}, \tag{B.2}$$

*where $c_1$ depends only on $C_P$. In other words, the FOSLS bilinear form*

$$([\tau, v], [\eta, s])_{\mathbb{X}_{\mathfrak{p}}} := (\tau - \mathfrak{p} \nabla v, \eta - \mathfrak{p} \nabla s)_{L^2(\Omega)} + (\operatorname{div} \tau, \operatorname{div} \eta)_{L^2(\Omega)}, \tag{B.3}$$

*is an equivalent inner product on $\mathbb{H}$, inducing the norm $\|[\tau, v]\|_{\mathbb{X}_{\mathfrak{p}}}^2 := ([\tau, v], [\tau, v])_{\mathbb{X}_{\mathfrak{p}}} = \|\mathcal{B}_{\mathfrak{p}}[\tau, v]\|_{L^2(\Omega; \mathbb{R}^{d+1})}^2$.*

*Proof.* Proofs of this result can, in essence, be found in the literature, see [17, 16]. For the convenience of the reader, we present a short version with an eye to the parameter dependence.

**Upper bound.**  Simply use triangle inequality and Cauchy-Schwarz to obtain

$$\|\mathcal{B}_{\mathfrak{p}}([\tau, v])\|_{\mathbb{L}_2} \leq \|\tau\|_{L_2(\Omega; \mathbb{R}^d)} + \|\operatorname{div} \tau\|_{L_2(\Omega)} + \|\mathfrak{p} \nabla v\|_{L_2(\Omega)} \leq \sqrt{2}\|\tau\|_{H(\operatorname{div};\Omega)} + \beta\|v\|_{H^1(\Omega)}$$
$$\leq (2 + \beta^2)^{1/2}\|[\tau, v]\|_{H(\operatorname{div};\Omega) \times H^1(\Omega)},$$

confirming the second part in (B.2).

**Lower bound.** Using integration by parts (since $\tau \cdot n|_{\Gamma_N} = 0$ and $v|_{\Gamma_D} = 0$), and the fact that, by assumption on $\mathfrak{P}$, $\alpha \|\nabla v\|_{L^2}^2 \le (\mathfrak{p}\nabla v, \nabla v)$, we obtain

$$
\|\nabla v\|_{L^2}^2 \le \alpha^{-1}(\mathfrak{p}\nabla v, \nabla v) = -\alpha^{-1}\big((\tau - \mathfrak{p}\nabla v, \nabla v) + (\operatorname{div}\tau, v)\big)
$$

$$
\le \alpha^{-1}\Big\{\|\tau - \mathfrak{p}\nabla v\|_{L^2}\|\nabla v\|_{L^2} + \|\operatorname{div}\tau\|_{L^2}\|v\|_{L^2}\Big\}
$$

$$
\le \alpha^{-1}\|\mathcal{B}_{\mathfrak{p}}([\tau, v])\|_{\mathbb{L}_2}\|v\|_{H^1(\Omega)} \le \alpha^{-1}(1 + C_P^2)^{1/2}\|\mathcal{B}_{\mathfrak{p}}([\tau, v])\|_{\mathbb{L}_2}\|\nabla v\|_{L_2},
$$

so that, in particular,

$$
\|\nabla v\|_{L_2} \le \alpha^{-1}(1 + C_P^2)^{1/2}\|\mathcal{B}_{\mathfrak{p}}([\tau, v])\|_{\mathbb{L}_2}. \tag{B.4}
$$

Thus, by Poincaré's inequality,

$$
\|v\|_{H^1}^2 \le (1 + C_P^2)\|\nabla v\|_{L_2}^2 \le \frac{(1 + C_P^2)^2}{\alpha^2}\|\mathcal{B}_{\mathfrak{p}}([\tau, v])\|_{\mathbb{L}_2}^2. \tag{B.5}
$$

Moreover, again by (B.4) and using $\|\mathfrak{p}\nabla v\|_{L_2} \le \beta\|\nabla v\|_{L_2}$,

$$
\|\tau\|_{H(\operatorname{div})} \le \|\tau\|_{L_2} + \|\operatorname{div}\tau\|_{L_2} \le \|\tau - \mathfrak{p}\nabla v\|_{L_2} + \|\mathfrak{p}\nabla v\|_{L_2} + \|\operatorname{div}\tau\|_{L_2}
$$

$$
\le (\sqrt{2} + A)\|\mathcal{B}_{\mathfrak{p}}([\tau, v])\|_{\mathbb{L}_2}, \quad \text{where } A := \frac{\beta(1 + C_P^2)^{1/2}}{\alpha}. \tag{B.6}
$$

Combining (B.5) and (B.6) yields the lower bound with

$$
c = \left(\Big(\sqrt{2} + \frac{\beta(1 + C_P^2)^{1/2}}{\alpha}\Big)^2 + \frac{(1 + C_P^2)^2}{\alpha^2}\right)^{-1/2},
$$

from which (B.1) and the first inequality in (B.2) follow.

To complete the proof, note that (B.1) already implies boundedness and injectivity of $\mathcal{B}_{\mathfrak{p}}$, so that the remainder of the assertion follows from the Open Mapping Theorem, once we have confirmed surjectivity of $\mathcal{B}_{\mathfrak{p}}$. To that end, note that we have incidentally shown that the bilinear form $B(\cdot, \cdot) : \mathbb{H} \times \mathbb{H} \to \mathbb{R}$, defined by $B_{\mathfrak{p}}([\sigma, u], [\tau, v]) := (\mathcal{B}_{\mathfrak{p}}([\sigma, u]), \mathcal{B}_{\mathfrak{p}}([\tau, v]))_\Omega$ is coercive. Since $\mathcal{B}_{\mathfrak{p}}$ is bounded, $B(\cdot, \cdot)$ is also bounded (with constant $\|\mathcal{B}_{\mathfrak{p}}\|_{\mathbb{H} \to \mathbb{L}_2}^2 \le 2 + \beta^2$. Lax-Milgram shows that the normal equations: for any $(y_1, y_2) \in \mathbb{L}_2$, find $[\sigma, u] \in \mathbb{H}$ such that

$$
([\sigma, u], [\tau, v])_{\mathbb{X}_{\mathfrak{p}}} = ([y_1, y_2], \mathcal{B}_{\mathfrak{p}}([\tau, v]))_\Omega \quad \forall [\tau, v] \in \mathbb{H},
$$

is well-posed and stable. Now suppose $\mathcal{B}_{\mathfrak{p}}$ is not surjective. Since $\mathcal{B}_{\mathfrak{p}}$ is closed, its range is then a closed subspace of $\mathbb{L}_2$. Then there exists a $y = [y_1, y_2] \in \mathbb{L}_2 \setminus \{0\}$ such that $(\mathcal{B}_{\mathfrak{p}}[\tau, v], y)_\Omega = 0$ for all $[\tau, v] \in \mathbb{H}$. On the other hand, we know from the above observation that there exists a unique $0 \ne [\sigma_y, u_y] \in \mathbb{H}$ such that

$$
(\mathcal{B}_{\mathfrak{p}}[\sigma_y, u_y], \mathcal{B}_{\mathfrak{p}}[\eta, w])_\Omega = (y, \mathcal{B}_{\mathfrak{p}}([\eta, w])_\Omega, \quad [\eta, w] \in \mathbb{H}.
$$

By assumption, the right hand side vanishes for all $[\eta, w] \in \mathbb{H}$, while for $[\eta, w] = [\sigma_y, u_y]$, we have $(\mathcal{B}_{\mathfrak{p}}[\sigma_y, u_y], \mathcal{B}_{\mathfrak{p}}[\sigma_y, u_y])_\Omega = \|\mathcal{B}_{\mathfrak{p}}[\sigma_y, u_y]\|_{\mathbb{L}_2}^2 \ge c^2\|[\sigma_y, u_y]\|_{H(\operatorname{div};\Omega) \times H^1(\Omega)}^2 > 0$, a contradiction. This contradiction confirms surjectivity and completes the proof. $\qquad\square$

## Appendix B.2  Linear elasticity

Assume is $\Omega$ is Lipschitz and $\Gamma_D$ has positive measure. Define for $[\underline{\tau}, \underline{v}] \in \mathbb{H} := H_{0, \Gamma_N}(\operatorname{div}; \Omega; \mathbb{R}^{d \times d}) \times H^1_{0, \Gamma_D}(\Omega; \mathbb{R}^d)$

$$
\mathcal{B}_{\mathfrak{p}}[\underline{\tau}, \underline{v}] := \begin{pmatrix} \mathcal{C}_{\mathfrak{p}}^{-1/2}(\underline{\tau} - \mathcal{C}_{\mathfrak{p}}\underline{\underline{\varepsilon}}(\underline{v})) \\ -\underline{\operatorname{div}}\,\underline{\underline{\tau}} \end{pmatrix} \in L_2(\Omega; \mathbb{R}^{d \times d}) \times L_2(\Omega; \mathbb{R}^d) =: \mathbb{L}_2.
$$

As before, suppressing at times the reference to domain and range in $L_2$-, $H^1$-, and $H(\operatorname{div})$-norms, we introduce the elasticity FOSLS bilinear form

$$
([\underline{\tau}, \underline{v}], [\underline{\eta}, \underline{s}])_{\mathbb{X}_{\mathfrak{p}}} := (\mathcal{C}_{\mathfrak{p}}^{-1/2}(\underline{\tau} - \mathcal{C}_{\mathfrak{p}}\underline{\underline{\varepsilon}}(\underline{v})), \mathcal{C}_{\mathfrak{p}}^{-1/2}(\underline{\eta} - \mathcal{C}_{\mathfrak{p}}\underline{\underline{\varepsilon}}(\underline{s})))_\Omega + (\underline{\operatorname{div}}\,\underline{\underline{\tau}}, \underline{\operatorname{div}}\,\underline{\underline{\eta}})_{L_2(\Omega)},
$$

and set $\|[\underline{\tau}, \underline{v}]\|_{\mathbb{X}_{\mathfrak{p}}}^2 := ([\underline{\tau}, \underline{v}], [\underline{\tau}, \underline{v}])_{\mathbb{X}_{\mathfrak{p}}} = \|\mathcal{B}_{\mathfrak{p}}[\underline{\tau}, \underline{v}]\|_{\mathbb{L}_2}^2$.

**Lemma 4.** *Assume the stiffness tensor $\mathcal{C}_{\mathfrak{p}}$ is uniformly symmetric positive definite (SPD) with bounds independent of $\mathfrak{p} \in \mathfrak{P}$, i.e., there exist positive constants $c_0, c_1$ such that*

$$c_0 \, \underline{\varepsilon}(\underline{w}) : \underline{\varepsilon}(\underline{w}) \le \underline{\varepsilon}(\underline{w}) : \mathcal{C}_{\mathfrak{p}} \, \underline{\varepsilon}(\underline{w}) \le c_1 \, \underline{\varepsilon}(\underline{w}) : \underline{\varepsilon}(\underline{w}) \quad \forall \underline{w} \in H^1(\Omega; \mathbb{R}^d).$$

*Then for all $[\underline{\tau}, \underline{v}] \in H_{0,\Gamma_N}(\mathrm{div}\,; \Omega; \mathbb{R}^{d \times d}) \times H^1_{0,\Gamma_D}(\Omega; \mathbb{R}^d)$,*

$$c \left( \|\underline{\tau}\|^2_{H(\mathrm{div})} + \|\underline{v}\|^2_{H^1} \right) \le \|[\underline{\tau}, \underline{v}]\|^2_{\mathbb{X}_{\mathfrak{p}}} \le C \left( \|\underline{\tau}\|^2_{H(\mathrm{div})} + \|\underline{v}\|^2_{H^1} \right), \tag{B.7}$$

*with constants $c, C > 0$ depending only on $c_0, c_1$ and the domain. Moreover, $\mathcal{B}_{\mathfrak{p}}$ is an isomorphism from $\mathbb{H} = H_{0,\Gamma_N}(\mathrm{div}\,; \Omega; \mathbb{R}^{d \times d}) \times H^1_{0,\Gamma_D}(\Omega; \mathbb{R}^d)$ onto $\mathbb{L}_2 = L_2(\Omega; \mathbb{R}^{d \times d}) \times L_2(\Omega; \mathbb{R}^d)$.*

*Proof.* Let $c_0, c_1 > 0$ be the uniform SPD bounds for $\mathcal{C}_{\mathfrak{p}}$ and let $C_P = C_P(\Omega, \Gamma_D)$ and $C_K = C_K(\Omega, \Gamma_D)$ be the Poincaré and Korn constants so that $\|\underline{v}\|_{L^2} \le C_P \|\nabla \underline{v}\|_{L^2}$ and $\|\nabla \underline{v}\|_{L^2} \le C_K \|\underline{\varepsilon}(v)\|_{L^2}$ for all $\underline{v} \in H^1_{0,\Gamma_D}(\Omega; \mathbb{R}^d)$ (e.g., [62]). The following arguments are analogous to the diffusion case, beginning with (B.7).

**Upper bound.** Triangle- and Cauchy Schwarz inequalities provide as before

$$\|\mathcal{B}_{\mathfrak{p}}[\underline{\tau}, \underline{v}]\|_{\mathbb{L}_2} \le \|\mathcal{C}_{\mathfrak{p}}^{-1/2} \underline{\tau}\|_{\mathbb{L}_2} + \|\mathcal{C}_{\mathfrak{p}}^{1/2} \underline{\varepsilon}(\underline{v})\|_{\mathbb{L}_2} + \|\underline{\mathrm{div}}\, \underline{\tau}\|_{\mathbb{L}_2} \le c_0^{-1} \sqrt{2} \|\underline{\tau}\|_{H(\mathrm{div})} + c_1 \|\underline{\mathrm{grad}}\, \underline{v}\|_{\mathbb{L}_2}$$

$$\le (2c_0 + c_1)^{1/2} \|[\underline{\tau}, \underline{v}]\|_{H(\mathrm{div}) \times H^1},$$

confirming the upper bound in (B.7).

**Lower bound.** Along the same lines as in the previous section, one can prove

$$\|\underline{\varepsilon}(v)\|_{L_2} \lesssim \frac{C_P}{c_0^{1/2}} \|\mathcal{B}_{\mathfrak{p}}[\underline{\tau}, \underline{v}]\|_{\mathbb{L}_2},$$

(and likewise for $\|\underline{\mathrm{grad}}\,(v)\|_{L_2}$) with an absolute proportionality constant to conclude further

$$\|\underline{v}\|_{H^1} \lesssim \frac{C_K C_P}{c_0^{1/2}} \|\mathcal{B}_{\mathfrak{p}}[\underline{\tau}, \underline{v}]\|_{\mathbb{L}_2}, \quad \underline{v} \in \mathbb{U} := H^1_{0,\Gamma_D}(\Omega; \mathbb{R}^d),$$

as well as

$$\|\underline{\tau}\|_{H(\mathrm{div})} \lesssim c_1^{1/2} \Big( 1 + \frac{c_1^{1/2} C_K C_P}{c_0^{1/2}} \Big) \|\mathcal{B}_{\mathfrak{p}}[\underline{\tau}, \underline{v}]\|_{\mathbb{L}_2}.$$

Assembling these bounds confirms the lower bound as well with an analogous dependence of $c$ on $(c_1/c_0)$, describing the square of the condition number of $\mathcal{B}_{\mathfrak{p}}$.

While (B.7) so far establishes bijectivity of $\mathcal{B}_{\mathfrak{p}}$ as a mapping from $\mathbb{H}$ to its range in $\mathbb{L}_2$, the same argument as in the previous section shows that this range is indeed all of $\mathbb{L}_2$, which by the Inverse Mapping Theorem completes the proof. $\qquad\square$

# Appendix C  Error estimates for finite element approximations

## Appendix C.1  Proof of Theorem 1 for the stationary diffusion problem

*Proof.* Again, for the convenience of the reader, we recall some relevant facts that can be found, for instance, in [62] concerning the FE error estimate for the first-order system least-squares (FOSLS) formulation of the diffusion problem. Recall the fiber-residual loss from (30)

$$\mathcal{L}([\tilde{\sigma}^{\circ}, \tilde{u}^{\circ}]; \mathfrak{p}) := \|\tilde{\sigma}^{\circ} - (\mathfrak{p}\nabla\tilde{u}^{\circ} + \mathfrak{p}\nabla w - z + f_1)\|^2_{L^2(\Omega)} + \|\mathrm{div}\, \tilde{\sigma}^{\circ} + f_2\|^2_{L^2(\Omega)}.$$

Let $S := [\mathfrak{p}\nabla w - z + f_1,\, f_2]$ denote the source vector. Then

$$\mathcal{L}([\tilde{\sigma}^\circ, \tilde{u}^\circ]; \mathfrak{p}) = \|\mathcal{B}_{\mathfrak{p}}[\tilde{\sigma}^\circ, \tilde{u}^\circ] - S\|^2_{L^2(\Omega;\mathbb{R}^{d+1})} = \|\tilde{\sigma}^\circ - (\mathfrak{p}\nabla\tilde{u}^\circ + \mathfrak{p}\nabla w - z + f_1)\|^2_{L^2(\Omega)} + \|\operatorname{div}\tilde{\sigma}^\circ + f_2\|^2_{L^2(\Omega)}.$$

Let $[\sigma^\circ, u^\circ] = [\sigma^\circ(\mathfrak{p}), u^\circ(\mathfrak{p})]$ denote the exact fiber solution, which satisfies $\mathcal{B}_{\mathfrak{p}}[\sigma^\circ, u^\circ] = S$ so that $\mathcal{L}([\sigma^\circ, u^\circ]; \mathfrak{p}) = 0$ (hence is the unique minimizer of the loss). For conforming FE spaces $\Sigma_h \subset H_{0,\Gamma_N}(\operatorname{div};\Omega)$ and $\mathbb{U}_h \subset H^1_{0,\Gamma_D}(\Omega)$, the Galerkin FOSLS solution $[\sigma^\circ_h, u^\circ_h] = [\sigma^\circ_h(\mathfrak{p}), u^\circ_h(\mathfrak{p})] \in \Sigma_h \times \mathbb{U}_h$ minimizes for each $\mathfrak{p} \in \mathfrak{P}$, $\mathcal{L}(\cdot; \mathfrak{p})$ over $\Sigma_h \times \mathbb{U}_h$ and satisfies the corresponding normal equations in $\Sigma_h \times \mathbb{U}_h$.

Since $[\sigma^\circ, u^\circ]$ solves the normal equation, one also has for any $[\tilde{\sigma}^\circ, \tilde{u}^\circ] \in \mathbb{H}$

$$\mathcal{L}([\tilde{\sigma}^\circ, \tilde{u}^\circ]; \mathfrak{p}) = \|\mathcal{B}_{\mathfrak{p}}[\tilde{\sigma}^\circ, \tilde{u}^\circ] - \mathcal{B}_{\mathfrak{p}}[\sigma^\circ, u^\circ]\|^2_{\mathbb{L}_2} = \|\mathcal{B}_{\mathfrak{p}}[\tilde{\sigma}^\circ - \sigma^\circ, \tilde{u}^\circ - u^\circ]\|^2_{\mathbb{L}_2} = \|[\tilde{\sigma}^\circ - \sigma^\circ, \tilde{u}^\circ - u^\circ]\|^2_{\mathbb{X}_{\mathfrak{p}}}.$$

Since $[\sigma^\circ_h, u^\circ_h]$ solves the normal equations in $\Sigma_h \times \mathbb{U}_h$, Galerkin orthogonality with respect to the inner product $(\cdot, \cdot)_{\mathbb{X}_{\mathfrak{p}}}$ and the resulting best-approximation property yields

$$\|[\sigma^\circ - \sigma^\circ_h, u^\circ - u^\circ_h]\|_{\mathbb{X}_{\mathfrak{p}}} = \min_{[\tau_h, v_h] \in \Sigma_h \times \mathbb{U}_h} \|[\sigma^\circ - \tau_h, u^\circ - v_h]\|_{\mathbb{X}_{\mathfrak{p}}}.$$

In what follows, to simplify notation, we occasionally write $L_2$ without specifying the domain and range, when this is clear from the context. Let $\Pi^{\mathrm{RT}}_k : H(\operatorname{div};\Omega) \to \Sigma_h$ denote the canonical Raviart–Thomas interpolant, (see [62], section 2.5.2, page 109) and $I_m : H^1(\Omega) \to \mathbb{U}_h$ an $H^1$ quasi-interpolant. Choosing $[\tau_h, v_h] = [\Pi^{\mathrm{RT}}_k \sigma^\circ, I_m u^\circ]$ and expanding the $\mathbb{X}_{\mathfrak{p}}$-norm, yields

$$\|[\sigma^\circ - \sigma^\circ_h, u^\circ - u^\circ_h]\|_{\mathbb{X}_{\mathfrak{p}}} \le \|[\sigma^\circ - \Pi^{\mathrm{RT}}_k \sigma^\circ, u^\circ - I_m u^\circ]\|_{\mathbb{X}_{\mathfrak{p}}}$$
$$\le \|\sigma^\circ - \Pi^{\mathrm{RT}}_k \sigma^\circ\|_{L_2} + \beta \|\nabla(u^\circ - I_m u^\circ)\|_{L_2} + \|\operatorname{div}(\sigma^\circ - \Pi^{\mathrm{RT}}_k \sigma^\circ)\|_{L_2}.$$

Standard RT and $H^1$ interpolation estimates on shape-regular meshes (see, e.g., [62, Proposition 2.5.4]) give, assuming regularity-orders $s_\sigma, s_{\mathrm{div}} \ge 0$ and $s_u \ge 1$ for $\sigma^\circ, \operatorname{div}\sigma^\circ, u^\circ$, respectively,

$$\|\sigma^\circ - \Pi^{\mathrm{RT}}_k \sigma^\circ\|_{L_2} \lesssim h^{\min(k+1, s_\sigma)} \|\sigma^\circ\|_{H^{s_\sigma}(\Omega)},$$
$$\|\operatorname{div}(\sigma^\circ - \Pi^{\mathrm{RT}}_k \sigma^\circ)\|_{L_2} \lesssim h^{\min(k+1, s_{\mathrm{div}})} \|\operatorname{div}\sigma^\circ\|_{H^{s_{\mathrm{div}}}(\Omega)},$$
$$\|\nabla(u^\circ - I_m u^\circ)\|_{L_2} \lesssim h^{\min(m, s_u - 1)} \|u^\circ\|_{H^{s_u}(\Omega)}.$$

where, generally $s_{\mathrm{div}} = s_\sigma - 1$. Note that these estimates are therefore relevant when $s_\sigma - 1$ is at least $k+1$, which takes advantage of the particular properties of the Raviart-Thomas finite element spaces.

Combining these bounds, using $\|\cdot\|_{\mathbb{X}_{\mathfrak{p}}} \approx \|\cdot\|_{H(\operatorname{div})\times H^1}$, we obtain the a priori bound stated in Theorem 1:

$$\mathcal{L}([\sigma^\circ_h, u^\circ_h]; \mathfrak{p}) = \|[\sigma^\circ - \sigma^\circ_h, u^\circ - u^\circ_h]\|^2_{\mathbb{X}_{\mathfrak{p}}}$$
$$\lesssim h^{2\min(k+1, s_\sigma)} \|\sigma^\circ\|^2_{H^{s_\sigma}} + h^{2\min(k+1, s_{\mathrm{div}})} \|\operatorname{div}\sigma^\circ\|^2_{H^{s_{\mathrm{div}}}} + h^{2\min(m, s_u - 1)} \|u^\circ\|^2_{H^{s_u}}.$$

In particular, when $s_\sigma, s_{\mathrm{div}} \ge k+1$, and $s_u \ge m+1$, it reduces to $\mathcal{L}([\sigma^\circ_h, u^\circ_h]; \mathfrak{p}) \lesssim h^{2(k+1)} + h^{2m}$. Thus, choosing $m = k+1$ balances the contributions, yielding $\mathcal{L}([\sigma^\circ_h, u^\circ_h]; \mathfrak{p}) \lesssim h^{2(k+1)}$. □

**Remark 3** (approximate lifts consistency). *In practice, in the evaluation of the FE loss $\mathcal{L}([\sigma^\circ_h, u^\circ_h]; \mathfrak{p})$, we replace the variational lifts $(w, z)$ by their $CG_m$ approximations $(w_h, z_h)$, leading to the approximate FE loss $\mathcal{L}_h([\sigma^\circ_h, u^\circ_h]; \mathfrak{p})$. It is straightforward to show (by $(a+b)^2 \le 2a^2 + 2b^2$) that*

$$\mathcal{L}_h([\sigma^\circ_h, u^\circ_h]; \mathfrak{p}) \le 2\mathcal{L}([\sigma^\circ_h, u^\circ_h]; \mathfrak{p}) + 2\|\delta\|^2_{L^2} \ \text{and} \ \mathcal{L}([\sigma^\circ_h, u^\circ_h]; \mathfrak{p}) \le 2\mathcal{L}_h([\sigma^\circ_h, u^\circ_h]; \mathfrak{p}) + 2\|\delta\|^2_{L^2},$$

*where the error $\delta := \mathfrak{p}\nabla(w - w_h) - (z - z_h)$ satisfies $\|\delta\|^2_{L^2} \lesssim h^{2m}$ [62, Proposition 2.5.4]) as the lifts are harmonic. This confirms that the approximate lifts do not degrade the asymptotic accuracy with $m \ge k+1$.*

## Appendix C.2 Error estimates of FE approximations for linear elasticity

**Theorem 4.** *Assume $\Omega$ is Lipschitz and $\mathcal{C}_\mathfrak{p}$ is uniformly bounded with bounds independent of $\mathfrak{p} \in \mathfrak{P}$. Let $\Sigma_h = (\mathrm{RT}_k^\circ)^d$ and $\mathbb{U}_h = (\mathrm{CG}_m^\circ)^d$ on a shape-regular mesh of size $h$, with $k \geq 0$, $m \geq 1$. Let $[\underline{\underline{\sigma}}^\circ, \underline{u}^\circ]$ be the exact solution and $[\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]$ the Galerkin solution in $\Sigma_h \times \mathbb{U}_h$. Then, for each fixed $\mathfrak{p} \in \mathfrak{P}$, and $[\underline{\underline{\sigma}}^\circ(\mathfrak{p}), \underline{u}^\circ(\mathfrak{p})] = [\underline{\underline{\sigma}}^\circ, \underline{u}^\circ]$, $[\underline{\underline{\sigma}}_h^\circ(\mathfrak{p}), \underline{u}_h^\circ(\mathfrak{p})] = [\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]$, one has the error–residual equivalence*

$$\mathcal{L}([\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]; \mathfrak{p}) \eqsim \|\underline{\underline{\sigma}}^\circ - \underline{\underline{\sigma}}_h^\circ\|_{H(\mathrm{div}\,;\Omega;\mathbb{R}^{d\times d})}^2 + \|\underline{u}^\circ - \underline{u}_h^\circ\|_{H^1(\Omega;\mathbb{R}^d)}^2, \tag{C.1}$$

*with equivalence constants, depending only on the uniform SPD bounds for $\mathcal{C}_\mathfrak{p}$ and the domain.*

*Moreover, if $\underline{\underline{\sigma}}^\circ \in H^{s_\sigma}(\Omega;\mathbb{R}^{d\times d})$, $\underline{\mathrm{div}}\,\underline{\underline{\sigma}}^\circ \in H^{s_{\mathrm{div}}}(\Omega;\mathbb{R}^d)$ with $s_\sigma, s_{\mathrm{div}} \geq 0$, and $\underline{u}^\circ \in H^{s_u}(\Omega;\mathbb{R}^d)$ with $s_u \geq 1$, then*

$$\mathcal{L}([\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]; \mathfrak{p}) \lesssim h^{2\min(k+1,s_\sigma)} \|\underline{\underline{\sigma}}^\circ\|_{H^{s_\sigma}}^2 + h^{2\min(k+1,s_{\mathrm{div}})} \|\underline{\mathrm{div}}\,\underline{\underline{\sigma}}^\circ\|_{H^{s_{\mathrm{div}}}}^2 + h^{2\min(m,s_u-1)} \|\underline{u}^\circ\|_{H^{s_u}}^2. \tag{C.2}$$

*In particular, if $\underline{\underline{\sigma}}^\circ \in H^{k+1}$, $\underline{\mathrm{div}}\,\underline{\underline{\sigma}}^\circ \in H^{k+1}$, and $\underline{u}^\circ \in H^{m+1}$, then $\mathcal{L}([\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]; \mathfrak{p}) \lesssim h^{2(k+1)} + h^{2m}$, and the balanced choice $m = k + 1$ yields the optimal scaling $\mathcal{L}([\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]; \mathfrak{p}) \lesssim h^{2(k+1)}$.*

*Proof.* By the same orthogonal projection argument in $(\cdot, \cdot)_{\mathbb{X}_\mathfrak{p}}$ defined in Appendix B.2, the Galerkin solution $[\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]$ satisfies the best-approximation property

$$\||[\underline{\underline{\sigma}}^\circ - \underline{\underline{\sigma}}_h^\circ, \underline{u}^\circ - \underline{u}_h^\circ]\||_{\mathbb{X}_\mathfrak{P}} = \min_{[\underline{\underline{\tau}}_h, \underline{v}_h] \in \Sigma_h \times \mathbb{U}_h} \||[\underline{\underline{\sigma}}^\circ - \underline{\underline{\tau}}_h, \underline{u}^\circ - \underline{v}_h]\||_{\mathbb{X}_\mathfrak{p}}.$$

Choosing interpolation operators $\Pi_h$ for $H(\mathrm{div}\,)$-conforming tensor fields and $I_h$ for $H^1$-conforming vector fields with standard approximation properties, and using the stability lemma, we obtain the a priori bound (C.2) whenever $\Sigma_h$ affords order-$k + 1$ approximation in $L_2$ for tensors and their divergences, and $\mathbb{U}_h$ affords order-$m$ in $H^1$. This holds, for instance, for componentwise $\mathrm{RT}_k$ spaces for stresses and $\mathrm{CG}_m$ for displacements; see [62, Ch. 2]. Balancing $m = k + 1$ yields the optimal rate $\mathcal{L}([\underline{\underline{\sigma}}_h^\circ, \underline{u}_h^\circ]; \mathfrak{p}) \lesssim h^{2(k+1)}$. $\qquad\square$

# Appendix D  Experimental details

## Appendix D.1  Visualization of the sparsity pattern of weight

In Figure D.10, we mark by blue dots the entries of the weight $W_\mathfrak{p}$ whose magnitudes exceed $10^{-10}$. We observe that the nonzero entries are predominantly concentrated along and near the main diagonal. In addition to this banded structure, finer diagonal patterns offset from the main diagonal are visible. These arise because each finite element contributes a structured local stencil, resulting in repeated small dense sub-blocks in the global matrix.
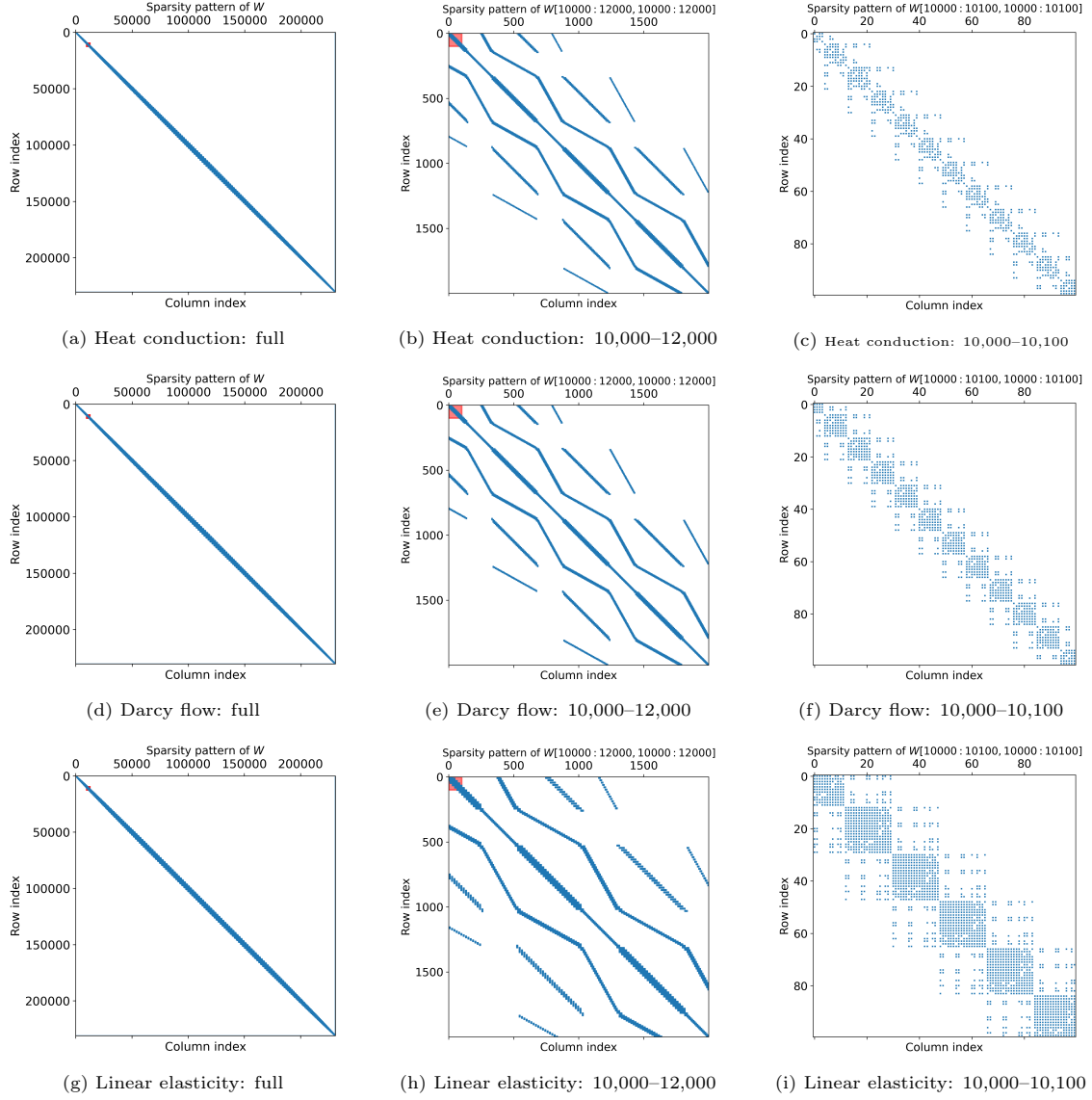
(a) Heat conduction: full  (b) Heat conduction: 10,000–12,000  (c) Heat conduction: 10,000–10,100

(d) Darcy flow: full  (e) Darcy flow: 10,000–12,000  (f) Darcy flow: 10,000–10,100

(g) Linear elasticity: full  (h) Linear elasticity: 10,000–12,000  (i) Linear elasticity: 10,000–10,100

Figure D.10: Sparsity patterns of the weights $W_{\mathfrak{p}}$ (threshold at $10^{-10}$) shown at three zoom levels: first column – full matrix $W_{\mathfrak{p}}$; second column – submatrix $W_{\mathfrak{p}}[10000 : 12000, 10000 : 12000]$ (highlighted by the red patch in the first column); third column – submatrix $W_{\mathfrak{p}}[10000 : 10100, 10000 : 10100]$ (highlighted by the red patch in the second column). Results are presented for three problem setups: heat conduction (top row), Darcy flow (middle row), and linear elasticity (bottom row).

41

## Appendix D.2  Visualization of reduced weight

In Figure D.11, we visualize the entries of a representative reduced weights $W_{\mathfrak{p}}^r$ for the three problem setups. The number of POD basis functions used is 128 for heat conduction, 512 for Darcy flow, and 512 for the elasticity problem.
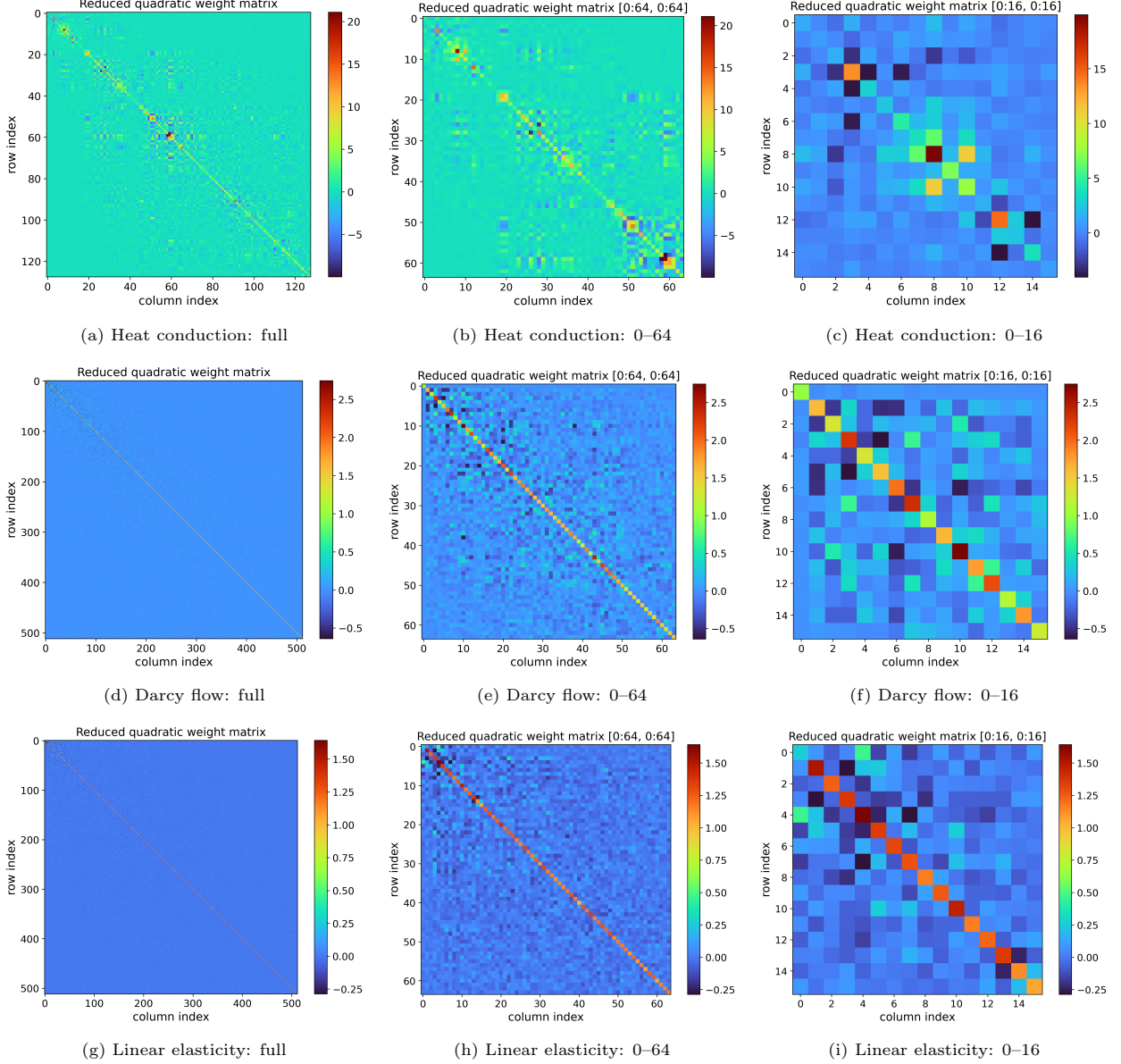


(a) Heat conduction: full      (b) Heat conduction: 0–64      (c) Heat conduction: 0–16

(d) Darcy flow: full      (e) Darcy flow: 0–64      (f) Darcy flow: 0–16

(g) Linear elasticity: full      (h) Linear elasticity: 0–64      (i) Linear elasticity: 0–16

Figure D.11: Reduced weights $W_{\mathfrak{p}}^r$ in (82) of heat conduction (top row), Darcy flow (middle row), and linear elasticity (bottom row) problems at three zoom levels, indicating dense matrix with dominant diagonal entries.

## Appendix D.3    Architecture and training details

We report the configurations of network architectures and training details of RBNO, PCA-Net and FNO in Tables D.3 to D.5 along with training time and inference time in Table D.6. For the PCA-Net and FNO, pointwise evaluations on grid points of high-fidelity solutions are used as training labels: 500 solutions for heat conduction, and 1000 for Darcy flow and linear elasticity. These numbers match those used to construct the POD basis in RBNO, which is the only stage where RBNO accesses high-fidelity solutions. All models are trained for 3000 iterations for diffusion and 6000 iterations for elasticity, with a batch size that fits in GPU memory. As in RBNO, training and validation losses are monitored at each iteration, and the model weights corresponding to the lowest validation loss are saved.

| Layer Type | Configuration | Output Shape (C, H, W) |
|---|---|---|
| **Stationary Diffusion** | | |
| Input | 1 channel image | (1, 129, 129) |
| Conv2d + act. | kernel=5, stride=2, out_ch=64 | (64, 63, 63) |
| Conv2d + act. | kernel=5, stride=2, out_ch=128 | (128, 30, 30) |
| Conv2d + act. | kernel=5, stride=3, out_ch=256 | (256, 9, 9) |
| Conv2d + act. | kernel=3, stride=1, out_ch=512 | (512, 7, 7) |
| Flatten | - | $(512{\times}7{\times}7)$ |
| Linear + act. | $512{\times}7{\times}7 \rightarrow 512$ | (512) |
| Linear | $512 \rightarrow$ output_dim | (output_dim) |
| # trainable parameters | 15.1 M | |
| **Linear Elasticity** | | |
| Input | 1 channel image | (1, 65, 129) |
| Conv2d + act. | kernel=5, stride=2, out_ch=64 | (64, 31, 63) |
| Conv2d + act. | kernel=5, stride=2, out_ch=128 | (128, 14, 30) |
| Conv2d + act. | kernel=5, stride=3, out_ch=256 | (256, 4, 9) |
| Conv2d + act. | kernel=3, stride=1, out_ch=512 | (512, 2, 7) |
| Flatten | - | $(512{\times}2{\times}7)$ |
| Linear + act. | $512{\times}2{\times}7 \rightarrow 512$ | (512) |
| Linear | $512 \rightarrow$ output_dim | (output_dim) |
| # trainable parameters | 6.1 M | |
| initialization | xavier_uniform | |
| activation | LeakyReLU (negative_slope=0.01) | |
| **Training details** | | |
| batch size | 1000 (diffusion), 500 (elasticity) | |
| optimizer | SOAP [51] (lr=5e-3, betas=(.95, .95), weight_decay=.01, precondition_frequency=5) | |
| lr scheduler | StepLR (step_size=50, gamma=0.95 (heat conduction), 0.9 (Darcy flow & elasticity)) | |

Table D.3: Architecture and training details for RBNO

|  | Dataset | | |
| --- | --- | --- | --- |
|  | Heat Conduction | Darcy Flow | Linear Elasticity |
| input dim | 20 | 512 | 512 |
| output dim | 128 | 512 | 512 |
| hidden dim | $(1024, 2048, 1024)$ | $(1024, 2048, 1024)$ | $(1024, 2048, 1024)$ |
| activation function | SELU | SELU | SELU |
| # trainable parameters | 4.4 M | 5.2 M | 5.2 M |
| **Training details** | | | |
| batch size | 500 | 1000 | 500 |
| AdamW [63] (lr, weight decay) | $(10^{-3}, 10^{-2})$ | $(10^{-3}, 10^{-2})$ | $(10^{-4}, 10^{-2})$ |
| StepLR (gamma, step size) | $(0.9, 50)$ | $(0.9, 50)$ | $(0.99, 100)$ |

Table D.4: Architecture and training details for PCA-Net

|  | Dataset | |
| --- | --- | --- |
|  | Stationary Diffusion | Linear Elasticity |
| number of modes | (32,32) | $(32, 32)$ |
| in channels | 3 | 3 |
| out channels | 3 | 6 |
| hidden channels | 128 | 128 |
| number of layers | 4 | 4 |
| lifting channel ratio | 2 | 2 |
| projection channel ratio | 2 | 2 |
| activation function | GELU | GELU |
| # trainable parameters | 35.9 M | 35.9 M |
| **Training details** | | |
| batch size | 10 | 10 |
| AdamW [63] (lr, weight decay) | $(10^{-3}, 10^{-2})$ | $(10^{-3}, 10^{-2})$ |
| StepLR (gamma, step size) | $(0.9, 50)$ | $(0.9, 50)$ |

Table D.5: Architecture and training details for FNO

| Problem | Method | Training time (min) | Inference time/sample (ms) |
| --- | --- | --- | --- |
| Heat Conduction | PCA-Net | 0.2 | 0.2 |
|  | FNO | 15 | 3.2 |
|  | RBNO (RB coef.) | 7 | 2.6 |
|  | RBNO (residual) | 15 | 2.6 |
|  | RBNO (both) | 19 | 2.6 |
| Darcy Flow | PCA-Net | 0.2 | 1.0 |
|  | FNO | 15 | 3.2 |
|  | RBNO (RB coef.) | 8 | 6.0 |
|  | RBNO (residual) | 14 | 6.0 |
|  | RBNO (both) | 19 | 6.0 |
| Linear Elasticity | PCA-Net | 1.8 | 1.0 |
|  | FNO | 16 | 1.6 |
|  | RBNO (RB coef.) | 10 | 6.0 |
|  | RBNO (residual) | 18 | 6.0 |
|  | RBNO (both) | 23 | 6.0 |

Table D.6: Comparisons of training and inference time for different methods.

## Appendix D.4  $H(\text{div}) \times H^1$ reconstruction from $\text{CG}_1 \times \text{CG}_1$ representations

Starting from a pointwise evaluation of the approximate solution $(\hat{\sigma}^\circ, \hat{u}^\circ)$ and its FE representation with $\text{CG}_1 \times \text{CG}_1$ elements, we construct an $H(\text{div}) \times H^1$-conforming approximation by projecting into the space $\mathbb{H}_h = \text{RT}_1^\circ \times \text{CG}_2^\circ$ in the $H(\text{div}) \times H^1$ norm, enforcing the appropriate homogeneous boundary conditions.

For the diffusion problems, we solve for $(\sigma^\circ, u^\circ) \in \mathbb{H}_h$:

$$(\sigma^\circ, \tau)_\Omega + (\text{div}\,\sigma^\circ, \text{div}\,\tau)_\Omega + (u^\circ, v)_\Omega + (\text{grad}\,u^\circ, \text{grad}\,v)_\Omega$$
$$= (\hat{\sigma}^\circ, \tau)_\Omega + (\text{div}\,\hat{\sigma}^\circ, \text{div}\,\tau)_\Omega + (\hat{u}^\circ, v)_\Omega + (\text{grad}\,\hat{u}^\circ, \text{grad}\,v)_\Omega, \quad (\tau, v) \in \mathbb{H}_h,$$

with $\sigma^\circ \cdot n = 0$ on $\Gamma_N$ and $u^\circ = 0$ on $\Gamma_D$.

For the linear elasticity problem, we solve for $(\underline{\underline{\sigma}}^\circ, \underline{u}^\circ) \in \mathbb{H}_h = (\text{RT}_1^\circ)^2 \times (\text{CG}_2^\circ)^2$:

$$(\underline{\text{div}}\,\underline{\underline{\sigma}}^\circ, \underline{\text{div}}\,\underline{\underline{\tau}}^\circ)_\Omega + (\underline{\underline{\sigma}}^\circ, \mathcal{C}_{\bar{\mathfrak{p}}}^{-1}\underline{\underline{\tau}}^\circ)_\Omega + (\underline{\underline{\varepsilon}}(\underline{u}^\circ), \mathcal{C}_{\bar{\mathfrak{p}}}\underline{\underline{\varepsilon}}(\underline{v}^\circ))_\Omega$$
$$= (\underline{\text{div}}\,\underline{\underline{\hat{\sigma}}}^\circ, \underline{\text{div}}\,\underline{\underline{\tau}}^\circ)_\Omega + (\underline{\underline{\hat{\sigma}}}^\circ, \mathcal{C}_{\bar{\mathfrak{p}}}^{-1}\underline{\underline{\tau}}^\circ)_\Omega + (\underline{\underline{\varepsilon}}(\underline{\hat{u}}^\circ), \mathcal{C}_{\bar{\mathfrak{p}}}\underline{\underline{\varepsilon}}(\underline{v}^\circ))_\Omega, \quad (\underline{\underline{\tau}}^\circ, \underline{v}^\circ) \in \mathbb{H}_h,$$

with $\underline{\underline{\sigma}}^\circ \cdot \underline{n} = (0,0)^\top$ on $\Gamma_N$ and $\underline{u}^\circ = (0,0)^\top$ on $\Gamma_D$.

We then compare the residual loss of the original FE solution with $\text{RT}_1 \times \text{CG}_2$ elements, its $\text{CG}_1 \times \text{CG}_1$ representation, and the reconstructed $\text{RT}_1 \times \text{CG}_2$ solution obtained via this projection. We also compare their relative errors in $\mathbb{L} = L_2 \times L_2$-norm, $L_2 \times H^1$-norm, and $\mathbb{H} = H(\text{div}) \times H^1$-norm. We report the corresponding results in Table D.7.

| Metric | Original $\text{RT}_1 \times \text{CG}_2$ | $\text{CG}_1 \times \text{CG}_1$ | Recons. $\text{RT}_1 \times \text{CG}_2$ |
|---|---|---|---|
| **Heat Conduction** | | | |
| $\|\sigma^\circ - (\mathfrak{p}\nabla u^\circ + \mathfrak{p}\nabla w - \mathbf{z} + \mathbf{f}_1)\|_{L^2}^2$ | $1.12 \times 10^{-4}$ | $8.26 \times 10^{-3}$ | $8.26 \times 10^{-3}$ |
| $\|\text{div}\,\sigma^\circ + f_2\|_{L^2}^2$ | $7.78 \times 10^{-10}$ | $2.07 \times 10^1$ | $2.07 \times 10^1$ |
| Relative error (L2-L2) | - | $1.71 \times 10^{-1}$ | $1.71 \times 10^{-1}$ |
| Relative error (L2-H1) | - | $1.39 \times 10^{-1}$ | $1.39 \times 10^{-1}$ |
| Relative error (Hdiv-H1) | - | 3.81 | 3.81 |
| **Darcy Flow** | | | |
| $\|\sigma^\circ - (\mathfrak{p}\nabla u^\circ + \mathfrak{p}\nabla w - \mathbf{z} + \mathbf{f}_1)\|_{L^2}^2$ | $6.13 \times 10^{-6}$ | $1.19 \times 10^{-3}$ | $1.19 \times 10^{-3}$ |
| $\|\text{div}\,\sigma^\circ + f_2\|_{L^2}^2$ | $4.00 \times 10^{-3}$ | 2.80 | 2.80 |
| Relative error (L2-L2) | - | $2.12 \times 10^{-3}$ | $2.12 \times 10^{-3}$ |
| Relative error (L2-H1) | - | $1.05 \times 10^{-2}$ | $1.05 \times 10^{-2}$ |
| Relative error (Hdiv-H1) | - | $1.40 \times 10^{-1}$ | $1.40 \times 10^{-1}$ |
| **Linear Elasticity** | | | |
| $\|\mathcal{C}^{-1/2}(\underline{\underline{\sigma}}^\circ + \underline{\underline{z}}) - \mathcal{C}^{1/2}(\underline{\underline{\varepsilon}}(\underline{u}^\circ + \underline{w}))\|_{L^2}^2$ | $7.75 \times 10^{-4}$ | $6.44 \times 10^{-3}$ | $7.48 \times 10^{-3}$ |
| $\|\underline{\text{div}}\,\underline{\underline{\sigma}}^\circ + \underline{f}\|_{L^2}^2$ | $6.77 \times 10^{-6}$ | $2.88 \times 10^1$ | $2.88 \times 10^1$ |
| Relative error (L2-L2) | - | $8.99 \times 10^{-3}$ | $8.49 \times 10^{-3}$ |
| Relative error (L2-H1) | - | $4.55 \times 10^{-2}$ | $4.58 \times 10^{-2}$ |
| Relative error (Hdiv-H1) | - | 2.99 | 2.99 |

Table D.7: Reconstruction quality for one test sample on the two residual loss terms and relative errors for the original $\text{RT}_1 \times \text{CG}_2$ solution, the intermediate $\text{CG}_1 \times \text{CG}_1$ representation, and the reconstructed $\text{RT}_1 \times \text{CG}_2$ from $\text{CG}_1 \times \text{CG}_1$ representation.

Across all three problems, we observe that the $\text{CG}_1 \times \text{CG}_1$ representation violates the divergence-related loss terms (second row in Table D.7 for each problem). Projecting this representation into the $\text{RT}_1 \times \text{CG}_2$ space does not reduce the residual loss or the errors in $\mathbb{H}$-norm (rigth column), even though the corresponding vectors of degrees of freedom remain close in the Euclidean sense, as visualized in Figure D.12. We have also tested DOLFINx's `interpolate` function and mesh refinement for transferring $\text{CG}_1 \times \text{CG}_1$ representations into $\text{RT}_1 \times \text{CG}_2$; the qualitative conclusions are unchanged.

(a) DoFs of $\sigma^\circ$ (heat conduction)

(b) DoFs of $u^\circ$ (heat conduction)

(c) DoFs of $\sigma^\circ$ (Darcy flow)

(d) DoFs of $u^\circ$ (Darcy flow)

(e) DoFs of $\sigma^\circ$ (Linear Elasticity)
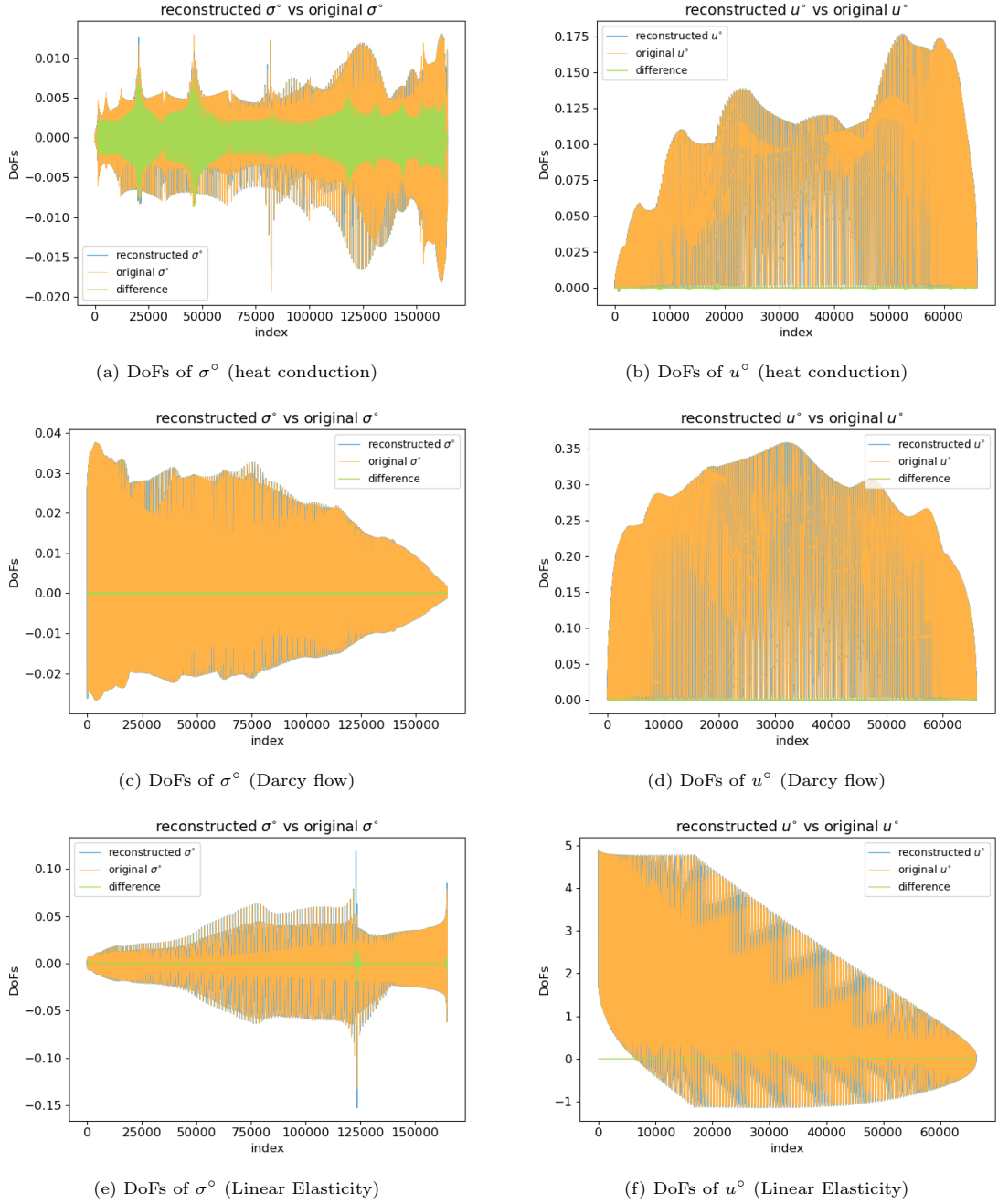
(f) DoFs of $u^\circ$ (Linear Elasticity)

Figure D.12: Comparison of DOFs and their reconstructed versions for the diffusion and elasticity problems.

## Appendix D.5 Additional results (expansion of Table 2)

In addition to the results shown in Table 2, we consider training with an MSE loss that matches the neural network predictions of the RB coefficients to their corresponding labels; this approach is denoted as RBNO (RB coef.). We also investigate whether adding this loss to the original RBNO training loss (the residual loss) during training can further reduce the error and loss over the test samples; this variant is denoted as RBNO (both).

For PCA-Net and FNO, the networks are trained to match pointwise evaluations at grid points to the

corresponding labels and output pointwise approximations at test time. To compare these approximations with the high-fidelity solution, we propose two approaches. First, we construct a $CG_1 \times CG_1$ representation of the approximation by mapping the pointwise evaluations to the degrees of freedom of the $CG_1 \times CG_1$ function. The first comparison method then evaluates this approximation against the $CG_1 \times CG_1$ representation of the high-fidelity solution (obtained by evaluating the solution on grid points and mapping to $CG_1 \times CG_1$ degrees of freedom). In the second method, we project the $CG_1 \times CG_1$ approximation into the $RT_1 \times CG_2$ space in the $\mathbb{H} = H(\text{div}) \times H^1$-norm and compare it with the high-fidelity solution (in the $RT_1 \times CG_2$ FE space). We indicate which reference function space is used, $[CG_1 \times CG_1]$ or $[RT_1 \times CG_2]$, in the results. The results for PCA-Net and FNO in table Table 2 are now denoted as PCA-Net $[CG_1 \times CG_1]$ and FNO $[CG_1 \times CG_1]$. The residual losses are evaluated using the corresponding FE representations of the approximations. In addition, we train PCA-Net to match the RB coefficients, rather than the PCA coefficients computed from solution evaluations at grid points. In this variant, we use the same basis as in RBNO to approximate the solution; this approach is denoted as PCA-Net (RB coef.).

We can see that RBNO with residual loss provides the best overall performance in terms of approximation accuracy and residual loss across all three problems. Adding the RB-coefficient MSE loss yields slight improvements for the diffusion problem, but in the elasticity problem, it even degrades performance. Standard data-driven models based on pointwise evaluations (FNO, PCA-Net) perform reasonably in $\mathbb{L}$-norm, but struggle with $\mathbb{H}$-norm and physics consistency (indicated by residual loss). While the PCA-Net trained on RB coefficient data produces a slightly larger $\mathbb{L}$-norm error than its pointwise counterpart (except for the heat conduction problem, where it is even lower), it generally achieves lower residual loss and $\mathbb{H}$-norm error (except for the $\mathbb{H}$-norm error in the linear elasticity problem); nevertheless, it is still outperformed by RBNO.

| Method | Relative error in $\mathbb{L}$-norm | Relative error in $\mathbb{H}$-norm | Residual loss |
|---|---|---|---|
| **Heat Conduction** | | | |
| PCA-Net $[CG_1 \times CG_1]$ | $1.27 \times 10^{-1}$ $(2.52 \times 10^{-2})$ | $3.32 \times 10^{-1}$ $(5.75 \times 10^{-2})$ | $1.84 \times 10^1$ $(3.01)$ |
| PCA-Net $[RT_1 \times CG_2]$ | $2.09 \times 10^{-1}$ $(1.99 \times 10^{-2})$ | $3.58$ $(2.95 \times 10^{-1})$ | $1.84 \times 10^1$ $(3.01)$ |
| PCA-Net (RB coef.) | $1.23 \times 10^{-1}$ $(2.52 \times 10^{-2})$ | $8.39 \times 10^{-2}$ $(1.65 \times 10^{-2})$ | $3.38 \times 10^{-2}$ $(2.66 \times 10^{-2})$ |
| FNO $[CG_1 \times CG_1]$ | $\mathbf{1.82 \times 10^{-2}}$ $(\mathbf{1.66 \times 10^{-3}})$ | $1.46 \times 10^{-1}$ $(1.53 \times 10^{-2})$ | $1.90 \times 10^1$ $(4.12)$ |
| FNO $[RT_1 \times CG_2]$ | $1.70 \times 10^{-1}$ $(1.57 \times 10^{-2})$ | $3.64$ $(3.98 \times 10^{-1})$ | $1.90 \times 10^1$ $(4.12)$ |
| RBNO (RB coef.) | $3.25 \times 10^{-2}$ $(1.98 \times 10^{-2})$ | $2.01 \times 10^{-2}$ $(1.25 \times 10^{-2})$ | $3.27 \times 10^{-3}$ $(2.50 \times 10^{-2})$ |
| RBNO (residual) | $2.43 \times 10^{-2}$ $(9.35 \times 10^{-3})$ | $1.86 \times 10^{-2}$ $(6.36 \times 10^{-3})$ | $5.33 \times 10^{-4}$ $(1.44 \times 10^{-3})$ |
| RBNO (both) | $2.25 \times 10^{-2}$ $(7.69 \times 10^{-3})$ | $\mathbf{1.71 \times 10^{-2}}$ $(\mathbf{5.38 \times 10^{-3}})$ | $\mathbf{4.55 \times 10^{-4}}$ $(\mathbf{2.67 \times 10^{-4}})$ |
| **Darcy Flow** | | | |
| PCA-Net $[CG_1 \times CG_1]$ | $1.79 \times 10^{-1}$ $(7.53 \times 10^{-2})$ | $1.13 \times 10^{-1}$ $(2.35 \times 10^{-2})$ | $2.59$ $(2.03)$ |
| PCA-Net $[RT_1 \times CG_2]$ | $1.79 \times 10^{-1}$ $(7.53 \times 10^{-2})$ | $1.25 \times 10^{-1}$ $(1.31 \times 10^{-2})$ | $2.59$ $(2.03)$ |
| PCA-Net (RB coef.) | $2.63 \times 10^{-1}$ $(1.10 \times 10^{-1})$ | $5.64 \times 10^{-2}$ $(2.01 \times 10^{-2})$ | $6.20 \times 10^{-1}$ $(9.40 \times 10^{-1})$ |
| FNO $[CG_1 \times CG_1]$ | $3.30 \times 10^{-2}$ $(1.18 \times 10^{-2})$ | $1.60 \times 10^{-1}$ $(3.81 \times 10^{-2})$ | $5.10$ $(2.34)$ |
| FNO $[RT_1 \times CG_2]$ | $\mathbf{3.25 \times 10^{-2}}$ $(\mathbf{1.18 \times 10^{-2}})$ | $1.89 \times 10^{-1}$ $(3.55 \times 10^{-2})$ | $5.10$ $(2.34)$ |
| RBNO (RB coef.) | $8.33 \times 10^{-2}$ $(3.39 \times 10^{-2})$ | $1.65 \times 10^{-2}$ $(5.59 \times 10^{-3})$ | $5.08 \times 10^{-2}$ $(6.36 \times 10^{-2})$ |
| RBNO (residual) | $6.39 \times 10^{-2}$ $(2.26 \times 10^{-2})$ | $1.23 \times 10^{-2}$ $(3.64 \times 10^{-3})$ | $2.72 \times 10^{-2}$ $(2.61 \times 10^{-2})$ |
| RBNO (both) | $6.11 \times 10^{-2}$ $(2.01 \times 10^{-2})$ | $\mathbf{1.22 \times 10^{-2}}$ $(\mathbf{3.72 \times 10^{-3}})$ | $\mathbf{2.66 \times 10^{-2}}$ $(\mathbf{2.38 \times 10^{-2}})$ |
| **Linear Elasticity** | | | |
| PCA-Net $[CG_1^2 \times CG_1^2]$ | $1.34 \times 10^{-1}$ $(1.03 \times 10^{-1})$ | $3.27 \times 10^{-1}$ $(1.14 \times 10^{-1})$ | $2.46 \times 10^1$ $(6.11)$ |
| PCA-Net $[RT_1^2 \times CG_2^2]$ | $1.34 \times 10^{-1}$ $(1.03 \times 10^{-1})$ | $2.62$ $(3.27 \times 10^{-1})$ | $2.46 \times 10^1$ $(6.11)$ |
| PCA-Net (RB coef.) | $2.90 \times 10^{-1}$ $(3.87 \times 10^{-2})$ | $6.12 \times 10^{-1}$ $(6.93 \times 10^{-2})$ | $1.51$ $(4.64 \times 10^{-1})$ |
| FNO $[CG_1^2 \times CG_1^2]$ | $\mathbf{1.17 \times 10^{-2}}$ $(\mathbf{3.34 \times 10^{-3}})$ | $4.30 \times 10^{-1}$ $(6.11 \times 10^{-2})$ | $2.82 \times 10^1$ $(6.45)$ |
| FNO $[RT_1^2 \times CG_2^2]$ | $2.22 \times 10^{-2}$ $(7.08 \times 10^{-3})$ | $2.78$ $(3.08 \times 10^{-1})$ | $2.81 \times 10^1$ $(6.45)$ |
| RBNO (RB coef.) | $2.95 \times 10^{-2}$ $(6.56 \times 10^{-3})$ | $5.32 \times 10^{-2}$ $(8.93 \times 10^{-3})$ | $1.02 \times 10^{-2}$ $(3.62 \times 10^{-3})$ |
| RBNO (residual) | $2.49 \times 10^{-2}$ $(5.89 \times 10^{-3})$ | $\mathbf{3.91 \times 10^{-2}}$ $(\mathbf{8.68 \times 10^{-3}})$ | $\mathbf{4.87 \times 10^{-3}}$ $(\mathbf{3.56 \times 10^{-3}})$ |
| RBNO (both) | $3.43 \times 10^{-2}$ $(8.74 \times 10^{-3})$ | $4.86 \times 10^{-2}$ $(1.69 \times 10^{-2})$ | $6.62 \times 10^{-3}$ $(2.01 \times 10^{-2})$ |

Table D.8: Comparison on the mean (and standard deviation, estimated over 500 test samples) of the relative errors in $\mathbb{L}_2 = L_2 \times L_2$-norm and $\mathbb{H} = H(\text{div}) \times H^1$-norm, and residual loss for PCA-Net, FNO and RBNO.

## Appendix D.6   Visualization of reference-prediction-difference

We visualize the reference–prediction difference for a random test sample in Figure D.13 (heat conduction), Figure D.14 (Darcy flow), and Figures D.15 and D.16 (elasticity). All functions are evaluated on a uniform grid, so the figures should be interpreted as showing the $\mathrm{CG}_1 \times \mathrm{CG}_1$ representation of the solutions.
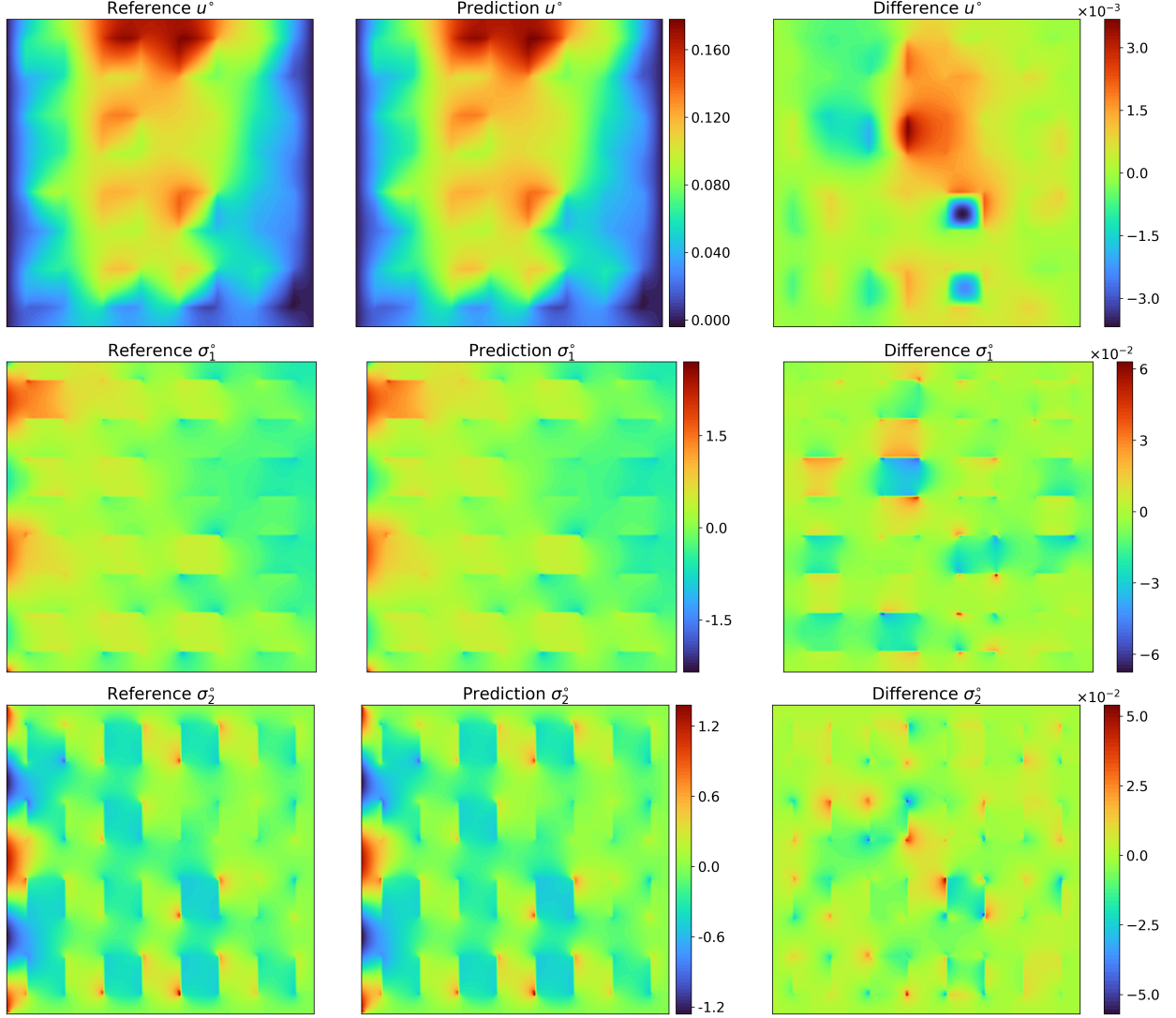


Figure D.13: Reference-prediction-difference (left-middle-right) of $u^\circ$ (top) and $\sigma^\circ = (\sigma_1^\circ, \sigma_2^\circ)$ (middle, bottom) [heat conduction].
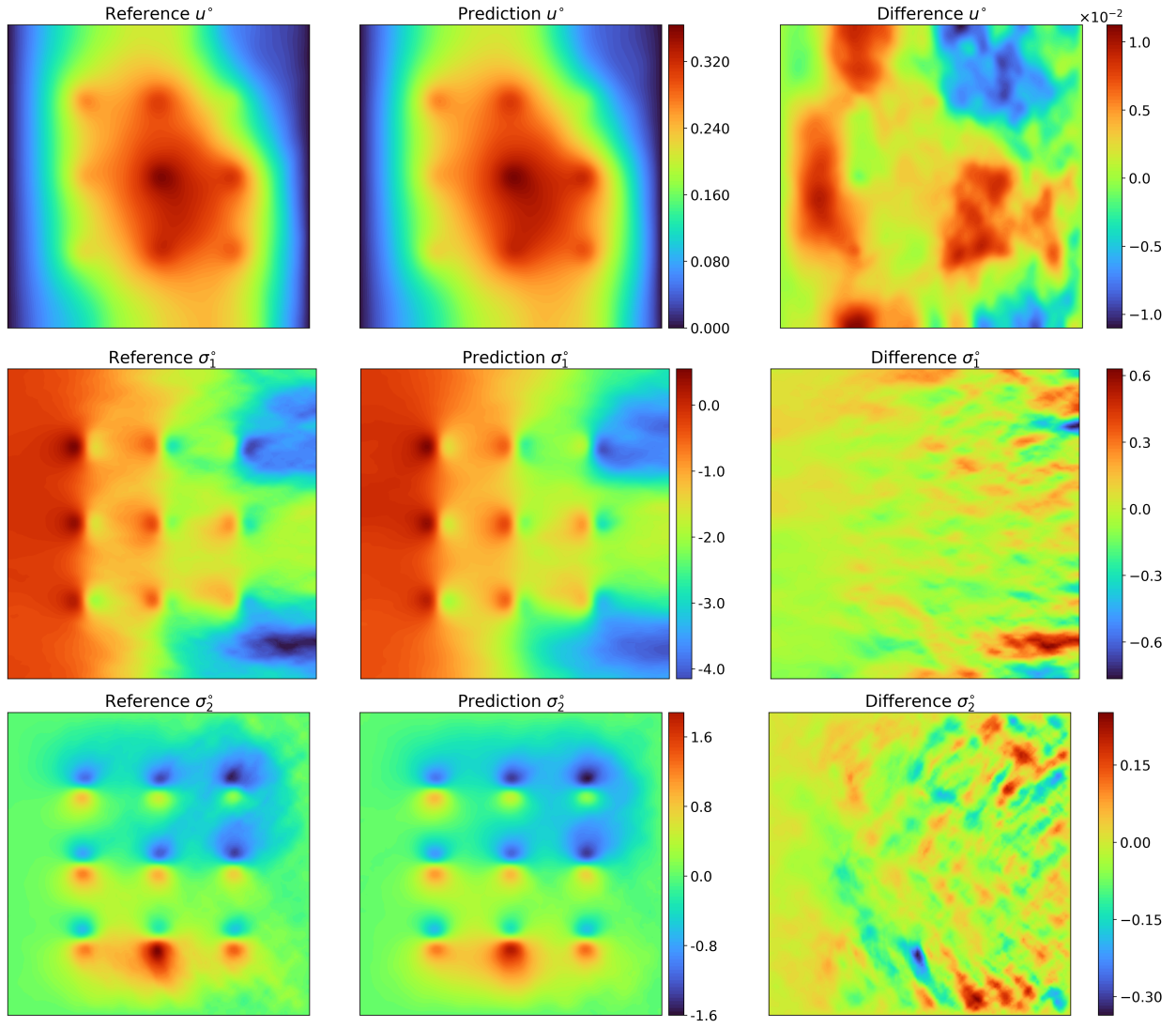
Figure D.14: Reference-approximation difference (left-middle-right) of $u^\circ$ (top) and $\sigma^\circ = (\sigma_1^\circ, \sigma_2^\circ)$ (middle, bottom) [Darcy flow].
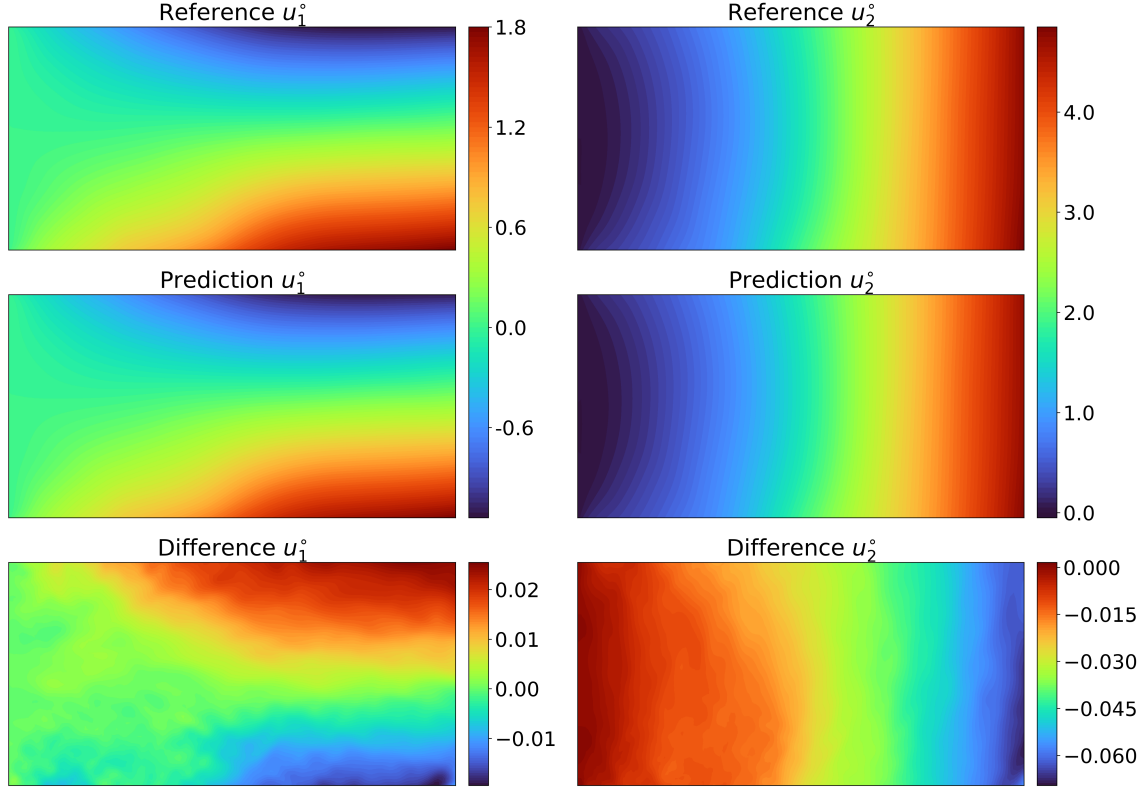
Figure D.15: Reference-prediction-difference (top-middle-bottom) of $u_1$ and $u_2$ (left and right) [linear elasticity].
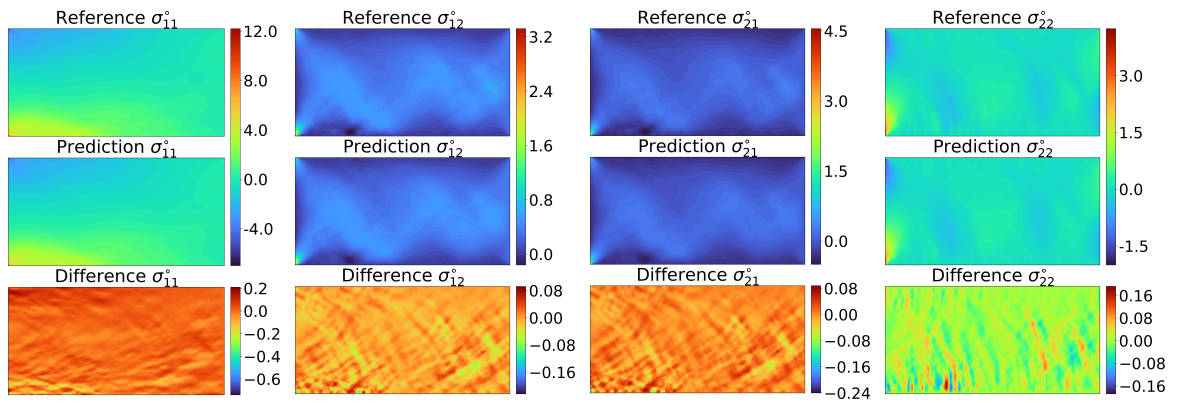


Figure D.16: Reference-prediction-difference (top-middle-bottom) of $\sigma_{11}, \sigma_{12}, \sigma_{21}$ and $\sigma_{22}$ (from left to right) [linear elasticity].