

---

# Provable In-Context Learning of Linear Systems and Linear Elliptic PDEs with Transformers

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Foundation models for natural language processing, empowered by the transformer  
2 architecture, exhibit remarkable *in-context learning* (ICL) capabilities: pre-trained  
3 models can adapt to a downstream task by only conditioning on few-shot prompts  
4 without updating the weights of the models. Recently, transformer-based founda-  
5 tion models also emerged as universal tools for solving scientific problems,  
6 including especially partial differential equations (PDEs). However, the theoretical  
7 underpinnings of ICL-capabilities of these models still remain elusive. This work  
8 develops rigorous error analysis for transformer-based ICL of the solution operators  
9 associated to a family of linear elliptic PDEs. Specifically, we show that a linear  
10 transformer defined by a linear self-attention layer can provably learn in-context to  
11 invert linear systems arising from the spatial discretization of the PDEs. We derive  
12 theoretical scaling laws for the proposed linear transformers in terms of the size of  
13 the spatial discretization, the number of training tasks, the lengths of prompts used  
14 during training and inference, under both the in-domain generalization setting and  
15 various settings of distribution shifts. Empirically, we validate the ICL-capabilities  
16 of transformers through extensive numerical experiments.

## 17 1 Introduction

18 Foundation models (FMs) for natural language processing (NLP), exemplified by ChatGPT Achiam  
19 et al. [2023], have demonstrated unprecedented power in text generation tasks. From an architectural  
20 perspective, the main novelty of these models is the use of transformer-based neural networks Vaswani  
21 et al. [2017], which are distinguished from feedforward neural networks by their self-attention layers.  
22 Those transformer-based FMs, pre-trained on a broad range of tasks with large amounts of data,  
23 exhibit remarkable transferability to diverse downstream tasks with limited data Brown et al. [2020].  
24 The success of of foundation models for NLP has recently sparked a large amount of work on building  
25 FMs in domain-specific scientific fields Batatia et al. [2023], Celaj et al. [2023], Méndez-Lucio et al.  
26 [2022]. Specifically, there is growing interest within the community of Scientific Machine Learning  
27 (SciML) in building scientific foundation models (SciFMs) to solve complex partial differential  
28 equations (PDEs) Subramanian et al. [2024], McCabe et al. [2023], Ye et al. [2024], Yang et al.  
29 [2023], Sun et al. [2024].

30 Traditional deep learning approaches for PDEs such as Physics-Informed Neural Networks Raissi  
31 et al. [2019] for learning solutions and neural operators Lu et al. [2019], Li et al. [2020] for learning  
32 solution operators need to be retrained from scratch for a different set of coefficients or different PDE  
33 systems. Instead, these SciFMs for PDEs, once pre-trained on large datasets of coefficients-solution  
34 pairs from multiple PDE systems, can be adapted to solving new PDE systems without training  
35 the model from scratch. Even more surprisingly, transformer-based FMs have demonstrated their  
36 **in-context learning (ICL)** capability in Achiam et al. [2023], Bubeck et al. [2023], Kirsch et al.

37 [2022] and in SciML Yang et al. [2023], Chen et al. [2024a], Yang and Osher [2024]: when given a  
38 prompt consisting of examples from a new learning task and a query, they are able to make correct  
39 predictions without updating their parameters. While the emergence of ICL has been deemed a  
40 paradigm shift in transformer-based FMs, its theoretical understandings remain underdeveloped.

41 The goal of this paper is to investigate the ICL capability of transformers for solving a class of  
42 linear elliptic PDEs and the associated linear systems. We are particularly interested in developing  
43 **neural scaling laws** that quantify the prediction risk of transformers as a function of the size of the  
44 training data, the model size, and other key parameters. Additionally, we aim to quantify the error  
45 incurred by **distribution shifts** between tasks and data used in pre-training and those in adaptation.  
46 As distribution shifts have been identified in Subramanian et al. [2024], McCabe et al. [2023], Ye  
47 et al. [2024], Yang et al. [2023] as a significant hurdle in the generalization capability of SciFMs, it is  
48 crucial to develop a rigorous theory of out-of-distribution generalization for SciFMs.

## 49 1.1 Main contributions.

50 We highlight our main contributions as follows:

- 51 • We formalize a framework for learning the solution operators of linear elliptic PDEs in-  
52 context. This is based on (1) reducing the infinite dimensional PDE problem into a problem  
53 of solving a finite dimensional linear system arising from spatial discretization of the PDE  
54 and (2) learning to invert the finite dimensional linear system in-context.
- 55 • We adopt transformers defined by single linear self-attention layers for ICL of the lin-  
56 ear systems and establish a non-asymptotic generalization bound of ICL in terms of the  
57 discretization size, the number of pre-training tasks, and the lengths of prompts used in  
58 pre-training and downstream tasks; see Theorem 1. This bound further enables us to prove  
59 an  $H^1$ -error bound for learning the solution of PDEs; see Theorem 2.
- 60 • We examine the prediction risk error that arises due to shifts in downstream task and  
61 covariate distributions. Specifically, we introduce a novel concept of task diversity and  
62 demonstrate that pre-trained transformers can generalize to out-of-distribution settings when  
63 the pre-training task distribution is diverse; see Theorem 3. Additionally, we provide several  
64 sufficient conditions under which task diversity holds; see Theorem 4.
- 65 • We demonstrate the ICL ability of linear transformers through several numerical experiments.

## 66 1.2 Related work

67 **ICL and FMs for PDE.** Several transformer-based FMs for solving PDEs have been developed  
68 in Subramanian et al. [2024], McCabe et al. [2023], Ye et al. [2024], Sun et al. [2024] where the  
69 pre-trained transformers are adapted to downstream tasks with fine-tuning on additional datasets. The  
70 work Yang et al. [2023], Yang and Osher [2024] study the in-context operator learning of differential  
71 equations where the adaption of the pre-trained model is achieved by only conditioning on new  
72 prompts. While these empirical work show great transferabilities of SciFMs for solving PDEs, their  
73 theoretical guarantees are largely open. To the best of our knowledge, this work is the first to derive  
74 the theoretical error bounds of transformers for learning linear elliptic PDEs in context.

75 **Theory of ICL for linear regression and other statistical models.** The work Garg et al. [2022]  
76 provides theoretical understanding of the ability of transformers in learning simple functions in-  
77 context. In the follow-up works Akyürek et al. [2022], Von Oswald et al. [2023], it is shown by  
78 explicit construction of attention matrices that linear transformers can implement a single step of  
79 gradient descent when given a new in-context linear regression task, and numerical experiments  
80 supported that trained transformer indeed implement gradient descent on unseen tasks. Several recent  
81 works Mahankali et al. [2024], Zhang et al. [2023], Ahn et al. [2024] extend the results of Von Oswald  
82 et al. [2023] by proving that one step of gradient descent is indeed optimal for learning linear models  
83 in-context. These works are further complemented by ICL guarantees for learning nonlinear functions  
84 Bai et al. [2024], Cheng et al. [2023], Kim et al. [2024] and for reinforcement learning problems Lin  
85 et al. [2023].

86 Among the aforementioned works, the settings of Zhang et al. [2023], Ahn et al. [2024], Chen et al.  
87 [2024b] are closest to us. Our theoretical bound on the population risk extends the results of Zhang

88 et al. [2023], Ahn et al. [2024] for the linear regression tasks to the tasks of inverting linear systems  
 89 that are associated to elliptic PDEs. Our main novelty is that our results apply to a much larger  
 90 class of task distributions, since our task matrices must respect the PDE structure. In particular,  
 91 this leads to new and nontrivial results regarding task distribution shifts, whereas the effect of task  
 92 distribution shifts is simple under the assumptions of the aforementioned works. We also provide  
 93 sample complexity bounds with respect to the number of pre-training tasks, which have not addressed  
 94 by the above works.

## 95 2 Problem set-up

### 96 2.1 In-context operator learning of linear elliptic PDEs

97 Consider the second-order strongly-elliptic PDE on a bounded Lipschitz domain  $\Omega \subseteq \mathbb{R}^{d_0}$ :

$$\begin{cases} \mathcal{L}_{a,V}u(x) := -\nabla \cdot (a(x)\nabla u(x)) + V(x)u(x) = f(x), & x \in \Omega \\ u(x) = 0, & x \in \partial\Omega. \end{cases} \quad (1)$$

98 where  $a \in L^\infty(\Omega)$  is strictly positive,  $V \in L^\infty(\Omega)$  is non-negative and  $f \in \mathcal{X}_f \subset L^2(\Omega)$ . By the  
 99 standard well-posedness of the elliptic PDE, the solution  $u \in \mathcal{X}_u \subset H_0^1(\Omega)$ . We are interested in  
 100 learning the linear solution operator  $\Psi : f \rightarrow u \in \mathcal{X}_u$  in context Yang et al. [2023]. More specifically,  
 101 at the training stage we are given a training dataset comprising  $N$  length- $n$  prompts of source-solution  
 102 pairs  $\{(f_i^j, u_i^j)_{i=1}^n\}_{j=1}^N$ , where  $\{f_i^j\} \stackrel{i.i.d.}{\sim} P_f$  for some distribution  $P_f$  on the space of functions  $f$ ,  
 103 and  $u_i^j$  are the solutions corresponding to  $f_i^j$  and parameters  $(a_j, V_j) \stackrel{i.i.d.}{\sim} P_a \times P_V$ , where  $P_a$  and  
 104  $P_V$  are distributions on the coefficient  $a$  and  $V$  respectively. An ICL model, after pre-trained on the  
 105 data above, is asked to predict the solution  $u$  for a new source term  $f$  conditioned on a new prompt  
 106  $(f_i, u_i)_{i=1}^m$  which may or may not have the same distribution as the training prompts. Further, the  
 107 prompt-length  $m$  in the downstream task may be different from the prompt-length  $n$  in the training.

108 While the ideal ICL problem above is stated for learning operators defined on infinite dimensional  
 109 function spaces, a practical ICL model (e.g. a transformer) can only operate on finite dimensional  
 110 data, which are typically observed in the form of finite dimensional projections or discrete evaluations.  
 111 To be more concrete, let  $\{\phi_k(x)\}_{k=1}^\infty$  be a basis on both  $\mathcal{X}_u$  and  $\mathcal{X}_f$ , and define a truncated base  
 112 set  $\Phi(x) := [\phi_1, \dots, \phi_d(x)]$  for some  $d < \infty$ . An approximate solution  $\tilde{u}$  to problem (1) can  
 113 be constructed in the framework of Galerkin method: we seek  $\tilde{u}(x) = \langle \mathbf{u}, \Phi(x) \rangle$  where  $\mathbf{u} \in \mathbb{R}^d$   
 114 solves the linear system  $A\mathbf{u} = \mathbf{f}$ , where the matrix  $A = (A_{ij}) \in \mathbb{R}^{d \times d}$  and the right hand side  
 115  $\mathbf{f} = (f_i) \in \mathbb{R}^d$  are defined by

$$A_{ij} = \langle \phi_j, \mathcal{L}_{a,V}\phi_i \rangle \text{ and } f_i = \langle f, \phi_i \rangle, i, j = 1, \dots, d. \quad (2)$$

116 As quantitative discretization error bounds of PDEs are well established, e.g. for finite element  
 117 methods Brenner and Scott [2007] and spectral methods Shen et al. [2011], this paper focuses on the  
 118 error analysis of in-context learning of the finite dimensional linear systems defined by the matrix  
 119 inversion  $A^{-1}$ , which will ultimately translate to estimation bounds for the PDEs.

### 120 2.2 ICL of linear systems

121 The consideration above reduces the original infinite dimensional in-context operator learning problem  
 122 to the finite dimensional ICL problem of solving linear systems. To keep the framework more general,  
 123 we make the following change of notations:  $\mathbf{f} \rightarrow \mathbf{y}$  and  $\mathbf{u} \rightarrow \mathbf{x}$ . An ICL model operates on a prompt  
 124 of  $n$  input-output pairs, denoted by  $S := \{(\mathbf{y}_i, \mathbf{x}_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}^d$  with  $\mathbf{x}_i = A^{-1}\mathbf{y}_i$  as well as a  
 125 new query input  $\mathbf{y}_{n+1} \in \mathbb{R}^d$ . Given multiple prompts, the model aims to predict  $\mathbf{x}_{n+1}$  corresponding  
 126 to the new independent query input  $\mathbf{y}_{n+1}$ . Unlike in supervised learning, each prompt the model  
 127 takes is drawn from a different data distribution. To be more precise, for  $j = 1, \dots, N$ , we assume  
 128 that the  $j$ -th prompt  $S^{(j)} := \{(\mathbf{y}_i^{(j)}, \mathbf{x}_i^{(j)})\}_{i=1}^n$  is generated from the sources  $\{\mathbf{y}_i^{(j)}\}_{i=1}^n \stackrel{i.i.d.}{\sim} p_{\mathbf{y}}$ ; the  
 129 solutions  $\mathbf{x}_i^{(j)}$  are associated to the  $j$ -th inversion task via  $\mathbf{x}_i^{(j)} = (A^{(j)})^{-1}\mathbf{y}_i^{(j)}$  where the matrices  
 130  $A^{(j)} \stackrel{i.i.d.}{\sim} p_A$ . Informed by task matrices derived from discretizations of PDEs as illuminated in (2),  
 131 we make the following assumption on the task distribution  $p_A$ .

132 **Assumption 1.** *The task distribution  $p_A$  is supported on the set of symmetric positive definite*  
 133 *matrices, and there exist constants  $c_A, C_A > 0$  such that the bounds  $c_A^{-1}\mathbf{I}_d \prec A \prec C_A\mathbf{I}_d$  hold for all*  
 134  *$A \in \text{supp}(p_A)$ . The source term  $\mathbf{y}$  follows a Gaussian distribution  $N(0, \Sigma)$ .*

135 Observe that Assumption 1 on  $A$  is very mild and holds for instance whenever the coefficient  $a$  is  
 136 strictly positive and  $V$  is non-negative and bounded. We will make repeated use of the bounds<sup>1</sup>

$$\|A^{-1}\|_{\text{op}} \leq c_A, \|A\|_{\text{op}} \leq C_A, p_A - \text{a.s.} \quad (3)$$

137 The Gaussian assumption on the covariate  $\mathbf{y}$  holds when we assume that the source term  $f$  of the PDE  
 138 is drawn from a Gaussian measure  $N(0, \Sigma_f)$ , where  $\Sigma_f : L^2(\Omega) \rightarrow L^2(\Omega)$  is bounded, in which  
 139 case the covariance matrix  $\Sigma$  is defined by  $\Sigma_{ij} = \langle \Sigma_f \phi_i, \phi_j \rangle_{L^2(\Omega)}$ .

### 140 2.3 Linear transformer architecture for linear systems

141 Inspired by the recent line of work on ICL of linear functions, we consider a linear transformer defined  
 142 by a single-layer linear self-attention layer for our ICL model. Following the standard convention, we  
 143 encode the data of each prompt into a prompt matrix

$$Z = \begin{bmatrix} \mathbf{y}_1 & \cdots & \mathbf{y}_n & \mathbf{y}_{n+1} \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n & 0 \end{bmatrix} \in \mathbb{R}^{D \times (n+1)}, \quad (4)$$

where  $D = 2d$ . For  $\tilde{P}, \tilde{Q} \in \mathbb{R}^{D \times D}$ , the linear self-attention module with parameters  $\tilde{\theta} = (\tilde{P}, \tilde{Q})$  is given by

$$\text{Attn}_{\tilde{\theta}}(Z) = Z + \frac{1}{n} \tilde{P} Z M Z^T \tilde{Q} Z,$$

where  $M = \begin{bmatrix} \mathbf{I}_n & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{(n+1) \times (n+1)}$  is a masking matrix to account for the asymmetry of the prompt matrix. Our definition of the self-attention module makes several simplifying assumptions compared to the standard definition in the literature, namely we merge the key and query matrices into a single matrix  $Q$  and we omit the softmax activation function. A transformer  $f_{\tilde{\theta}}$  predicts a new label  $\mathbf{x}$  for the downstream task by reading out the  $\mathbf{x}$ -component from the self-attention output, i.e.

$$f_{\tilde{\theta}}(Z) := [\text{Attn}_{\tilde{\theta}}(Z)]_{d+1:D, n+1} = \sum_{j=1}^d \langle \mathbf{e}_{d+j}, \text{Attn}_{\tilde{\theta}}(Z) \mathbf{e}_{n+1} \rangle \mathbf{e}_{d+j},$$

where  $\mathbf{e}_i$  denotes the  $i^{\text{th}}$  standard basis vector. Since the output of the transformer only reads out the last  $d$  entries on the bottom right of the output of the self-attention layer, many blocks in  $\tilde{P}$  and  $\tilde{Q}$  do not actually play a role in the prediction defined by the transformer. More precisely, similar to Von Oswald et al. [2023], Zhang et al. [2023], Ahn et al. [2024], if we set  $\tilde{P} = \begin{bmatrix} 0 & 0 \\ 0 & P \end{bmatrix}$  and  $\tilde{Q} = \begin{bmatrix} Q & 0 \\ 0 & 0 \end{bmatrix}$  with  $P, Q \in \mathbb{R}^{d \times d}$ , then output of the transformer can be re-written in a compact form: with  $\theta = (P, Q)$ ,

$$\text{TF}_{\theta}(Z) = P A^{-1} Y_n Q \mathbf{y},$$

144 where  $Y_n := \frac{1}{n} \sum_{k=1}^n \mathbf{y}_k \mathbf{y}_k^T$  denotes the empirical covariance matrix associated to the in-context  
 145 examples. We work with this simplified parameterization for the remainder of our theoretical analysis.

### 146 2.4 Generalization of ICL

147 Our goal is to find the attention matrices  $P$  and  $Q$  that minimize the *population risk* functional

$$\mathcal{R}_n(P, Q; n) = \mathbb{E} \left[ \left\| \text{TF}_{\theta}(Z) - A^{-1} \mathbf{y} \right\|^2 \right] = \mathbb{E} \left[ \left\| P A^{-1} Y_n Q \mathbf{y} - A^{-1} \mathbf{y} \right\|^2 \right], \quad (5)$$

<sup>1</sup>Most of our estimates involve bounds on the norm of  $A^{-1}$ , since it represents the 'solution operator' of the PDE. However, for technical reasons, we also require a bound on the norm of  $A$ .

148 where the expectation is taken over  $A \sim p_A$ ,  $\{\mathbf{y}, \mathbf{y}_1, \dots, \mathbf{y}_n\} \sim N(0, \Sigma)^{\otimes n+1}$ . Since we do not  
 149 have access to the distribution on tasks,  $P$  and  $Q$  are instead trained by minimizing the corresponding  
 150 empirical risk functional defined on  $N$  tasks:

$$\mathcal{R}_{n,N}(P, Q) = \frac{1}{N} \sum_{i=1}^N \left\| P A_i^{-1} Y_n^{(i)} Q \mathbf{y}_i - A_i^{-1} \mathbf{y}_i \right\|^2, \quad (6)$$

151 where  $\{A_i\} \stackrel{i.i.d.}{\sim} p_A$ ,  $\{\mathbf{y}_i\} \stackrel{i.i.d.}{\sim} N(0, \Sigma)$ , and  $Y_n^{(i)}$  is the empirical covariance matrix associated to  
 152 the in-context examples  $\{\mathbf{y}_1^{(i)}, \dots, \mathbf{y}_n^{(i)}\}$  which are (jointly) independent from  $\mathbf{y}_i$ .

153 A pre-trained transformer is expected to make predictions on a downstream task that consists of a new  
 154 length- $m$  prompt  $\{(\mathbf{y}_i, \mathbf{x}_i)\}_{i=1}^m = \{(\mathbf{y}_i, (A')^{-1} \mathbf{y}_i)\}_{i=1}^m$  and a new test sample  $\mathbf{y}$ , where the input  
 155 samples  $\{\mathbf{y}_i\}_{i=1}^n \cup \{\mathbf{y}\} \sim P'_y$  and the matrix  $A' \sim P'_A = N(0, \Sigma')$ . Our primary interest is to bound  
 156 the generalization performance (measured by the prediction risk) of the pre-trained transformer for  
 157 the downstream task in two different scenarios.

158 • **In-domain generalization:** The distributions of tasks and of prompt data in the pre-training are  
 159 the same as these in the downstream task ( $P_y = P'_y$  and  $P_A = P'_A$ ). Thus in-domain generalization  
 160 measures the testing performance on unseen samples in the downstream task that do not appear in the  
 161 training samples. The in-domain generalization error is defined by

$$\mathcal{R}_m(P, Q; m) = \mathbb{E}_{A \sim p_A, (y_1, \dots, y_m, y) \sim N(0, \Sigma)^{\otimes (m+1)}} \left[ \left\| P A^{-1} Y_m Q \mathbf{y} - A^{-1} \mathbf{y} \right\|^2 \right]. \quad (7)$$

162 • **Out-of-domain (OOD) generalization:** The distributions of tasks or within-task data in the  
 163 pre-training are different from those in the downstream task, i.e.  $P_y \neq P'_y$  or  $P_A \neq P'_A$ . Specifically,  
 164 the OOD-generalization error with respect to the task distribution shift is defined by

$$\mathcal{R}_m^{P'_A}(P, Q; m) = \mathbb{E}_{A' \sim p'_A, (y_1, \dots, y_m, y) \sim N(0, \Sigma)^{\otimes (m+1)}} \left[ \left\| P (A')^{-1} Y_m Q \mathbf{y} - (A')^{-1} \mathbf{y} \right\|^2 \right]. \quad (8)$$

165 We also define the OOD-generalization error with respect to the covariate distribution shift by

$$\mathcal{R}_m^{\Sigma'}(P, Q; m) = \mathbb{E}_{A \sim p_A, (y_1, \dots, y_m, y) \sim N(0, \Sigma')^{\otimes (m+1)}} \left[ \left\| P A^{-1} Y_m Q \mathbf{y} - A^{-1} \mathbf{y} \right\|^2 \right]. \quad (9)$$

166 Notice that the prompt length  $m$  in the prediction risk need not equal the prompt length  $n$  in the  
 167 model pre-training. We are particularly interested in quantifying the scaling laws of the generalization  
 168 errors for the pre-trained transformer as the amount of data increases to infinity, i.e.  $N, n, m \uparrow \infty$ .

### 169 3 Theoretical results

#### 170 3.1 Error bounds for in-domain generalization of learning linear systems

171 Our first result studies the generalization ability of the transformer obtained by empirical risk  
 172 minimization over a set of norm-constrained transformers, where the error is measured by the  
 173 prediction risk  $\mathcal{R}_m$ .

**Theorem 1.** Let  $\hat{\theta} = (P_N, Q_N) \in \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_{n,N}(\theta)$ , where  $\|\theta\| := \max(\|P\|_{op}, \|Q\|_{op})$ .  
 Then for  $n$  sufficiently large and  $m \leq n$ , we have with probability  $\geq 1 - \frac{1}{\operatorname{poly}(N)}$ ,

$$\mathcal{R}_m(\hat{\theta}) \lesssim \frac{1}{m} + \frac{1}{n^2} + \frac{d^2}{\sqrt{N}},$$

174 where the implicit constants depend on  $M$ , the data covariance  $\Sigma$ , and the task distribution  $p_A$ , and  
 175 we have omitted factors which are polylog in  $N$ .

176 The precise statement of Theorem 1 is given in Appendix B, where we discuss what happens when  
 177  $m > n$ . We refer to  $m \leq n$  as the practical regime, since it is commonly satisfied by large pre-trained  
 178 transformers. Notice that the prompt lengths during training and testing contribute different rates  
 179 to the overall sample complexity bound, with the sequence length  $n$  during training contributing  
 180 an  $O(n^{-2})$  rate while the sequence length  $m$  at inference contributing an  $O(m^{-1})$  rate; a similar  
 181 phenomenon was observed in [Zhang et al., 2023, Theorem 4.2] for in-context linear regression.

### 182 3.2 Error bounds for in-domain generalization of learning elliptic PDEs

183 Building upon Theorem 1, we proceed to bound the ICL-generalization error for learning the elliptic  
 184 PDE (1). Our next result provides a rather general upper bound on the ICL-generalization error for  
 185 the PDE solution in terms of the spatial discretization error of the PDE and the ICL-generalization  
 186 error in learning the finite linear systems associated to the discretization. The discretization error is  
 187 typically fully determined by the number  $d$  of basis functions used in the Galerkin projection. The  
 188 second term is bounded by Theorem 1. In the following result, let  $u$  denote the solution to the elliptic  
 189 PDE specified by (1). We write  $u_d$  for a discrete approximation to  $u$  with the mesh size  $h$  and we  
 190 write  $\widehat{u}_d$  for the approximate solution obtained by solving a discrete linear system with a pre-trained  
 191 transformer.

**Theorem 2.** *Let  $\Phi'$  be the stiffness matrix defined by  $\Phi'_{ij} = (\phi'_i, \phi'_j)_{L^2(\Omega)}$  and let  $\Phi$  be the mass matrix defined by  $\Phi_{ij} = (\phi_i, \phi_j)_{L^2(\Omega)}$ . Assume that both matrices are symmetric and positive definite. Then,*

$$\mathbb{E}\|u - \widehat{u}_d\|_{H^1(\Omega)}^2 \lesssim \mathbb{E}\|u - u_d\|_{H^1(\Omega)}^2 + (1 + \lambda_{\max}(\Phi^{-1/2}\Phi'\Phi^{-1/2})) \cdot \mathcal{R}_m(\widehat{\theta}),$$

192 where  $\widehat{\theta}$  is a minimizer of the empirical risk defined in Theorem 1 and  $\lambda_{\max}(\cdot)$  denotes the largest  
 193 eigenvalue of a symmetric positive definite matrix.

194 Theorem 2 bounds the in-domain generalization error of ICL of the PDE as a sum of the discretization  
 195 error of the PDE solver and the statistical error of learning the linear system associated to the  
 196 discretization of the PDE. It is worth-noting that there is a trade-off between the two terms; the  
 197 first term decreases as the number of basis functions (or fineness of the mesh) increases, while  
 198 the prefactor  $\lambda_{\max}(\Phi^{-1/2}\Phi'\Phi^{-1/2})$  in the second term can grow as the number of basis functions  
 199 tends to infinity. The abstract bound established in Theorem 2 is agnostic to the choice of PDE  
 200 discretization. We show in Appendix C how this result can be used to derive an explicit error estimate  
 201 for the ICL in the context of a  $P^1$ -finite element discretization of the PDE in one dimension.

### 202 3.3 OOD-generalization under task distribution shift

203 Let  $\widehat{\theta}$  denote the minimizer of the empirical risk  $\mathcal{R}_{n,N}$  over the bounded set  $\{\|\theta\| \leq M\}$  for some  
 204  $M > 0$ , and recall that the training tasks (modeled by  $A$ ) are drawn from a distribution  $p_A$ . Let  
 205  $p'_A$  denote the distribution of  $A$  in the downstream tasks, and let  $\mathcal{R}_m, \mathcal{R}'_m$  be the prediction risk  
 206 functionals defined as in (8) where the expectations over tasks are taken with respect to  $p_A$  and  
 207  $p'_A$  respectively. We would like to bound the quantity  $\mathcal{R}'_m(\widehat{\theta})$ , which represents the test error of  
 208 the trained transformer under a shift on the task distribution. We say that a pre-trained model  $\widehat{\theta}$   
 209 **achieves OOD generalization** if its population risk with respect to the downstream task distribution  
 210  $p'_A$  converges to zero in probability:  $\lim_{(m,n,N) \rightarrow \infty} \mathcal{R}'_m(\widehat{\theta}) \xrightarrow{P} 0$ . In order to state our results on  
 211 OOD generalization, we first introduce the following 'infinite-context' variant of the in-domain  
 212 denoted by  $\mathcal{R}_\infty$ :

$$\mathcal{R}_\infty(\theta) = \mathbb{E}_{A \sim p_A} [\|(PA^{-1}\Sigma Q - A^{-1})\Sigma^{1/2}\|_F^2]. \quad (10)$$

213 We also define an OOD-generalization risk  $\mathcal{R}'_\infty$  similar to above with  $p_A$  replaced by  $p'_A$ . We denote  
 214 by  $\mathcal{M}_\infty$  and  $\mathcal{M}'_\infty$  the sets of minimizers of  $\mathcal{R}_\infty$  and  $\mathcal{R}'_\infty$  respectively. We are now able to define the  
 215 key notion of task diversity.

216 **Definition 1.** *The pre-training task distribution  $p_A$  is **diverse** relative to the downstream task  
 217 distribution  $p'_A$  if  $\mathcal{M}_\infty \subseteq \mathcal{M}'_\infty$ .*

218 The importance of task diversity has been observed in the prior work Tripuraneni et al. [2020] for  
 219 transfer learning. Our notion of diversity differs from the previous notion in that we compare the  
 220 set of minimizers of population losses instead of the loss values. Theorem 3 below shows that the  
 221 task diversity, in the sense of Definition 1, is sufficient for the pre-trained transformer to achieve  
 222 OOD-generalization.

**Theorem 3.** *Let  $p_A$  and  $p'_A$  denote the pre-training and downstream task distributions respectively, and suppose  $p_A$  is diverse relative to  $p'_A$ . Then, with  $\widehat{\theta} \in \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_{n,N}(\theta)$ , we have*

$$\mathcal{R}'_m(\widehat{\theta}) \lesssim \mathcal{R}_m(\widehat{\theta}) + \frac{d(p_A, p'_A)}{m} + \operatorname{dist}(\widehat{\theta}, \mathcal{M}_\infty)^2,$$

223 where  $d(p_A, p'_A)$  is a discrepancy between the pre-training and downstream task distributions that  
 224 satisfies  $d(p_A, p'_A) = 0$  if  $p_A = p'_A$ .

225 The precise definition of the discrepancy  $d(p_A, p'_A)$  is technical and can be found in the statement of  
 226 Lemma 2 in the appendix. The OOD generalization error is bounded by a sum of three terms: the  
 227 in-domain generalization error, the task-shift error, and the model error, the latter of which is captured  
 228 by  $\text{dist}(\hat{\theta}_n, \mathcal{M}_\infty)$ . A salient feature of Theorem 3, compared to the prior ICL-generalization bound  
 229 Mroueh [2023] under distribution shift, is that the task-shift error inherits a factor of  $m^{-1}$ , which  
 230 elucidates the robustness of transformers under shifts in the task distribution. Theorem 3 also extends  
 231 the prior OOD-generalization result of ICL for linear regression Zhang et al. [2023] to learning linear  
 232 systems. However, unlike in the linear regression setting, the set of minimizers of the population  
 233 risk in the linear system setting can vary substantially when the task distribution changes, we need  
 234 the training tasks to be sufficient diverse compared to the downstream tasks in order to control the  
 235 additional model error due to the change of the minimizers; see Appendix D for more details. We  
 236 also note that Proposition 4 in the appendix shows that the minimizers of the empirical risk converge  
 237 in probability to the minimizers of  $\mathcal{R}_\infty$ , thus guaranteeing that the bound in Theorem 3 is  $o_P(1)$ .

238 Since task diversity is sufficient to achieve OOD generalization, it is natural to ask what conditions on  
 239  $p_A$  and  $p'_A$  guarantee task diversity. The following result provides two sufficient conditions. We refer  
 240 the readers to Appendix D for additional discussions on task diversity. To state the result, we recall that  
 241 the notion of the centralizer  $\mathcal{C}(\mathcal{S})$  of a subset  $\mathcal{S} \subseteq \mathbb{R}^{d \times d}$ :  $\mathcal{C}(\mathcal{S}) = \{P \in \mathbb{R}^{d \times d} : PS = SP \forall S \in \mathcal{S}\}$ .

242 **Theorem 4.** *Let  $p_A, p'_A$  be two distributions on the matrices  $A$  that satisfy Assumption 1. Then*

- 243 1. *If  $\text{supp}(p'_A) \subseteq \text{supp}(p_A)$ , then  $p_A$  is diverse relative to  $p'_A$ .*
- 244 2. *Define  $\mathcal{S}(p_A) := \{A_1 A_2^{-1} : A_1, A_2 \in \text{supp}(p_A)\}$ . If  $\mathcal{C}(\mathcal{S}(p_A)) = \{c\mathbf{I} : c \in \mathbb{R}\}$ , then  $p_A$  is*  
 245 *diverse relative to any distribution  $p'_A$ .*

246 The first statement of Theorem 4 is a natural one: it says that the pre-training task distribution is  
 247 diverse whenever the downstream task distribution is a 'subset' of it, in the sense of supports. The  
 248 second condition is particularly interesting because it implies OOD-generalization (by Theorem 3)  
 249 regardless of the downstream task distribution. The second condition based on the centralizer of the  
 250 set  $\mathcal{S}(p_A)$  is less obvious, but heuristically it enforces that the support of  $p_A$  must be large enough  
 251 that the only matrices which can commute with all pairwise products in  $\mathcal{S}(p_A)$  are scalars. Our  
 252 empirical results suggest that the task distributions associated to elliptic PDE problems are diverse.

### 253 3.4 OOD-generalization under covariate distribution shift

254 We now study the OOD-generalization error due to the distribution shift with respect to the Gaussian  
 255 covariates  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ , i.e., the vectors at which a task matrix  $A$  is evaluated. The next proposition  
 256 provides a quantitative upper bound for the generalization error in terms of the discrepancy between  
 257 the covariance matrices. To simplify the proof, we use a Frobenius norm bound on the empirical risk  
 258 minimizer. However, this choice of norm is not essential to the result.

259 **Theorem 5.** *Let  $\Sigma = W\Lambda W^T$  and  $\tilde{\Sigma} = \tilde{W}\tilde{\Lambda}\tilde{W}^T$  be the covariance matrices of Gaussian covariates*  
 260 *used in the training and testing respectively. Let  $(\hat{P}, \hat{Q})$  be minimizers of the empirical risk associated*  
 261 *to covariates sampled from  $N(0, \Sigma)$  and take  $M > 0$  such that  $\max(\|\hat{P}\|_F, \|\hat{Q}\|_F) \leq M$ . Then*

$$\mathcal{R}_m^\Sigma(\hat{P}, \hat{Q}) \lesssim \mathcal{R}_m^\Sigma(\hat{P}, \hat{Q}) + \|\Sigma - \tilde{\Sigma}\|_{op} + \frac{1}{m}\|W - \tilde{W}\|_{op}.$$

262 where the implicit constants depend on  $M, \Sigma, \tilde{\Sigma}$ , and the constant  $c_A$  defined in Assumption 1.

263 Theorem 5 states that the OOD-generalization error with respect to the covariate distribution shift  
 264 is Lipschitz stable with respect to changes in the covariance matrix. However, unlike the case of  
 265 task distribution shift, the covariate distribution shift error can not be mitigated by increasing the  
 266 prompt-length in the downstream task; see also Figure 3. A similar phenomenon was observed in  
 267 Zhang et al. [2023].

268 **4 Numerical experiments**

269 **4.1 In-domain generalization**

270 We first investigate numerically the neural scaling law of the transformer model for solving the linear  
 271 system associated to the Galerkin discretization of the elliptic PDE (1) in the setting of in-domain  
 272 generalization. More precisely, we consider the one dimensional elliptic PDE  $(-\Delta + V(x))u(x) =$   
 273  $f(x)$  on  $\Omega = [0, 1]$  with Dirichlet boundary condition. We assume that the source  $f \sim N(0, \mathbb{I})$ ,  
 274 where  $\mathbb{I}$  denotes the identity operator. We discretize the PDE using Galerkin projection under  $d$  sine  
 275 bases. Further we assume that the potential  $V$  is uniform random field that is obtained by dividing the  
 276 domain into  $2d + 1$  sub-intervals and in each cell independently, the potential takes values uniformly  
 277 in  $[1, 2]$ . In Figure 1: A-C, we demonstrate the empirical scaling law of the linear transformer for  
 278 learning the discrete linear system by showing the log-log plots of the  $\ell^2$ -errors as functions of the  
 279 number of pre-training tasks  $N$ , the sequence length  $n$  during training and the sequence length  $m$   
 280 at inference. These numerical results suggest that the decaying rates of the prediction errors are  
 281  $O(N^{-\frac{1}{2}})$ ,  $O(n^{-2})$  and  $O(m^{-1})$  respectively, which agree with the rates predicted in Theorem 1 in  
 282 the practical regime  $m \leq n$ . We also demonstrate the ICL-generalization error for learning the PDE  
 283 solutions. Figure 1:D shows that prediction error increases as  $d$  increases indicating that ICL of the  
 284 linear system becomes harder as the  $d$  increases.

285 Figure 2:B shows the  $H^1$ -error curve between the numerical solution predicted by the ICL-model  
 286 and the ground-truth as a function of the number of bases  $d$ , while fixing the prompt-lengths and the  
 287 number of tasks. The U-shaped curve indicates the trade-off between the dimension of the discrete  
 288 problem and the amount of data. More details on the experiment set-ups can be found in Appendix H.

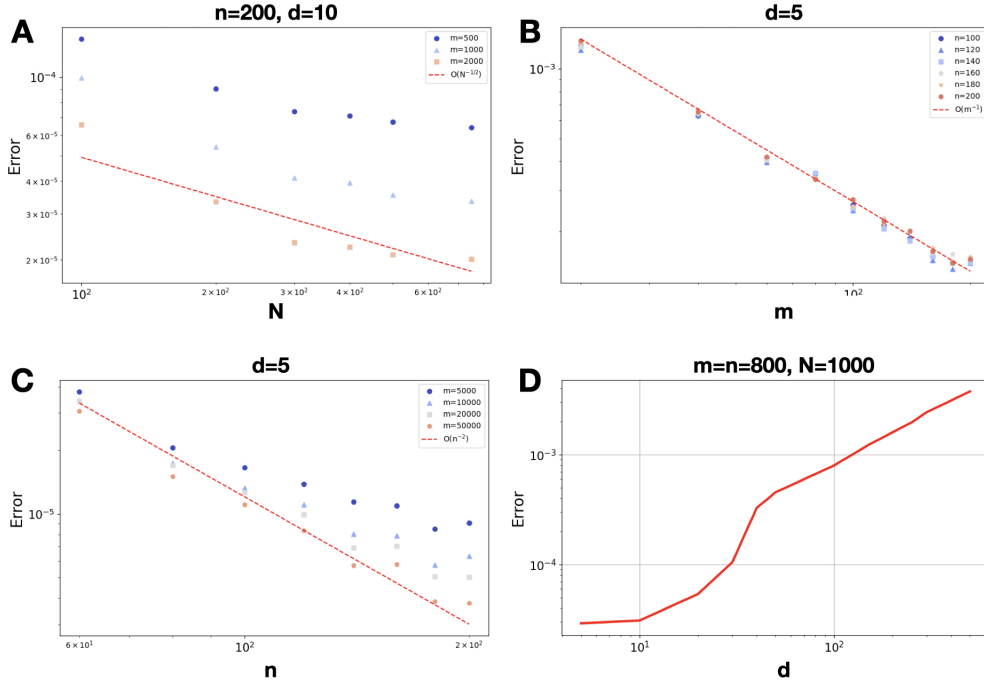


Figure 1: The figures A-D show the log-log plots for the  $\ell^2$ -error of learning the linear system associated to the PDE discretization with respect to the number of tasks  $N$ , the prompt length  $n$  during training, the prompt length  $m$  during inference, and the dimension  $d$  of the linear system.

289 **4.2 Out-of-domain generalization**

290 **Task shifts.** We validate the ICL-capability of pre-trained transformers for learning the linear systems  
 291 and PDEs under task distribution shifts. Specifically, for the PDE (1) in one dimension, we consider  
 292 the task distribution shifts in  $a$  and  $V$  exclusively. To sample  $a(x)$ , we write  $a(x) = e^{b(x)}$ , where  
 293  $b(x)$  is sampled from a centered normal distribution with covariance operator  $-(\Delta + \tau\mathbb{I})^{-\alpha}$ , for



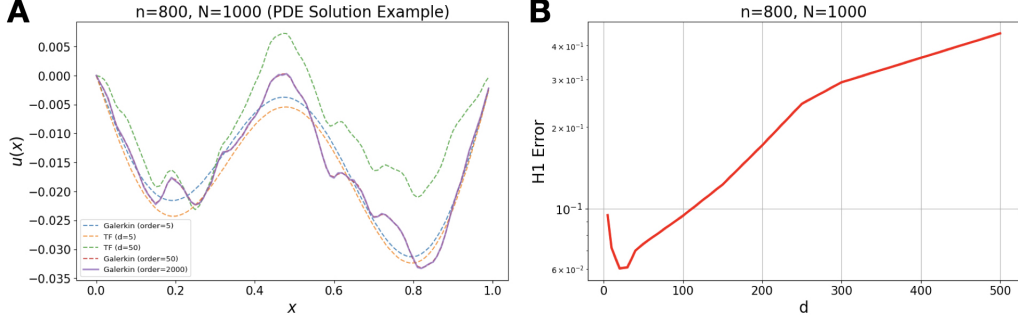


Figure 2: The left plot shows the PDE solution defined by the pre-trained transformer with the reference solution, obtained by Galerkin’s method with 2000 basis functions. The right plot shows the  $H^1$ -error between the solution predicted by the transformer and reference solution with respect to the number of Galerkin basis functions  $d$ .

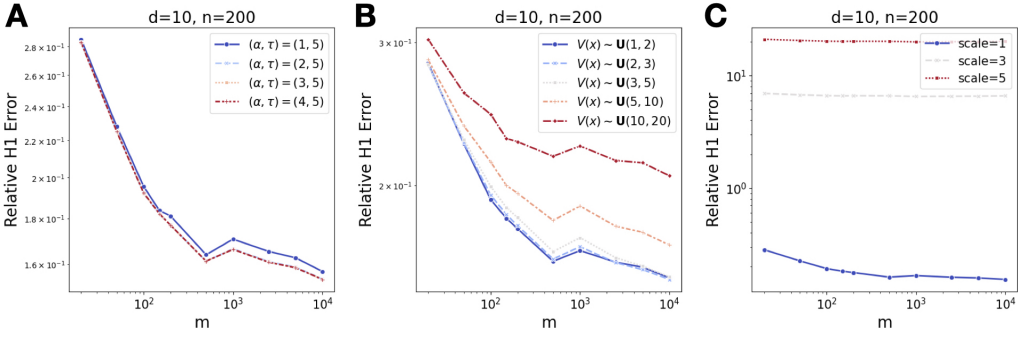


Figure 3: Figures A, B show the relative  $H^1$ -error under shifts on  $a(x)$  and  $V(x)$  respectively. Figure C shows the relative  $H^1$ -error under the covariate shift on the source term  $f$ .

294  $\alpha, \tau > 0$ . During training, we set  $\alpha = 3$  and  $\tau = 5$ , and during inference, we vary the values  
 295 of  $\alpha$  and  $\tau$  according to Figure 3: A. We assume the potential  $V$  is piecewise constant on  $2d + 1$   
 296 subintervals and that the value of  $V$  on each cell is drawn according to the uniform distribution on  
 297  $[a, b]$ . During training, we set  $[a, b] = [1, 2]$ , and we vary the values of  $[a, b]$  at inference according to  
 298 Figure 3: B. For further details on the experimental setup, see in Appendix H. Figure 3: A shows that  
 299 the pre-trained transformer can perform equally well on tasks on smoother  $a$  but perform slightly  
 300 worse on tasks with less regular  $a$ . Figure 3: B shows the OOD-generalization errors increase as the  
 301 distribution shift in  $V$  becomes stronger, but they decrease as the context length at inference increases,  
 302 as predicted by Theorem 3.

303 **Covariate shifts.** Finally, we test the performance of the pre-trained transformer under covariate  
 304 distribution shifts. Specifically, we train the model to solve the PDE (1), where the source term  
 305  $f \sim N(0, C)$  for  $C = (-\Delta + c\mathbb{I})^{-\beta}$ , where  $c, \beta > 0$  are fixed. Then, at inference, we consider  
 306 solving the same PDE, but where the source term is defined by  $N(0, 3C)$  or  $N(0, 5C)$ . Figure 3  
 307 show that the pre-trained transformers are not robust to covariate distribution shifts. We refer to  
 308 Figure 5 in the appendix for additional numerical results for the covariant shifts in  $c$  and  $\beta$ .

## 309 5 Conclusion

310 In this work, we studied the ability of a transformer characterized by a single linear self-attention  
 311 layer to in-context learn the solution operator of a linear elliptic PDE. We characterized the role of  
 312 the number of pre-training task, the number of in-context examples during pre-training and testing,  
 313 the mesh size, and various distribution shifts on the PDE coefficients in the overall PDE recovery  
 314 error. We also provided thorough numerical experiments to demonstrate our theory. There are several  
 315 natural extensions of this work, such as to nonlinear and time-dependent PDE problems. In these  
 316 more complex settings, it is crucial to characterize the role that depth and nonlinearity play in the  
 317 ability of transformers to approximate the PDE solution. We leave these directions to future work.

## 318 References

- 319 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
320 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical  
321 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 322 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
323 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing  
324 systems*, 30, 2017.
- 325 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
326 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
327 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 328 Ilyes Batatia, Philipp Benner, Yuan Chiang, Alin M Elena, Dávid P Kovács, Janosh Riebesell,  
329 Xavier R Advincula, Mark Asta, William J Baldwin, Noam Bernstein, et al. A foundation model  
330 for atomistic materials chemistry. *arXiv preprint arXiv:2401.00096*, 2023.
- 331 Albi Celaj, Alice Jiexin Gao, Tammy TY Lau, Erle M Holgersen, Alston Lo, Varun Lodaya, Christo-  
332 pher B Cole, Robert E Denroche, Carl Spickett, Omar Wagih, et al. An rna foundation model  
333 enables discovery of disease mechanisms and candidate therapeutics. *bioRxiv*, pages 2023–09,  
334 2023.
- 335 Oscar Méndez-Lucio, Christos Nicolaou, and Berton Earnshaw. Mole: a molecular foundation model  
336 for drug discovery. *arXiv preprint arXiv:2211.02657*, 2022.
- 337 Shashank Subramanian, Peter Harrington, Kurt Keutzer, Wahid Bhimji, Dmitriy Morozov, Michael W  
338 Mahoney, and Amir Gholami. Towards foundation models for scientific machine learning: Charac-  
339 terizing scaling and transfer behavior. *Advances in Neural Information Processing Systems*, 36,  
340 2024.
- 341 Michael McCabe, Bruno Régaldó-Saint Blancard, Liam Holden Parker, Ruben Ohana, Miles Cranmer,  
342 Alberto Bietti, Michael Eickenberg, Siavash Golkar, Geraud Krawezik, Francois Lanusse, et al.  
343 Multiple physics pretraining for physical surrogate models. *arXiv preprint arXiv:2310.02994*,  
344 2023.
- 345 Zhanhong Ye, Xiang Huang, Leheng Chen, Hongsheng Liu, Zidong Wang, and Bin Dong. Pdeformer:  
346 Towards a foundation model for one-dimensional partial differential equations. *arXiv preprint  
347 arXiv:2402.12652*, 2024.
- 348 Liu Yang, Siting Liu, Tingwei Meng, and Stanley J Osher. In-context operator learning with data  
349 prompts for differential equation problems. *Proceedings of the National Academy of Sciences*, 120  
350 (39):e2310142120, 2023.
- 351 Jingmin Sun, Yuxuan Liu, Zecheng Zhang, and Hayden Schaeffer. Towards a foundation model  
352 for partial differential equation: Multi-operator learning and extrapolation. *arXiv preprint  
353 arXiv:2404.12355*, 2024.
- 354 Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A  
355 deep learning framework for solving forward and inverse problems involving nonlinear partial  
356 differential equations. *Journal of Computational physics*, 378:686–707, 2019.
- 357 Lu Lu, Pengzhan Jin, and George Em Karniadakis. Deeponet: Learning nonlinear operators for  
358 identifying differential equations based on the universal approximation theorem of operators. *arXiv  
359 preprint arXiv:1910.03193*, 2019.
- 360 Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew  
361 Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential equations.  
362 *arXiv preprint arXiv:2010.08895*, 2020.
- 363 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,  
364 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio  
365 Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with GPT-4.  
366 *arXiv:2303.12712*, 2023.

- 367 Louis Kirsch, James Harrison, Jascha Sohl-Dickstein, and Luke Metz. General-purpose in-context  
368 learning by meta-learning transformers. *arXiv preprint arXiv:2212.04458*, 2022.
- 369 Wuyang Chen, Jialin Song, Pu Ren, Shashank Subramanian, Dmitriy Morozov, and Michael W  
370 Mahoney. Data-efficient operator learning via unsupervised pretraining and in-context learning.  
371 *arXiv preprint arXiv:2402.15734*, 2024a.
- 372 Liu Yang and Stanley J Osher. Pde generalization of in-context operator networks: A study on 1d  
373 scalar nonlinear conservation laws. *arXiv preprint arXiv:2401.07364*, 2024.
- 374 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn  
375 in-context? a case study of simple function classes. *Advances in Neural Information Processing*  
376 *Systems*, 35:30583–30598, 2022.
- 377 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning  
378 algorithm is in-context learning? investigations with linear models. In *The Eleventh International*  
379 *Conference on Learning Representations*, 2022.
- 380 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,  
381 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In  
382 *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023.
- 383 Arvind V Mahankali, Tatsunori Hashimoto, and Tengyu Ma. One step of gradient descent is provably  
384 the optimal in-context learner with one layer of linear self-attention. In *The Twelfth International*  
385 *Conference on Learning Representations*, 2024.
- 386 Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context.  
387 *arXiv preprint arXiv:2306.09927*, 2023.
- 388 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement  
389 preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing*  
390 *Systems*, 36, 2024.
- 391 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:  
392 Provable in-context learning with in-context algorithm selection. *Advances in neural information*  
393 *processing systems*, 36, 2024.
- 394 Xiang Cheng, Yuxin Chen, and Suvrit Sra. Transformers implement functional gradient descent to  
395 learn non-linear functions in context. *arXiv preprint arXiv:2312.06528*, 2023.
- 396 Juno Kim, Tai Nakamaki, and Taiji Suzuki. Transformers are minimax optimal nonparametric  
397 in-context learners. *arXiv preprint arXiv:2408.12186*, 2024.
- 398 Licong Lin, Yu Bai, and Song Mei. Transformers as decision makers: Provable in-context reinforce-  
399 ment learning via supervised pretraining. *arXiv preprint arXiv:2310.08566*, 2023.
- 400 Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head  
401 softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv preprint*  
402 *arXiv:2402.19442*, 2024b.
- 403 Susanne Brenner and Ridgway Scott. *The Mathematical Theory of Finite Element Methods*, volume 15.  
404 Springer Science & Business Media, 2007.
- 405 Jie Shen, Tao Tang, and Li-Lian Wang. *Spectral methods: algorithms, analysis and applications*,  
406 volume 41. Springer Science & Business Media, 2011.
- 407 Nilesh Tripuraneni, Michael Jordan, and Chi Jin. On the theory of transfer learning: The importance  
408 of task diversity. *Advances in neural information processing systems*, 33:7852–7862, 2020.
- 409 Youssef Mroueh. Towards a statistical theory of learning to learn in-context with transformers. In  
410 *NeurIPS 2023 Workshop Optimal Transport and Machine Learning*, 2023.
- 411 Frank Cole and Yulong Lu. Score-based generative models break the curse of dimensionality in  
412 learning a family of sub-gaussian probability distributions. *arXiv preprint arXiv:2402.08082*, 2024.

- 413 Jiyoung Park, Ian Pelakh, and Stephan Wojtowytsch. Minimum norm interpolation by perceptrons:  
 414 Explicit regularization and implicit bias. *Advances in Neural Information Processing Systems*, 36,  
 415 2023.
- 416 Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cam-  
 417 bridge university press, 2019.
- 418 Richard M Dudley. The sizes of compact subsets of hilbert space and continuity of gaussian processes.  
 419 *Journal of Functional Analysis*, 1(3):290–330, 1967.
- 420 Alexandre Ern and Jean-Luc Guermond. *Theory and practice of finite elements*, volume 159. Springer,  
 421 2004.
- 422 Daniele Boffi. Finite element approximation of eigenvalue problems. *Acta numerica*, 19:1–120,  
 423 2010.
- 424 Michael E Sander, Raja Giryes, Taiji Suzuki, Mathieu Blondel, and Gabriel Peyré. How do trans-  
 425 formers perform in-context autoregressive learning? *arXiv preprint arXiv:2402.05787*, 2024.
- 426 Gilbert Strang. *Introduction to linear algebra*. SIAM, 2022.
- 427 Mark Rudelson and Roman Vershynin. Hanson-wright inequality and sub-gaussian concentration.  
 428 2013.
- 429 Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.

## 430 A Notation

431 Before delving into the proofs of our main results, we briefly go over all relevant notation:

- 432 • Physical dimension of PDE problem:  $d_0$
- 433 • Dimension of task matrix for ICL:  $d$
- 434 • Task matrix for ICL:  $A$
- 435 • Covariates for ICL:  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$
- 436 • Prompt matrix for ICL:  $Z$
- 437 • Empirical covariance matrix of  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ :  $Y_n$
- 438 • Distribution on tasks:  $p_A$
- 439 • Upper bound on largest eigenvalue of  $A^{-1}$  over  $\text{supp}(p_A)$ :  $c_A$
- 440 • Covariance operator of the distribution on  $L^2(\Omega)$ -valued covariates:  $\Sigma_f$
- 441 • Covariance matrix of the distribution on  $\mathbb{R}^d$ -valued covariates:  $\Sigma$
- 442 • Parameters of transformer:  $\theta = (P, Q)$
- 443 • Prediction of the transformer with parameters  $\theta$ :  $\text{TF}_\theta(Z)$
- 444 • Population risk for training:  $\mathcal{R}_n$
- 445 • Population risk for inference:  $\mathcal{R}_m$
- 446 • Empirical risk:  $\mathcal{R}_{n,N}$
- 447 • "Infinite-context" population risk:  $\mathcal{R}_\infty$
- 448 • Number of context examples per prompt during training:  $n$
- 449 • Number of context examples per prompt during inference:  $m$
- 450 • Number of pre-training tasks:  $N$

451 **B Proofs for Subsection 3.1**

452 In this section we prove Theorem 1, which controls the (in-distribution) generalization error for  
 453 in-context learning of linear systems in terms of the context length during training, the context length  
 454 during inference, and the number of pre-training tasks. Before the proof, we present a more precise  
 455 statement of the theorem.

456 **Theorem 6** (Theorem 1, precise version). *Let  $\hat{\theta} = (P_N, Q_N) \in \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_{n,N}(\theta)$ , where*  
 457  $\|\theta\| := \max(\|P\|_{op}, \|Q\|_{op})$ . *Then for  $n$  sufficiently large, we have with probability  $\geq 1 - \operatorname{poly}(N)$*

$$\begin{aligned} \mathcal{R}_m(\hat{\theta}) \lesssim & \frac{(c_A^2 + d)\operatorname{Tr}(\Sigma)}{m} + \frac{c_A^2 C_A^4 \|\Sigma\|_{op}^2 \|\Sigma^{-1}\|_{op}^2 \left(1 + \operatorname{Tr}_{\Sigma}(\mathbb{E}_{A \sim p_A}[A^{-2}])\right)^2 \operatorname{Tr}(\Sigma)}{n^2} \\ & + \frac{d^2 c_A^2 \|\Sigma\|_{op}^2 \max(1, \|\Sigma^{-1}\|_{op})^4}{\sqrt{N}} \\ & + \max(1, \|\Sigma^{-1}\|_{op})^4 c_A^2 \max(\operatorname{Tr}(\Sigma), \|\Sigma\|_{op}^2) \operatorname{Tr}(\Sigma) \left| \frac{1}{n} - \frac{1}{m} \right|, \end{aligned} \quad (11)$$

458 where we have omitted factors which are polylog in  $N$ .

459 **Remark 1.** *We would like to comment on the possible suboptimality of the bound (11). Specifically,*  
 460 *the last term on the right side of (11), which we term the "context mismatch error", is mainly due to*  
 461 *our proof strategy and can likely be removed with a refined analysis. This term is not observed in our*  
 462 *numerical experiments; see Figure 1. In the practical<sup>2</sup> regime where the length of the testing prompts*  
 463 *is less than that of the training prompts (i.e.  $m \leq n$ ), we have  $|\frac{1}{n} - \frac{1}{m}| \leq \frac{1}{m}$ , and hence the context-*  
 464 *mismatch error is absorbed into the  $O(\frac{1}{m})$  term, leading to the following overall generalization*  
 465 *bound*

$$\mathcal{R}_m(\hat{\theta}) \lesssim \frac{1}{m} + \frac{1}{n^2} + \frac{1}{\sqrt{N}}. \quad (12)$$

466 *Proof of Theorem 1. **Step 1 - error decomposition:*** Throughout the proof, we use the notation  
 467  $\theta = (P, Q)$  and  $\|\theta\| = \max(\|P\|_{op}, \|Q\|_{op})$ . Write  $\ell(A, Y_n, \mathbf{y}; \theta) = \|(PA^{-1}Y_nQ - A^{-1})\mathbf{y}\|^2$ , so  
 468 that the risk functionals can be expressed as

$$\mathcal{R}_n(\theta) = \mathbb{E}_{A, Y_n, \mathbf{y}} \ell(A, Y_n, \mathbf{y}; \theta), \quad \mathcal{R}_{n,N}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(A_i, Y_n^{(i)}, \mathbf{y}_i; \theta).$$

469 Let us introduce an auxiliary parameters  $t > 0$  – to be specified precisely at the end of the proof –  
 470 and define the events

$$\mathcal{A}_t(Y_n, \mathbf{y}) = \left\{ \|\mathbf{y}\| \leq \sqrt{\operatorname{Tr}(\Sigma)} + t, \|Y_n\|_{op} \leq \|\Sigma\|_{op} \left(1 + t + \sqrt{\frac{d}{n}}\right) \right\}.$$

471 Define the truncated loss function as  $\ell^{R,t}(A, Y_n, \mathbf{y}; \theta) = \ell(A, Y_n, \mathbf{y}; \theta) \cdot \mathbb{1}_{\{\mathcal{A}_{R,t}\}}(Y_n, \mathbf{y})$ , and let  
 472  $\mathcal{R}_n^t$ ,  $\mathcal{R}_{n,N}^t$ , and  $\mathcal{R}_m^t$  denote the associated truncated risk functionals. Further, let  $\theta^*$  denote a  
 473 fixed parameter, to be specified later on. We decompose the generalization error into a sum of  
 474 approximation error, statistical error conditioned on the data being bounded, and truncation error that

<sup>2</sup>The performance of GPTs is known to deteriorate when the test sequence length exceeds the train sequence length; Zhang et al. [2023] conjectures this phenomenon to be the result of positional encoding.

475 leverages the tail decay of the data distribution. In more detail, we have

$$\mathcal{R}_m(\hat{\theta}) = \left(\mathcal{R}_m(\hat{\theta}) - \mathcal{R}_m^t(\hat{\theta})\right) + \left(\mathcal{R}_m^t(\hat{\theta}) - \mathcal{R}_{m,N}^t(\hat{\theta})\right) + \left(\mathcal{R}_{m,N}^t(\hat{\theta}) - \mathcal{R}_{m,N}^t(\theta^*)\right) \quad (13)$$

$$+ \left(\mathcal{R}_{m,N}^t(\theta^*) - \mathcal{R}_m^t(\theta^*)\right) + \left(\mathcal{R}_m^t(\theta^*) - \mathcal{R}_m(\theta^*)\right) + \mathcal{R}_m(\theta^*) \quad (14)$$

$$\leq \sup_{\|\theta\| \leq M} \left(\mathcal{R}_m(\theta) - \mathcal{R}_m^t(\theta)\right) + 2 \sup_{\|\theta\| \leq M} \left|\mathcal{R}_m^t(\theta) - \mathcal{R}_{m,N}^t(\theta)\right| \quad (15)$$

$$+ \left(\mathcal{R}_{m,N}^t(\hat{\theta}) - \mathcal{R}_{m,N}^t(\theta^*)\right) + \inf_{\|\theta^*\| \leq M} \mathcal{R}(\theta^*). \quad (16)$$

476 where we discarded the nonpositive term  $\left(\mathcal{R}^t(\theta^*) - \mathcal{R}(\theta^*)\right)$ . This decomposition mimics the standard  
 477 decomposition of generalization error into approximation and statistical errors, with an additional  
 478 term that arises from truncating the data. Similar techniques have recently been used in Cole and Lu  
 479 [2024] and Park et al. [2023]. There is one more technical detail to be addressed. We would like to say  
 480 that the term  $\left(\mathcal{R}_{m,N}^t(\hat{\theta}) - \mathcal{R}_{m,N}^t(\theta^*)\right)$  is nonpositive with high probability, as a consequence of the  
 481 minimality of  $\hat{\theta}$ . However, the parameter  $\hat{\theta}$  is a minimizer of the empirical risk  $\mathcal{R}_{n,N}$  corresponding  
 482 to the context length  $n$  during training, as opposed to the empirical risk  $\mathcal{R}_{m,N}$  corresponding to the  
 483 context length  $m$  during inference. However, it is easy to see that the following bound holds

$$\mathcal{R}_{m,N}^t(\hat{\theta}) - \mathcal{R}_{m,N}^t(\theta^*) \leq 2 \sup_{\|\theta\| \leq M} \left(\mathcal{R}_{m,N}^t(\theta) - \mathcal{R}_m^t(\theta)\right) + 2 \sup_{\|\theta\| \leq M} \left(\mathcal{R}_{n,N}^t(\theta) - \mathcal{R}_n^t(\theta)\right) \quad (17)$$

$$+ \sup_{\|\theta\| \leq M} \left(\mathcal{R}_m(\theta) - \mathcal{R}_m^t(\theta)\right) + \sup_{\|\theta\| \leq M} \left(\mathcal{R}_n(\theta) - \mathcal{R}_n^t(\theta)\right) \quad (18)$$

$$+ 2 \sup_{\|\theta\| \leq M} \left|\mathcal{R}_m(\theta) - \mathcal{R}_n(\theta)\right| + \left(\mathcal{R}_{n,N}^t(\hat{\theta}) + \mathcal{R}_{n,N}^t(\theta^*)\right). \quad (19)$$

484 Plugging the estimate 17 into the bound from 13 gives the final bound

$$\mathcal{R}_m(\hat{\theta}) \leq 2 \underbrace{\sup_{\|\theta\| \leq M} \left(\mathcal{R}_m - \mathcal{R}_m^t\right)(\theta) + \sup_{\|\theta\| \leq M} \left(\mathcal{R}_n - \mathcal{R}_n^t\right)(\theta)}_{\text{data truncation error}} \quad (20)$$

$$+ 4 \underbrace{\sup_{\|\theta\| \leq M} \left(\mathcal{R}_m^t - \mathcal{R}_{m,N}^t\right)(\theta) + 2 \sup_{\|\theta\| \leq M} \left(\mathcal{R}_n^t - \mathcal{R}_{n,N}^t\right)(\theta)}_{\text{statistical error}} \quad (21)$$

$$+ 2 \underbrace{\sup_{\|\theta\| \leq M} \left|\mathcal{R}_m(\theta) - \mathcal{R}_n(\theta)\right|}_{\text{context mismatch error}} + \underbrace{\left(\mathcal{R}_{n,N}^t(\hat{\theta}) - \mathcal{R}_{n,N}^t(\theta^*)\right)}_{\leq 0 \text{ w.h.p.}} + \underbrace{\mathcal{R}_m(\theta^*)}_{\text{approx. error}} \quad (22)$$

$$= I + II + III + IV + V. \quad (23)$$

485 The plan of action is to bound term  $I$  using the tail decay of the data and term  $II$  using tools  
 486 from empirical process theory; term  $III$  is controlled via Lemma 12; term  $IV$  can be shown to be  
 487 nonpositive with high-probability, and term  $V$ , the approximation error, is controlled by Proposition  
 488 1.

489 **Step 2 - bounding the truncation error:** By Lemma 7 and Example 6.2 in Wainwright [2019], when  
 490  $\mathbf{y} \sim N(0, \Sigma)$  and  $Y_n$  is the empirical covariance of iid samples from  $N(0, \Sigma)$  we have

$$P(\mathcal{A}_t^c(Y_n, \mathbf{y})) \leq \exp\left(-\frac{nt^2}{2}\right) + \exp\left(-\frac{t^2}{C\|\Sigma\|_{\text{op}}}\right)$$

for some universal constant  $C > 0$ . Therefore, for any  $\|\theta\| \leq M$ , we can apply the Cauchy-Schwarz inequality to obtain

$$\begin{aligned} \mathcal{R}_m(\theta) - \mathcal{R}_m^t(\theta) &= \mathbb{E}\|(PA^{-1}Y_m Q - A^{-1})\mathbf{y}\|^2 \cdot \mathbb{1}\{\mathcal{A}_{R,t}^c(Y_m, \mathbf{y})\} \\ &\leq \left(\mathbb{E}\|(PA^{-1}Y_m Q - A^{-1})\mathbf{y}\|^4\right)^{1/2} \cdot \mathbb{P}\left(\mathcal{A}_{R,t}^c(Y_m, \mathbf{y})\right)^{1/2} \\ &\leq c_A^2 \left(M^2 \left(\mathbb{E}\|Y_n\|_{\text{op}}^4\right)^{1/2} + 1\right) \left(\mathbb{E}\|\mathbf{y}\|^4\right)^{1/2} \cdot \sqrt{\exp\left(-\frac{mt^2}{2}\right) + \exp\left(-\frac{t^2}{C\|\Sigma\|_{\text{op}}}\right)}. \end{aligned}$$

491 This shows that the truncation error is quite mild, since  $R$  and  $t$  can be made large – in fact, we will  
 492 see that the generalization error depends only poly-logarithmically on  $R$ . Analogous bounds hold for  
 493  $\sup_{\|\theta\| \leq M} \left(\mathcal{R}_n - \mathcal{R}_n^t\right)(\theta)$ .

**Step 3 - Reduction to bounded data:** Note that by the union bound,

$$\mathcal{B}_{N,t} := \bigcap_{i=1}^N \mathcal{A}_t(Y_n^{(i)}, \mathbf{y}_i)$$

satisfies

$$\mathbb{P}(\mathcal{B}_{N,R,t}) \geq 1 - N \left( \exp\left(-\frac{nt^2}{2}\right) + \exp\left(-\frac{t^2}{C\|\Sigma\|_{\text{op}}}\right) \right).$$

Moreover, on the event  $\mathcal{B}_{N,t}$ , we have  $\ell(\cdot; \theta) = \ell^{R,t}(\cdot; \theta)$ , and hence  $\hat{\theta} = \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_N^t(\theta)$ . Therefore, if we restrict attention to the event  $\mathcal{B}_{N,R,t}$ , we may assume boundedness of the data, which is crucial to proving statistical error bounds, and the error term

$$IV = \left(\mathcal{R}_N^t(\hat{\theta}) - \mathcal{R}_N^t(\theta^*)\right)$$

494 is nonpositive by the minimality of  $\mathcal{R}_N^t(\hat{\theta})$ . For the remainder of the proof, we assume that the event  
 495  $\mathcal{B}_{N,R,t}$  holds, i.e., all expectations taken are conditioned on the event  $\mathcal{B}_{N,R,t}$ .

496 **Step 4 - bounding the statistical error:** The statistical error is measured by

$$\begin{aligned} &\sup_{\|\theta\| \leq M} \left| \mathcal{R}_n^t(\theta) - \mathcal{R}_{n,N}^t(\theta) \right| \\ &= \sup_{\|\theta\| \leq M} \left| \mathbb{E}_{A, Y_n, \mathbf{y}} \|(PA^{-1}Y_n Q - A^{-1})\mathbf{y}\|^2 - \frac{1}{N} \sum_{i=1}^N \|(PA_i^{-1}Y_n^{(i)} Q - A_i^{-1})\mathbf{y}_i\|^2 \right|, \end{aligned}$$

497 where the expectations over  $Y_n$  and  $\mathbf{y}$  are over truncated versions of their original distributions. By a  
 498 standard symmetrization argument, we have

$$\begin{aligned} &\sup_{\|\theta\| \leq M} \left| \mathbb{E}_{A, Y_n, \mathbf{y}} [\|(PA^{-1}Y_n Q - A^{-1})\mathbf{y}\|^2] - \frac{1}{N} \sum_{i=1}^N \|(PA_i^{-1}Y_n^{(i)} Q - A_i^{-1})\mathbf{y}_i\|^2 \right| \\ &\leq 2\mathbb{E}_{A_i, Y_n^{(i)}, \mathbf{y}_i} \mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \|(PA_i^{-1}Y_n^{(i)} Q - A_i^{-1})\mathbf{y}_i\|^2 \\ &= 2\mathbb{E}_{A_i, Y_n^{(i)}, \mathbf{y}_i} \mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \left( \|PA_i^{-1}Y_n^{(i)} Q \mathbf{y}_i\|^2 + \|A_i^{-1}\mathbf{y}_i\|^2 - 2\langle PA_i^{-1}Y_n^{(i)} Q \mathbf{y}_i, A_i^{-1}\mathbf{y}_i \rangle \right) \\ &\leq 2\mathbb{E}_{A_i, Y_n^{(i)}, \mathbf{y}_i} \mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \|PA_i^{-1}Y_n^{(i)} Q \mathbf{y}_i\|^2 \\ &+ 4\mathbb{E}_{A_i, Y_n^{(i)}, \mathbf{y}_i} \mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \langle PA_i^{-1}Y_n^{(i)} Q \mathbf{y}_i, A_i^{-1}\mathbf{y}_i \rangle, \end{aligned}$$

499 where the last inequality follows from the triangle inequality, noting that the term  $\sum_{i=1}^N \epsilon_i \|A_i^{-1}\mathbf{y}_i\|^2$   
 500 is independent of  $\theta$  and hence vanishes in the expectation over  $\epsilon_i$ . Now, define the function classes

$$\Theta_1(M) = \{(A, Y_n, \mathbf{y}) \mapsto \|PA^{-1}Y_n Q \mathbf{y}\|^2 : \|\theta\| \leq M\},$$

$$\Theta_2(M) = \{(A, Y_n, \mathbf{y}) \mapsto \langle PA^{-1}Y_n Q \mathbf{y}, A^{-1}\mathbf{y} \rangle : \|\theta\| \leq M\}.$$

501 By Dudley's integral theorem Dudley [1967], it holds that

$$\mathbb{E}_{A_i, Y_n^{(i)}, \mathbf{y}_i} \mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \|PA_i^{-1} Y_n^{(i)} Q \mathbf{y}_i\|^2 \leq \inf_{\epsilon > 0} \frac{12\sqrt{2}}{\sqrt{N}} \int_{\epsilon}^{D_1(M)} \sqrt{\log \mathcal{N}(\Theta_1(M), \|\cdot\|_N, \tau)} d\tau, \quad (24)$$

where  $\mathcal{N}(\Theta_1(M), \|\cdot\|_N, \tau)$  is the  $\tau$ -covering number of the function class  $\Theta_1(M)$  with respect to the metric induced by the empirical  $L^2$  norm  $\|F\|_N^2 = \frac{1}{N} \sum_{i=1}^N F(A_i, Y_n^{(i)}, \mathbf{y}_i)^2$  and

$$D_1(M) = \sup_{\|\theta\| \leq M} \left\| \|PA^{-1} Y_n Q \mathbf{y}\|^2 \right\|_N.$$

502 Note the bound

$$\begin{aligned} D_1(M)^2 &= \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \|PA_i^{-1} Y_n^{(i)} Q \mathbf{y}_i\|^4 \\ &\leq \frac{1}{N} \sum_{i=1}^N M^8 c_A^4 \|\Sigma\|_{\text{op}}^4 \left(1 + t + \sqrt{\frac{d}{n}}\right)^4 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^4 \end{aligned}$$

503 and hence  $D_1(M) \leq M^4 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}}\right)^2 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^2$ . Similarly, for  $\theta_1 = (P_1, Q_1)$ ,  $\theta_2 =$   
504  $(P_2, Q_2)$ , with  $\|\theta_1\|, \|\theta_2\| \leq M$ , we have

$$\begin{aligned} \|\theta_1 - \theta_2\|_N^2 &= \frac{1}{N} \sum_{i=1}^N \|(P_1 - P_2)A_i^{-1} Y_n^{(i)} (Q_1 - Q_2) \mathbf{y}_i\|^4 \\ &\leq 16M^4 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}}\right)^2 R^2 \cdot \frac{1}{N} \sum_{i=1}^N \|(P_1 - P_2)A_i^{-1} Y_n^{(i)} (Q_1 - Q_2)\|^2 \\ &\leq M^4 c_A^4 \|\Sigma\|_{\text{op}}^4 \left(1 + t + \sqrt{\frac{d}{n}}\right)^4 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^4 \cdot \max\left(\|P_1 - P_2\|_{\text{op}}^2, \|Q_1 - Q_2\|_{\text{op}}^2\right). \end{aligned}$$

This shows that the metric induced by  $\|\cdot\|_N$  is dominated by the metric  $d(\theta_1, \theta_2) = \max\left(\|P_1 - P_2\|_{\text{op}}, \|Q_1 - Q_2\|_{\text{op}}\right)$ , up to a factor of  $M^2 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}}\right)^2 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^2$ . The covering number of the set  $\{\|\theta\| \leq M\}$  in the metric  $d(\cdot, \cdot)$  is well-known, from which we conclude that

$$\log \mathcal{N}(\Theta_1(M), \|\cdot\|_N, \tau) \leq 2d^2 \log \left( M^2 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + \frac{2}{\tau}\right) \right).$$

505 Optimizing over the choice of  $\epsilon$  in Equation 24, this proves that

$$\begin{aligned} \mathbb{E}_{A_i, Y_n^{(i)}, \mathbf{y}_i} \mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \|PA_i^{-1} Y_n^{(i)} Q \mathbf{y}_i\|^2 & \quad (25) \\ = O\left(\frac{d^2 M^4 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}}\right)^2 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^2}{\sqrt{N}}\right), & \quad (26) \end{aligned}$$

where  $O(\cdot)$  omits factors that are logarithmic in  $N$ . An analogous argument proves a bound of the same order on the quantity

$$\mathbb{E}_{A_i, Y_n^{(i)}, \mathbf{y}_i} \mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \langle PA_i^{-1} Y_n^{(i)} Q \mathbf{y}_i, A_i^{-1} \mathbf{y}_i \rangle,$$

which in turn bounds the statistical error

$$\sup_{\|\theta\| \leq M} \left| \mathcal{R}_n^t(\theta) - \mathcal{R}_{n,N}^t(\theta) \right|$$



by the right-hand side of Equation 25. The same argument proves in analogous bound on the statistical error term

$$\sup_{\|\theta\| \leq M} \left| \mathcal{R}_m^t(\theta) - \mathcal{R}_{m,N}^t(\theta) \right|,$$

506 where  $n$  is replaced by  $m$  in the bound of Equation 25.

**Step 5: Bounding the context mismatch error** The context mismatch error satisfies the bound

$$\sup_{\|\theta\| \leq M} \left| \mathcal{R}_m(\theta) - \mathcal{R}_n(\theta) \right| \leq 2M^4 c_A^2 \max(\text{Tr}(\Sigma), \|\Sigma\|_{\text{op}}^2) \text{Tr}(\Sigma) \left| \frac{1}{n} - \frac{1}{m} \right|.$$

507 The proof of this fact is deferred to Lemma 12.

**Step 6 - Approximation error:** It remains to bound the approximation error term  $\mathcal{R}(\theta^*)$ . From Proposition 1, we have

$$\mathcal{R}_m(\theta^*) \leq \frac{c_A^2 \text{Tr}(\Sigma)}{m} + \frac{c_A^6 \|\Sigma^{-1}\|_{\text{op}}^2 \|\Sigma\|_{\text{op}}^6 \text{Tr}(\Sigma)}{n^2} + O\left(\frac{1}{mn}\right)$$

508 for an appropriate choice of  $\theta^*$ , where  $C_1$  and  $C_2$  depend only on the task and data distributions.  
509 Moreover, upon inspection of the proof of Proposition 1, we see that the  $\theta^* = (\mathbf{I}_d, Q_n)$  that attains  
510 this error is an  $O(1/n)$ -perturbation of the pair  $(\mathbf{I}_d, \Sigma^{-1})$ . Thus, if  $n$  is sufficiently large, we are  
511 guaranteed that  $\theta^*$  belongs in the set  $\{\|\theta\| \leq M\}$  for  $M \geq 2 \max(1, \|\Sigma^{-1}\|_{\text{op}})$ .

512 **Step 7 - Balancing error terms:** Putting everything together and applying the error decomposition  
513 from step 1, we have shown that <sup>3</sup>

$$\begin{aligned} \mathcal{R}_m(\hat{\theta}) &\lesssim c_A^2 \left( M^2 \mathbb{E}[\|Y_n\|_{\text{op}}^4]^{1/2} + 1 \right) \mathbb{E}[\|\mathbf{y}\|^4]^{1/2} \cdot \sqrt{\exp\left(-\frac{nt^2}{2}\right) + \exp\left(-\frac{t^2}{C\|\Sigma\|_{\text{op}}}\right)} \\ &\quad + \frac{d^2 M^4 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}}\right)^2 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^2}{\sqrt{N}} + \frac{2\text{Tr}(\mathbb{E}[A^{-2}]\Sigma)}{n}, \end{aligned}$$

with probability at least

$$1 - N \left( \exp\left(-\frac{nt^2}{2}\right) + \exp\left(-\frac{t^2}{C\|\Sigma\|_{\text{op}}}\right) \right).$$

For a fixed  $p > 0$ , we choose  $t$  such that

$$\left( \exp\left(-\frac{nt^2}{2}\right) + \exp\left(-\frac{t^2}{C\|\Sigma\|_{\text{op}}}\right) \right) = \frac{1}{N^{p+1}}.$$

514 It is clear that such a  $t$  satisfies  $t \lesssim \sqrt{p \log(N)}$ . For such a  $t$ , we have, omitting universal constants  
515 and  $\log(N)$  factors, that

$$\begin{aligned} \mathcal{R}_m(\hat{\theta}) &\lesssim \frac{c_A^2 \text{Tr}(\Sigma)}{m} + \frac{c_A^6 \|\Sigma^{-1}\|_{\text{op}}^2 \|\Sigma\|_{\text{op}}^6 \text{Tr}(\Sigma)}{n^2} + \sqrt{p} \frac{c_A^2 \left( M^2 \mathbb{E}[\|Y_n\|_{\text{op}}^4]^{1/2} + 1 \right) \mathbb{E}[\|\mathbf{y}\|^4]^{1/2}}{N} \\ &\quad + \frac{d^2 M^4 c_A^2 \|\Sigma\|_{\text{op}}^2}{\sqrt{N}} + M^4 c_A^2 \max(\text{Tr}(\Sigma), \|\Sigma\|_{\text{op}}^2) \text{Tr}(\Sigma) \left| \frac{1}{n} - \frac{1}{m} \right|, \quad \text{w.p.} \geq 1 - \frac{2}{N^p}. \end{aligned}$$

516 We omit the third term from the final bound, since, asymptotically, it is dominated by the fourth  
517 term.  $\square$

We now present an important preliminary result, which gives an upper bound on  $\inf_{\theta} \mathcal{R}_m(\theta)$ , the minimal risk achieved by a transformer in the infinite-task limit. To motivate our result, we

<sup>3</sup>For simplicity, we have omitted the terms from the truncation and statistical errors which depend on  $m$ , as they do not change the order of the final bound with respect to  $m$ ,  $n$ , or  $N$ .

first observe that for  $\theta = (P, Q)$ , the output of the transformer  $\text{TF}_\theta$  at a prompt  $Z$  of length  $m$  corresponding to a task matrix  $A$  is

$$\text{TF}_\theta(Z) = P \left( \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \mathbf{y}_i^T \right) Q \mathbf{y}.$$

Since  $\mathbf{x}_i = A^{-1} \mathbf{y}_i$ , we can equivalently write the prediction of the transformer as

$$\text{TF}_\theta(Z) = P A^{-1} Y_m Q \mathbf{y},$$

518 where  $Y_m = \frac{1}{m} \sum_{i=1}^m \mathbf{y}_i \mathbf{y}_i^T$  is the empirical covariance associated to the context vectors  
 519  $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$ . Note that if we set  $P = Id$  and  $Q = \Sigma^{-1}$  to be the inverse of the data covariance  
 520 matrix, then for sufficiently large  $m$  we have  $\text{TF}_\theta(Z) \approx A^{-1} \mathbf{y}$ . This suggests that the transformer can  
 521 learn to solve linear systems in a way that is extremely robust to shifts in the distribution on the task  
 522 matrices. We note that similar choices of attention matrices have been studied in the linear regression  
 523 setting (Ahn et al. [2024], Zhang et al. [2023]). Our result essentially employs the parameterization  
 524  $P = \mathbf{I}_d$  and  $Q = \Sigma^{-1}$ , but with an additional bias term to account for the fact that the sequence  
 525 length  $n$  during training may differ from the sequence length  $m$  during inference.

Before stating our result precisely, let us define  $B := \mathbb{E}_{A \sim p_A} [A^{-2}]$ . In addition, recall the weighted trace of a matrix  $K$  with respect to the covariance  $\Sigma = W \Lambda W^T$  defined by

$$\text{Tr}_\Sigma(K) := \sum_{i=1}^d \sigma_i^2 \langle K \varphi_i, \varphi_i \rangle,$$

526 where  $\sigma_1^2, \dots, \sigma_d^2$  are the eigenvalues of  $\Sigma$  and  $\varphi_i = W e_i$  are the eigenvectors. Note that the weighted  
 527 trace is independent of the choice of eigenbasis.

**Proposition 1.** *With*

$$Q_n = B \left( \frac{n+1}{n} \Sigma B + \frac{\text{Tr}_\Sigma(B)}{n} \Sigma \right)^{-1},$$

*we have*

$$\mathcal{R}_m(\mathbf{I}_d, Q_n) \leq \frac{(c_A^2 + d) \text{Tr}(\Sigma)}{m} + \frac{c_A^2 C_A^4 \|\Sigma\|_{\text{op}}^2 \|\Sigma^{-1}\|_{\text{op}}^2 \left(1 + \text{Tr}_\Sigma(B)\right)^2 \text{Tr}(\Sigma)}{n^2} + O\left(\frac{1}{mn}\right).$$

528 *Proof.* By Lemma 8, we can write  $Q_n = \Sigma^{-1} + \frac{1}{n} K$ , where

$$\|K\|_{\text{op}} \leq \|\Sigma^{-1}\|_{\text{op}} \|\Sigma\|_{\text{op}} \left(1 + \text{Tr}_\Sigma(B)\right) C_A^2. \quad (27)$$

529 It follows that

$$\begin{aligned} \mathcal{R}_m(\mathbf{I}_d, Q_n) &= \mathbb{E}_{A, Y_m} [\text{Tr}(A^{-1} (Y_m Q_n - \mathbf{I}_d) \Sigma (Q_n Y_m - \mathbf{I}_d) A^{-1})] \\ &= \mathbb{E}_{Y_m} [\text{Tr}(B (Y_m Q_n - \mathbf{I}_d) \Sigma (Q_n^T Y_m - \mathbf{I}_d))], \quad B := \mathbb{E}[A^{-2}] \\ &= \text{Tr}(B \Sigma) + \mathbb{E}_{Y_m} [\text{Tr}(B Y_m Q_n \Sigma Q_n^T Y_m)] - \text{Tr}(B \Sigma Q_n \Sigma) - \text{Tr}(B \Sigma Q_n^T \Sigma) \\ &= \text{Tr}(B \Sigma) + \text{Tr}(B \Sigma Q_n \Sigma Q_n \Sigma) - \text{Tr}(B \Sigma Q_n \Sigma) - \text{Tr}(B \Sigma Q_n^T \Sigma) \\ &\quad + \frac{1}{m} \left( \text{Tr}(B \Sigma Q_n \Sigma Q_n^T \Sigma) + \text{Tr}_\Sigma(Q_n \Sigma Q_n^T) \text{Tr}(B \Sigma) \right) \end{aligned}$$

530 where the last equality follows from Lemma 4. Writing  $Q_n = \Sigma^{-1} + \frac{1}{n} K$  and doing some simplifying  
 531 algebra, we find that

$$\begin{aligned} \mathcal{R}_m(\mathbf{I}_d, Q_n) &= \frac{1}{m} \left( \text{Tr}((B + \text{Tr}_\Sigma(\Sigma^{-1} \mathbf{I}_d) \Sigma)) + \frac{1}{n^2} \text{Tr}(B \Sigma K \Sigma K^T \Sigma) \right) + O\left(\frac{1}{mn}\right) \\ &= \frac{1}{m} \left( \text{Tr}((B + d \mathbf{I}_d) \Sigma) \right) + \frac{1}{n^2} \text{Tr}(B \Sigma K \Sigma K^T \Sigma) + O\left(\frac{1}{mn}\right), \end{aligned}$$

where we used the fact that  $\text{Tr}_\Sigma(\Sigma^{-1}) = d$ . Using the bound on the norm of  $K$  stated in Equation 27, and the fact that  $\|B\|_{\text{op}} \leq c_A^2$ , we have

$$\text{Tr}(B \Sigma K \Sigma K^T \Sigma) \leq c_A^2 C_A^4 \|\Sigma\|_{\text{op}}^2 \|\Sigma^{-1}\|_{\text{op}}^2 \left(1 + \text{Tr}_\Sigma(B)\right)^2 \text{Tr}(\Sigma).$$

Similarly, the bound

$$\text{Tr}((B + d\mathbf{I}_d)\Sigma) \leq (c_A^2 + d)\text{Tr}(\Sigma)$$

holds. We conclude that

$$\mathcal{R}_m(\mathbf{I}_d, Q_n) \leq \frac{(c_A^2 + d)\text{Tr}(\Sigma)}{m} + \frac{c_A^2 C_A^4 \|\Sigma\|_{\text{op}}^2 \|\Sigma^{-1}\|_{\text{op}}^2 (1 + \text{Tr}_\Sigma(B))^2 \text{Tr}(\Sigma)}{n^2} + O\left(\frac{1}{mn}\right).$$

532

□

533 To justify our ansatz for upper bounding the approximation error (i.e., how the matrix  $Q_n$  in Proposi-  
534 tion 1 was chosen), we introduce the following lemma.

**Lemma 1.** *The minimizer of the functional  $Q \mapsto \mathcal{R}_n(\mathbf{I}_d, Q)$  is given by*

$$Q_n = B \left( \frac{n+1}{n} \Sigma B + \frac{\text{Tr}_\Sigma(B)}{n} \Sigma \right)^{-1},$$

535 where  $B = \mathbb{E}[A^{-2}]$  and  $\text{Tr}_\Sigma(\cdot)$  denotes the  $\Sigma$ -weighted trace.

*Proof.* Let us recall the definition of the population risk functional

$$\mathcal{R}(\mathbf{I}_d, Q) = \mathbb{E} \left[ \left\| A^{-1} (Y_n Q - I) y \right\|^2 \right],$$

536 where  $Y_n := \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^T$  denotes the empirical covariance of  $\{\mathbf{y}_i\}_{i=1}^n$ . Note that, conditioned on  
537  $A$  and  $\{\mathbf{y}_i\}_{i=1}^n$ ,  $A^{-1} (Y_n Q - I) y$  is a centered Gaussian random vector with covariance  $A^{-1} (Y_n Q -$   
538  $I) \Sigma (Q Y_n - I) A^{-1}$ . In addition, since the task and data distributions are independent, we can replace  
539 the task by its expectation. It therefore holds that

$$\mathbb{E} \left[ \left\| A^{-1} (Y_n Q - I) y \right\|^2 \right] = \mathbb{E}_{Y_n} \left[ \text{Tr} \left( B (Y_n Q - I) \Sigma (Q^T Y_n - I) \right) \right].$$

Since this is a convex functional of  $Q$ , it suffices to characterize the critical point. Taking the derivative, we find that the critical point equation for the risk is

$$\nabla_Q \mathcal{R}(\mathbf{I}_d, Q) = \mathbb{E}_{Y_n} [\Sigma Q^T Y_n B Y_n + Y_n B Y_n Q \Sigma] - 2 \Sigma B \Sigma = 0.$$

Using Lemma 4 to compute the expectation, we further rewrite the critical point equation as

$$\left( \frac{n+1}{n} B \Sigma + \frac{\text{Tr}_\Sigma(B)}{n} \Sigma \right) Q + Q^T \left( \frac{n+1}{n} \Sigma B + \frac{\text{Tr}(\Sigma)}{n} \Sigma \right) = 2B.$$

540 This equation is solved by the matrix  $Q_n$  defined in the statement of the Lemma. □

## 541 C Proofs and additional results for Subsection 3.2

542 In this section, we present a proof of Theorem 2 and provide an example of the PDE recovery error  
543 bound when the spatial discretization is defined by a  $P^1$ -finite element method.

*Proof of Theorem 2.* By the triangle inequality, we have

$$\mathbb{E} \left[ \|u - \widehat{u}_d\|_{H^1(\Omega)}^2 \right] \leq 2\mathbb{E} \left[ \|u - u_d\|_{H^1(\Omega)}^2 \right] + 2\mathbb{E} \left[ \|u_d - \widehat{u}_d\|_{H^1(\Omega)}^2 \right].$$

544 Notice that  $\mathbb{E} \left[ \|\mathbf{u}_d - \widehat{\mathbf{u}}_d\|_{L^2(\Omega)}^2 \right] = \mathcal{R}_m(\widehat{\theta})$ , where  $\widehat{\theta}$  is as defined in the statement of Theorem 1. The  
545 desired estimate therefore follows, provided we can bound  $\mathbb{E} \left[ \|u_d - \widehat{u}_d\|_{H^1(\Omega)}^2 \right]$  by a multiple of

546  $\mathbb{E} \left[ \|u_d - \widehat{u}_d\|_{L^2(\Omega)}^2 \right]$ . For any  $g = \sum_{k=1}^d c_k \phi_k \in \text{span}\{\phi_k\}_{k=1}^d$ , we have

$$\begin{aligned} \|g\|_{H^1(\Omega)}^2 &= \|g\|_{L^2(\Omega)}^2 + \left\| \sum_{k=1}^d c_k \phi'_k(x) \right\|_{L^2(\Omega)}^2 \\ &= c^T (\Phi + \Phi') c \\ &= \tilde{c} (\mathbf{I}_d + \Phi^{-1/2} \Phi' \Phi^{-1/2}) \tilde{c} \\ &\leq (1 + \lambda_{\max}(\Phi^{-1/2} \Phi' \Phi^{-1/2})) \|\tilde{c}\|^2 \\ &= (1 + \lambda_{\max}(\Phi^{-1/2} \Phi' \Phi^{-1/2})) \|g\|_{L^2(\Omega)}^2, \end{aligned}$$

where  $\tilde{c} = \Phi c$ . We conclude that

$$\mathbb{E} \left[ \|u_d - \widehat{u}_d\|_{H^1(\Omega)}^2 \right] \leq (1 + \lambda_{\max}(\Phi^{-1/2} \Phi' \Phi^{-1/2})) \cdot \mathbb{E} \left[ \|u_d - \widehat{u}_d\|_{L^2(\Omega)}^2 \right] = 2 \max_{1 \leq k \leq d} \|\phi_k\|_{H^1(\Omega)}^2 \cdot \mathcal{R}_m(\widehat{\theta}),$$

547 and therefore that

$$\mathbb{E} \left[ \|u - \widehat{u}_d\|_{H^1(\Omega)}^2 \right] \lesssim \mathbb{E} \left[ \|u - u_d\|_{H^1(\Omega)}^2 \right] + (1 + \lambda_{\max}(\Phi^{-1/2} \Phi' \Phi^{-1/2})) \cdot \mathcal{R}_m(\widehat{\theta}).$$

548

□

**Example 1** (PDE recovery error with FEM discretization in 1D). Consider the elliptic PDE (1) on a unit interval  $\Omega = [0, 1]$ . Let  $\mathcal{I}_k = [(k-1)j, kh]$  for  $0 \leq k \leq d$  be the uniform mesh on  $\Omega$ , where  $h = d^{-1}$  is the mesh size. Let  $P_1^h(\Omega)$  be the linear finite element space spanned by the  $P_1$ -finite element base functions  $\{\phi_k\}_{k=0}^d$ . Let  $\mathbf{u}_h \in P_1^h(\Omega)$  denote the  $P_1$ -finite element approximation of the solution  $u$ . Suppose that Assumption 1 holds for the task distributions  $P_a, P_V$  and assume further that  $a(x) \in C^1(\Omega)$   $P_a$ -a.s and  $V \in C(\Omega)$   $P_V$ -a.s. Then by classical regularity estimates for elliptic PDEs, the solution  $u \in H^2(\Omega)$  and satisfies  $\|u\|_{H^2(\Omega)} \lesssim \|f\|_{L^2(\Omega)}$  up to a universal constant. Moreover, by Theorem 3.16 in Ern and Guermond [2004], the FEM-solution  $u_d$  satisfies the discretization error estimate

$$\|u - u_d\|_{H^1(\Omega)} \lesssim h \|u\|_{H^2(\Omega)}.$$

It follows that

$$\mathbb{E} \left[ \|u - u_d\|_{H^1(\Omega)}^2 \right] \lesssim h^2 \mathbb{E} [\|u\|_{H^2(\Omega)}^2] \lesssim h^2 \mathbb{E} [\|f\|_{L^2(\Omega)}^2] = h^2 \text{Tr}(\Sigma_f),$$

where  $\Sigma_f : L^2(\Omega) \rightarrow L^2(\Omega)$  is the covariance operator of  $f \sim P_f$ . In addition, it can be shown that for piecewise linear FEM on 1D, the stiffness and mass matrices satisfy  $\lambda_{\max}(\Phi^{-1/2} \Phi' \Phi^{-1/2}) \lesssim h^{-2}$  (see e.g. equation (2.4) of Boffi [2010]). By Theorem 2, we conclude that in the practical regime that  $m \leq n$ , the PDE recovery error of the transformer is bounded by

$$\mathbb{E} \left[ \|u - \widehat{u}_h\|_{H^1(\Omega)}^2 \right] \lesssim h^2 + \frac{1}{h^2} \left( \frac{1}{m} + \frac{C_A^4 \|\Sigma^{-1}\|_{op}^2}{n^2} + \frac{d^2 \|\Sigma^{-1}\|_{op}^4}{\sqrt{N}} \right).$$

549 Note that the terms  $\|\Sigma^{-1}\|_{op}$  and  $C_A^4$  depend on the number of Galerkin basis functions  $d$ . For the  
550 matrix  $A$  corresponding to the FEM discretization, it can be shown that  $C_A \lesssim h^{-2}$ . In addition,  
551 when the covariance operator of the random source is given by  $\Sigma_f = (-\Delta + I)^{-\alpha}$  for some  $\alpha > 0$   
552 which controls the smoothness of the source term, it follows from the inverse inequalities [Ern and  
553 Guermond, 2004, Lemma 12.1] that  $\|\Sigma^{-1}\|_{op} \lesssim h^{2\alpha}$ . Inserting this estimate to above leads to the  
554 final PDE recovery bound in terms of the mesh size  $h$

$$\mathbb{E} \left[ \|u - \widehat{u}_h\|_{H^1(\Omega)}^2 \right] \lesssim h^2 + \frac{1}{h^2 m} + \frac{1}{h^{10+4\alpha} n^2} + \frac{1}{h^{4+8\alpha} \sqrt{N}}, \quad (28)$$

555 or equivalently in terms of the number of Galerkin basis functions  $d$

$$\mathbb{E} \left[ \|u - \widehat{u}_h\|_{H^1(\Omega)}^2 \right] \lesssim \frac{1}{d^2} + \frac{d^2}{m} + \frac{d^{10+4\alpha}}{n^2} + \frac{d^{4+8\alpha}}{\sqrt{N}}. \quad (29)$$

556 Here, we have hidden all constants from the estimate of Theorem 1 that do not depend on the  
557 dimension  $d$ .

## 558 D Proofs and additional results for Subsection 3.3

559 We first state a more general version of Theorem 3, which does not assume that the pre-training task  
560 distribution is diverse relative to the downstream task distribution.

561 **Theorem 7.** *Let  $p_A$  and  $p'_A$  denote the pre-training and downstream task distributions respectively  
562 and assume both satisfy Assumption 1. Let  $\mathcal{M}_\infty(p_A)$  and  $\mathcal{M}_\infty(p'_A)$  denote the minimizers of  $\mathcal{R}_\infty$   
563 and  $\mathcal{R}'_\infty$  respectively, and let  $\hat{\theta} \in \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_{n,N}(\theta)$  denote the empirical risk minimizer. Then*

$$\mathcal{R}'_m(\hat{\theta}) \lesssim \mathcal{R}_m(\hat{\theta}) + \frac{d(p_A, p'_A)}{m} + \operatorname{dist}(\hat{\theta}, \mathcal{M}_\infty(p_A))^2 + \operatorname{dist}(\hat{\theta}, \mathcal{M}_\infty(p'_A))^2,$$

564 where  $d(p_A, p'_A)$  is a distance between the distributions  $p_A$  and  $p'_A$ , and the implicit constants depend  
565 on  $M$ ,  $\Sigma$ , and the constant  $c_A$  defined in Assumption 1.

566 Notice that Theorem 3 is a direct consequence of Theorem 7, because the assumption that  $p_A$  is  
567 diverse relative to  $p'_A$  implies that for any  $\theta$ ,  $\operatorname{dist}(\theta, \mathcal{M}_\infty(p'_A)) \leq \operatorname{dist}(\theta, \mathcal{M}_\infty(p_A))$ . The fourth  
568 term in the bound of Theorem 7, corresponding to  $\operatorname{dist}(\hat{\theta}, \mathcal{M}_\infty(p'_A))^2$ , is novel to the best of our  
569 knowledge, and it motivates the definition of task diversity. It highlights the hardness of learning  
570 general linear systems in-context, compared to learning linear regression models Zhang et al. [2023]  
571 or linear systems corresponding to diagonal matrices Chen et al. [2024b].

572 *Proof of Theorem 7.* Recall that  $\hat{\theta} \in \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_{n,N}(\theta)$  is the ERM. Let  $\theta_* = (P_*, Q_*)$  denote  
573 a projection of  $\hat{\theta}$  onto the set  $\mathcal{M}_\infty$  and let  $\theta'_* = (P'_*, Q'_*)$  denote a projection of  $\hat{\theta}$  onto  $\mathcal{M}_\infty$ . Let  
574  $\epsilon_1 = \|\hat{\theta} - \theta_*\|$  and  $\epsilon_2 = \|\hat{\theta} - \theta'_*\|$ . Then we have the error decomposition

$$\mathcal{R}'_m(\hat{\theta}) = \mathcal{R}_m(\hat{\theta}) + (\mathcal{R}'_m(\hat{\theta}) - \mathcal{R}'_m(\theta'_*)) + (\mathcal{R}_m(\theta'_*) - \mathcal{R}_m(\theta_*)) + (\mathcal{R}_m(\theta_*) - \mathcal{R}_m(\hat{\theta}))$$

575 Taking the infimum over all projections  $\theta_*$  and  $\theta'_*$  of  $\hat{\theta}$  onto  $\mathcal{M}_\infty(p_A)$  and  $\mathcal{M}_\infty(p'_A)$ , followed by  
576 the supremum over  $\hat{\theta}$  in  $\{\|\theta\| \leq M\}$ , we arrive at the bound

$$\begin{aligned} \mathcal{R}'_m(\hat{\theta}) \leq & \mathcal{R}_m(\hat{\theta}) + \sup_{\|\hat{\theta}\| \leq M} \inf_{\theta_*, \theta'_*} |\mathcal{R}_m(\theta_*) - \mathcal{R}'_m(\theta'_*)| + \sup_{\|\theta_1\|, \|\theta_2\| \leq M, \|\theta_1 - \theta_2\| \leq \epsilon_2} |\mathcal{R}_m(\theta_1) - \mathcal{R}_m(\theta_2)| \\ & + \sup_{\|\theta_1\|, \|\theta_2\| \leq M, \|\theta_1 - \theta_2\| \leq \epsilon_1} |\mathcal{R}'_m(\theta_1) - \mathcal{R}'_m(\theta_2)|. \end{aligned}$$

The second and third terms can be bounded using a simple Lipschitz continuity estimate. Note that  
for  $m$  sufficiently large and  $\theta = (P, Q)$  with  $\|\theta\| \leq M$ , we have

$$\|(PA^{-1}Y_mQ - A^{-1})\Sigma^{1/2}\|_F^2 \lesssim c_A^2(1 + \|\Sigma\|_{\text{op}}M^2)^2\operatorname{Tr}(\Sigma)$$

for any  $A \in \operatorname{supp}(p_A)$ . It follows that

$$R_m(\theta) = \mathbb{E}_{A \sim p_A, Y_m} [\|(PA^{-1}Y_mQ - A^{-1})\Sigma^{1/2}\|_F^2]$$

is  $O(c_A^2(1 + \|\Sigma\|_{\text{op}}M^2)^2\operatorname{Tr}(\Sigma))$ -Lipschitz on  $\{\|\theta\| \leq M\}$ . We therefore have

$$\sup_{\|\theta_1\|, \|\theta_2\| \leq M, \|\theta_1 - \theta_2\| \leq \epsilon_1} |\mathcal{R}_m(\theta_1) - \mathcal{R}_m(\theta_2)| \lesssim (c_A^2(1 + \|\Sigma\|_{\text{op}}M^2)^2\operatorname{Tr}(\Sigma)) \epsilon_1^2.$$

An analogous bound holds for  $\sup_{\|\theta_1\|, \|\theta_2\| \leq M, \|\theta_1 - \theta_2\| \leq \epsilon_2} |\mathcal{R}'_m(\theta_1) - \mathcal{R}'_m(\theta_2)|$ , since the test dis-  
tribution  $p'_A$  is also assumed to satisfy Assumption 1. To bound the term  $|\mathcal{R}_m(\theta_*) - \mathcal{R}'_m(\theta'_*)|$ , we  
recall by Lemma 5 that for any  $\theta = (P, Q)$ ,

$$\mathcal{R}_m(\theta) = \mathcal{R}_\infty(\theta) + \frac{1}{m} \mathbb{E}_{A \sim p_A} [\operatorname{Tr}(PA^{-1}\Sigma Q \Sigma Q^T \Sigma A^{-1}P^T) + \operatorname{Tr}_\Sigma(Q \Sigma Q^T) \operatorname{Tr}(PA^{-1}\Sigma A^{-1}P^T)]$$

and

$$\begin{aligned} \mathcal{R}'_m(\theta) = & \mathcal{R}'_\infty(\theta) + \frac{1}{m} \mathbb{E}_{A \sim p'_A} [\operatorname{Tr}(P(A')^{-1}\Sigma Q \Sigma Q^T \Sigma (A')^{-1}P^T) \\ & + \operatorname{Tr}_\Sigma(Q \Sigma Q^T) \operatorname{Tr}(P(A')^{-1}\Sigma (A')^{-1}P^T)]. \end{aligned}$$

577 In particular, since  $\theta_* \in \operatorname{argmin}_\theta \mathcal{R}_\infty(\theta)$  and  $\theta'_* \in \operatorname{argmin}_\theta \mathcal{R}'_\infty(\theta)$ , and each functional achieves 0 as  
 578 its minimum value, we have

$$\begin{aligned} |\mathcal{R}_m(\theta_*) - \mathcal{R}'_m(\theta'_*)| &\leq \frac{1}{m} \left| \mathbb{E}_{A \sim p_A} [\operatorname{Tr}(P_* A^{-1} \Sigma Q_* \Sigma Q_*^T \Sigma A^{-1} P_*^T)] \right. \\ &\quad + \operatorname{Tr}_\Sigma(Q_* \Sigma Q_*^T) \operatorname{Tr}(P_* A^{-1} \Sigma A^{-1} P_*^T)] \\ &\quad - \mathbb{E}_{A \sim p'_A} [\operatorname{Tr}(P'_*(A')^{-1} \Sigma Q'_* \Sigma (Q'_*)^T \Sigma (A')^{-1} (P'_*)^T)] \\ &\quad \left. + \operatorname{Tr}_\Sigma(Q'_* \Sigma (Q'_*)^T) \operatorname{Tr}(P'_*(A')^{-1} \Sigma (A')^{-1} (P'_*)^T) \right] \\ &=: \frac{1}{m} \left| \mathbb{E}_{A \sim p_A} [f(A; \theta_*)] - \mathbb{E}_{A' \sim p'_A} [f(A'; \theta'_*)] \right|. \end{aligned}$$

579 It follows that

$$\begin{aligned} \sup_{\|\hat{\theta}\| \leq M} \inf_{\theta_*, \theta'_*} |\mathcal{R}_m(\theta_*) - \mathcal{R}'_m(\theta'_*)| &\leq \frac{1}{m} \sup_{\|\hat{\theta}\| \leq M} \inf_{\theta_*, \theta'_*} \left| \mathbb{E}_{A \sim p_A} [f(A; \theta_*)] - \mathbb{E}_{A' \sim p'_A} [f(A'; \theta'_*)] \right| \\ &=: \frac{1}{m} d(p_A, p'_A), \end{aligned}$$

580 where, again, the infimum is taken over all  $\theta_* \in \operatorname{argmin}_{\theta \in \mathcal{M}_\infty(p_A)} \|\theta - \hat{\theta}\|^2$  and  $\theta'_* \in$   
 581  $\operatorname{argmin}_{\theta' \in \mathcal{M}_\infty(p'_A)} \|\theta' - \hat{\theta}\|^2$ . Combining the estimates for each individual term in the error de-  
 582 composition, we obtain the final bound in the statement of Theorem 7. The fact that the bound  
 583 we have obtained tends to zero as the sample size  $(m, n, N) \rightarrow \infty$  follows from examination of  
 584 each term in the estimate: the in-domain generalization error  $\mathcal{R}_m(\hat{\theta})$  tends to zero in probability by  
 585 Theorem 1, the term  $\frac{d(p_A, p'_A)}{m}$  is deterministic and tends to zero as  $m \rightarrow \infty$ , and  $\operatorname{dist}(\hat{\theta}, \mathcal{M}_\infty)$  tends  
 586 to zero as  $N$  and  $n$  tend to infinity, respectively, by Proposition 4.  $\square$

587 The discrepancy  $d(p_A, p'_A)$  defined in the proof of Theorem 3 may not be a metric, but, crucially,  
 588 it satisfies  $d(p_A, p_A) = 0$ . This ensures that the error term due to distribution shift in Theorem 3  
 589 vanishes when the pre-training and downstream tasks coincide. We give a simple proof of this fact  
 590 below.

**Lemma 2.** *Let*

$$d(p_A, p'_A) = \sup_{\|\hat{\theta}\| \leq M} \inf_{\theta_*, \theta'_*} \left| \mathbb{E}_{A \sim p_A} [f(A; \theta_*)] - \mathbb{E}_{A' \sim p'_A} [f(A'; \theta'_*)] \right|,$$

where the infimum is taken over all projections  $\theta_*$  and  $\theta'_*$  of  $\hat{\theta}$  onto the sets  $\mathcal{M}_\infty(p_A)$  and  $\mathcal{M}_\infty(p'_A)$   
 respectively, and

$$f(A; \theta) = \operatorname{Tr}(P A^{-1} \Sigma Q \Sigma Q^T \Sigma A^{-1} P^T) + \operatorname{Tr}_\Sigma(Q \Sigma Q^T) \operatorname{Tr}(P A^{-1} \Sigma A^{-1} P^T), \quad \theta = (P, Q).$$

591 Then  $d(p_A, p'_A) = 0$  if  $p_A = p'_A$ .

*Proof.* Note that we can upper bound  $d(p_A, p_A)$  by

$$d(p_A, p_A) \leq \sup_{\|\hat{\theta}\| \leq M} \inf_{\theta_*} \left| \mathbb{E}_{A \sim p_A} [f(A; \theta_*)] - \mathbb{E}_{A \sim p_A} [f(A; \theta_*)] \right|,$$

where the infimum is now taken only over all projections  $\theta_*$  of  $\hat{\theta}$  onto  $\mathcal{M}_\infty(p_A)$ . Clearly we have

$$\left| \mathbb{E}_{A \sim p_A} [f(A; \theta_*)] - \mathbb{E}_{A \sim p_A} [f(A; \theta_*)] \right| = 0$$

592 for all  $\theta_*$ , hence  $d(p_A, p_A) \leq 0$ . Since  $d(p_A, p_A)$  is clearly non-negative, we conclude that  
 593  $d(p_A, p_A) = 0$ .  $\square$

594 The next proposition gives a characterization of the minimizers of the functionals  $\mathcal{R}_\infty$  and  $\mathcal{R}'_\infty$ .  
 595 Apart from being interesting in its own right, it is a key tool to prove Theorem 4.

596 **Proposition 2.** Fix a task distribution  $p_A$  satisfying Assumption 1. Then  $\theta = (P, Q)$  is a minimizer  
 597 of  $\mathcal{R}_\infty$  if and only if  $P$  commutes with all elements of the set  $\{A_1 A_2^{-1} : A_1, A_2 \in \text{supp}(p_A)\}$  and  $Q$   
 598 is given by  $Q = \Sigma^{-1} A_0 P^{-1} A_0^{-1}$  for any  $A_0 \in \text{supp}(p_A)$ .

*Proof of Proposition 2.* Recall that

$$\mathcal{R}_\infty(\theta) = \mathbb{E}_{A \sim p_A} [\|(PA^{-1}\Sigma Q - A^{-1})\Sigma^{1/2}\|_F^2], \theta = (P, Q),$$

and  $\mathcal{M}_\infty(p_A) = \text{argmin}_\theta \mathcal{R}_\infty(\theta)$ . Let us first prove that for any  $p_A$  satisfying Assumption 1,  $\theta \in \mathcal{M}_\infty(p_A)$  if and only if  $PA^{-1}\Sigma Q = A^{-1}$  for all  $A \in \text{supp}(p_A)$ . Let us first observe that the minimum value of  $\mathcal{R}_\infty$  is 0 - this is attained, for instance, at  $P = \mathbf{I}_d$  and  $Q = \Sigma^{-1}$ . It is clear that if the equality  $PA^{-1}\Sigma Q = A^{-1}$  holds over the support of  $p_A$ , then  $\mathbb{E}_{A \sim p_A} [\|(PA^{-1}\Sigma Q - A^{-1})\Sigma^{1/2}\|_F^2] = 0$ . Conversely, suppose  $(P, Q)$  satisfies  $\mathbb{E}_{A \sim p_A} [\|(PA^{-1}\Sigma Q - A^{-1})\Sigma^{1/2}\|_F^2] = 0$ . Fixing  $A_0 \in \text{supp}(p_A)$  and  $\epsilon > 0$ , let  $p_{A, \epsilon}(A_0)$  denote the normalized restriction of  $p_A$  to the ball of radius  $\epsilon$  centered about  $A_0$ . Then the equality  $\mathbb{E}_{A \sim p_A} [\|(PA^{-1}\Sigma Q - A^{-1})\Sigma^{1/2}\|_F^2] = 0$  implies that

$$\mathbb{E}_{A \sim p_{A, \epsilon}(A_0)} [\|(PA^{-1}\Sigma Q - A^{-1})\Sigma^{1/2}\|_F^2] = 0$$

599 for each  $\epsilon > 0$ . Since  $p_{A, \epsilon}(A_0)$  converges weakly to the Dirac measure centered at  $A_0$ , we have  
 600 that  $\|(PA_0^{-1}\Sigma Q - A_0^{-1})\Sigma^{1/2}\|_F^2 = 0$ , and hence that  $PA_0^{-1}\Sigma Q = A_0^{-1}$ . As  $A_0$  was arbitrary, this  
 601 concludes the first part of the proof.

602 Now, suppose  $\theta = (P, Q)$  is a minimizer of  $\mathcal{R}_\infty$ . By the previous argument, this is equivalent to the  
 603 system of equations  $PA^{-1}\Sigma Q = A^{-1}$  holding simultaneously for all  $A \in \text{supp}(p_A)$ . In particular,  
 604 for any fixed  $A_0 \in \text{supp}(p_A)$ , the equation  $PA_0^{-1}\Sigma Q = A_0^{-1}$  can be solved for  $Q$ , yielding  
 605  $Q = \Sigma^{-1} A_0 P^{-1} A_0^{-1}$ . Since the matrix  $Q$  is constant, this implies that the function  $A \mapsto AP^{-1}A^{-1}$   
 606 is a constant on the support of  $p_A$ . We have therefore shown that the minimizers of  $\mathcal{R}_\infty$  can be  
 607 completely characterized as  $\{(P, \Sigma^{-1} A_0 P^{-1} A_0^{-1}) : P \in \mathbb{R}^{d \times d}\}$ , where  $A_0$  is any element of  
 608  $\text{supp}(p_A)$ . In addition, the fact that the function  $A \mapsto AP^{-1}A^{-1}$  is constant on the support of  $p_A$   
 609 implies that  $P$  commutes with all products of the form  $\{A_1 A_2^{-1} : A_1, A_2 \in \text{supp}(p_A)\}$ .  $\square$

610 We now give a proof of Theorem 4.

611 *Proof of Theorem 4.* 1) This is a direct corollary of Proposition 2.

612 2) Let  $\theta_* = (P_*, Q_*)$  be a minimizer of  $\mathcal{R}_\infty$ . Then Proposition 2 implies that  $P_* \in \mathcal{C}(\mathcal{S}(p_A))$ . Since  
 613 the centralizer of  $\mathcal{S}(p_A)$  is trivial by assumption, this implies that  $P_* = c\mathbf{I}_d$  for some  $c \in \mathbb{R} \setminus \{0\}$ .  
 614 Using the characterization of minimizers of  $\mathcal{R}_\infty$  derived in Proposition 2, we have that  $Q_*$  solves the  
 615 equation  $cA^{-1}\Sigma Q_* = A^{-1}$  for all  $A \in \text{supp}(p_A)$ , and therefore  $Q = c^{-1}\Sigma^{-1}$ .  $\square$

616 The proof of Theorem 4 implies that if  $\text{supp}(p_A)$  satisfies the condition that the centralizer of  
 617  $\{A_1 A_2^{-1} : A_1, A_2 \in \text{supp}(p_j)\}$  is trivial, then all minimizers of  $\mathcal{R}_\infty$  are of the form  $\{(P, Q) =$   
 618  $(c\mathbf{I}_d, c^{-1}\Sigma^{-1}) : c \neq 0\}$ . In this case, it is worth noting that the discrepancy on task distributions  
 619  $d(p_A, p'_A)$  defined in Theorem 3 admits a much simpler expression. We state this result as a Corollary  
 620 below.

**Corollary 1.** Under the assumption that the pre-training task distribution  $p_A$  satisfies the centralizer condition

$$\mathcal{C}(\{A_1 A_2^{-1} : A_1, A_2 \in \text{supp}(p_j)\}) = \{c\mathbf{I}_d : c \in \mathbb{R}\},$$

the out-of-distribution generalization error admits the more tractable expression

$$\mathcal{R}'_m(\hat{\theta}) = \mathcal{R}_m(\hat{\theta}) + \frac{(d+1) |\text{Tr}((\mathbb{E}_{A \sim p_A}[A^{-2}] - \mathbb{E}_{A' \sim p'_A}[(A')^{-2}]) \Sigma)|}{m} + \text{dist}(\hat{\theta}, \mathcal{M}_\infty(p_A))^2.$$

621 In particular, the second term, reflecting the discrepancy between  $p_A$  and  $p'_A$ , depends only on the  
 622 second moments of  $A^{-1}$  and  $(A')^{-1}$ .

*Proof.* By combining Theorems 3 and 4, we immediately derive the bound on the out-of-distribution generalization error

$$\mathcal{R}'_m(\hat{\theta}) = \mathcal{R}_m(\hat{\theta}) + \frac{d(p_A, p'_A)}{m} + \text{dist}(\hat{\theta}, \mathcal{M}_\infty(p_A))^2,$$

where the distance  $d(p_A, p'_A)$  is given by

$$d(p_A, p'_A) = |\mathcal{R}_m(\theta_*) - \mathcal{R}'_m(\theta_*)|,$$

and  $\theta_*$  is defined as the projection of  $\hat{\theta}$  onto the  $\mathcal{M}_\infty(p_A)$ . Under our assumptions, we have  $\mathcal{M}_\infty(p_A) = \{(c\mathbf{I}_d, c^{-1}\Sigma^{-1}) : c \in \mathbb{R} \setminus \{0\}\}$ , and applying Lemma 6 to compute  $\mathcal{R}_m(\theta_*)$  and  $\mathcal{R}'_m(\theta_*)$ , we obtain

$$d(p_A, p'_A) = (d+1) |\text{Tr}((\mathbb{E}_{A \sim p_A}[A^{-2}] - \mathbb{E}_{A' \sim p'_A}[(A')^{-2}]) \Sigma)|.$$

623

□

624 To conclude this section, we investigate the diversity of task distributions whose support consists of  
 625 simultaneously diagonalizable matrices. The simultaneous-diagonalizability of task matrices has been  
 626 used as a key assumption in the existing theoretical analysis of in-context learning of linear systems  
 627 (Chen et al. [2024b]) and in the in-context learning of linear dynamical systems (Sander et al. [2024]).  
 628 In addition, it is also relevant to the PDE setting: if the diffusion coefficient  $a(x)$  and potential  $V(x)$   
 629 are both constant,  $a(x) \equiv a_0$ ,  $V(x) \equiv v_0$ , then the solution operator of the corresponding elliptic  
 630 PDE is given by  $(-a_0\Delta + v_0\mathbf{I})^{-1}$ , whose diagonalization is independent of the constants  $a_0$  and  $v_0$ .  
 631 It is therefore natural to ask whether such a task distribution is diverse in the sense of Definition 1.

632 **Proposition 3.** *Let  $p_A$  and  $p'_A$  denote the pre-training and downstream task distributions, and*  
 633 *suppose that the matrices in  $\text{supp}(p_A)$  are simultaneously diagonalizable for a common orthogonal*  
 634 *matrix  $U$ . Suppose additionally that there exist matrices  $A_1, A_2 \in \text{supp}(p_A)$  and  $A'_1 A'_2 \in \text{supp}(p'_A)$*   
 635 *such that  $A_1 A_2^{-1}$  and  $A'_1 (A'_2)^{-1}$  have no repeated eigenvalues.*

- 636 1. *If  $\text{supp}(p'_A)$  is also simultaneously diagonalizable with respect to  $U$ , then  $p_A$  is diverse*  
 637 *relative to  $p'_A$ .*
- 638 2. *If there exist matrices  $A'_3, A'_4 \in \text{supp}(p'_A)$  such that  $A'_3 (A'_4)^{-1}$  is not diagonalizable with*  
 639 *respect to  $U$ , then  $p_A$  is not diverse relative to  $p'_A$ .*

640 Proposition 3 reveals that a simultaneously-diagonalizable task distribution cannot achieve out-of-  
 641 distribution generalization under arbitrary shifts in the downstream task distribution; namely the  
 642 downstream task distribution must also be simultaneously diagonalizable in the same basis. However,  
 643 it also shows that, provided the pre-training and downstream task distributions are simultaneously  
 644 diagonalizable, pre-trained transformers can generalize under arbitrary shifts on the distribution shifts  
 645 on the eigenvalues of the task matrices. This provides a precise characterization of the diversity of a  
 646 simultaneously diagonalizable task distribution.

647 Before proving Proposition 3, we first introduce a preliminary lemma.

**Lemma 3.** *Let  $p_A$  be a task distribution satisfying Assumption 1. Suppose that the support of  $p_A$  is*  
*simultaneously diagonalizable with a common orthogonal diagonalizing matrix  $U \in \mathbb{R}^{d \times d}$ . Assume*  
*in addition that there exist  $A_1, A_2 \in \text{supp}(p_A)$  such that  $A_1 A_2^{-1}$  has distinct eigenvalues. Then*  
 *$\mathcal{M}_\infty(p_A) = \Theta_{U, \Sigma}$ , where*

$$\Theta_{U, \Sigma} := \{(P, \Sigma^{-1} P^{-1}) : P = U D U^T, D = \text{diag}(\lambda_1, \dots, \lambda_d)\}.$$

648 *Proof.* By Proposition 2, a parameter  $(P, Q)$  belongs to  $\mathcal{M}_\infty(p_A)$  if and only if  $P$  commutes  
 649 with all products of the form  $\{A_i A_j^{-1} : A_i, A_j \in \text{supp}(p_A)\}$ , in which case  $Q$  is defined by  
 650  $Q = \Sigma^{-1} A_0 P^{-1} A_0^{-1}$  for any  $A_0 \in \text{supp}(p_A)$ . Let  $A_1, A_2 \in \text{supp}(p_A)$  be as defined in the statement  
 651 of the lemma. Since  $P$  and  $A_1 A_2^{-1}$  are commuting diagonalizing matrices and  $A_1 A_2^{-1}$  has no repeated  
 652 eigenvalues (Strang [2022]), they must be simultaneously diagonalizable. This implies that  $P$  is  
 653 diagonal in the basis  $U$ , and hence  $Q$  is given by  $Q = \Sigma^{-1} A_0 P^{-1} A_0^{-1} = \Sigma^{-1} P^{-1}$ . □

654 *Proof of Proposition 3.* For 1), if the support of  $p'_A$  is also simultaneously diagonalizable with respect  
 655 to  $U$ , then Lemma 3 implies that  $\mathcal{M}_\infty(p_A) = \mathcal{M}_\infty(p'_A) = \Theta_{U, \Sigma}$ , where  $\Theta_{U, \Sigma}$ , where  $\Theta_{U, \Sigma}$  is as  
 656 defined in the statement of Lemma 3. This proves that if the support of  $p'_A$  is also simultaneously  
 657 diagonalizable with respect to  $U$ , then  $p_A$  is diverse.



658 For 2), we must find a minimizer of  $\mathcal{R}_\infty$  which is not a minimizer of  $\mathcal{R}'_\infty$ . Consider the parameter  
659  $\theta = (P, \Sigma^{-1}P^{-1})$ , where  $P = UDU^T$  for  $D$  an invertible diagonal matrix with no repeated entries.  
660 By Lemma 3,  $\theta$  is a minimizer of  $\mathcal{R}_\infty$ . Let  $A'_3, A'_4 \in \text{supp}(p'_A)$  be such that  $A'_3(A'_4)^{-1}$  is not  
661 diagonalizable with respect to  $U$ . Since  $A'_3(A'_4)^{-1}$  and  $P$  are not simultaneously diagonalizable and  
662  $P$  has no repeated eigenvalues (Strang [2022]),  $P$  does not commute with  $A'_3(A'_4)^{-1}$ . By Proposition  
663 2,  $\theta$  is therefore not a minimizer of  $\mathcal{R}'_\infty$ , completing the proof.  $\square$

## 664 E Proofs for Subsection 3.4

665 We begin by stating a more formal version of Theorem 5 where the constants are more explicit.

666 **Theorem 8.** *Let  $\Sigma = W\Lambda W^T$  and  $\tilde{\Sigma} = \tilde{W}\tilde{\Lambda}\tilde{W}^T$  be two covariance matrices, let  $(\hat{P}, \hat{Q})$  be*  
667 *minimizers of the empirical risk when the in-context examples follow the distribution  $N(0, \Sigma)$  and*  
668 *take  $M > 0$  such that  $\max(\|\hat{P}\|_F, \|\hat{Q}\|_F) \leq M$ . Then*

$$\begin{aligned} \mathcal{R}_m^{\tilde{\Sigma}}(\hat{P}, \hat{Q}) &\lesssim \mathcal{R}_m^{\Sigma}(\hat{P}, \hat{Q}) + c_A^2 M^4 \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}})^2 \|\Sigma - \tilde{\Sigma}\|_{\text{op}} \\ &\quad + \frac{1}{m} \cdot c_A^2 M^4 \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}})^2 \text{Tr}(\tilde{\Sigma}) \left( \|\Sigma - \tilde{\Sigma}\|_{\text{op}} + \|\Lambda - \tilde{\Lambda}\|_1 + \|W - \tilde{W}\|_{\text{op}} \right). \end{aligned}$$

Theorem 5 then follows from Theorem 8 by bounding  $\|\Lambda - \tilde{\Lambda}\|_1 \lesssim \|\Sigma - \tilde{\Sigma}\|_{\text{op}}$ , merging the term

$$\frac{1}{m} \cdot c_A^2 M^4 \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}})^2 \text{Tr}(\tilde{\Sigma}) \left( \|\Sigma - \tilde{\Sigma}\|_{\text{op}} + \|\Lambda - \tilde{\Lambda}\|_1 \right)$$

669 into the second term, and omitting the constant factors.

670 *Proof of Theorem 8.* By the triangle inequality, we have

$$\mathcal{R}_m^{\tilde{\Sigma}}(\hat{P}, \hat{Q}) \leq \mathcal{R}_m^{\Sigma}(\hat{P}, \hat{Q}) + \sup_{\|P\|_{\text{op}}, \|Q\|_{\text{op}} \leq M} \left| \mathcal{R}_m^{\tilde{\Sigma}}(P, Q) - \mathcal{R}_m^{\Sigma}(P, Q) \right|. \quad (30)$$

671 It therefore suffices to bound the second term. From the proof of Proposition 1, we know that

$$\mathcal{R}_m^{\Sigma}(P, Q) = \mathbb{E}_A \left[ \frac{m+1}{m} \text{Tr}(PA^{-1}\Sigma Q\Sigma Q^T \Sigma A^{-1}P^T) + \frac{\text{Tr}_{\Sigma}(Q\Sigma Q^T)}{m} \text{Tr}(PA^{-1}\Sigma A^{-1}P^T) \right] \quad (31)$$

$$+ \mathbb{E}_A \left[ \text{Tr}(A^{-1}\Sigma A^{-1}) - \text{Tr}(PA^{-1}\Sigma Q\Sigma A^{-1}) - \text{Tr}(A^{-1}\Sigma Q^T \Sigma A^{-1}P^T) \right]. \quad (32)$$

672 Similarly, we have

$$\mathcal{R}_m^{\tilde{\Sigma}}(P, Q) = \mathbb{E}_A \left[ \frac{m+1}{m} \text{Tr}(PA^{-1}\tilde{\Sigma} Q\tilde{\Sigma} Q^T \tilde{\Sigma} A^{-1}P^T) + \frac{\text{Tr}_{\tilde{\Sigma}}(Q\tilde{\Sigma} Q^T)}{m} \text{Tr}(PA^{-1}\tilde{\Sigma} A^{-1}P^T) \right] \quad (33)$$

$$+ \mathbb{E}_A \left[ \text{Tr}(A^{-1}\tilde{\Sigma} A^{-1}) - \text{Tr}(PA^{-1}\tilde{\Sigma} Q\tilde{\Sigma} A^{-1}) - \text{Tr}(A^{-1}\tilde{\Sigma} Q^T \tilde{\Sigma} A^{-1}P^T) \right]. \quad (34)$$

673 We seek to bound the difference  $\left| \mathcal{R}_m^{\Sigma}(\theta) - \mathcal{R}_m^{\tilde{\Sigma}}(\theta) \right|$  by bounding the respective differences of each  
674 term appearing in the expressions for  $\mathcal{R}_m^{\Sigma}$  and  $\mathcal{R}_m^{\tilde{\Sigma}}$ . By a simple applications of Hölder's inequality  
675 and the triangle inequality, we see that

$$\begin{aligned} \mathbb{E}_A \text{Tr}(PA^{-1}(\Sigma Q\Sigma - \tilde{\Sigma} Q\tilde{\Sigma})A^{-1}) &\leq \mathbb{E}_A \|A^{-1}PA^{-1}\|_F \|\Sigma Q\Sigma - \tilde{\Sigma} Q\tilde{\Sigma}\|_F \\ &\leq c_A^2 \|P\|_F \left( \|(\Sigma - \tilde{\Sigma})Q\Sigma\|_F + \|\tilde{\Sigma}Q(\Sigma - \tilde{\Sigma})\|_F \right) \\ &\leq c_A^2 \|P\|_F \left( \|Q\Sigma\|_F + \|\tilde{\Sigma}Q\|_F \right) \|\Sigma - \tilde{\Sigma}\|_{\text{op}} \\ &\leq 2c_A^2 \|P\|_F \|Q\|_F \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}}) \|\Sigma - \tilde{\Sigma}\|_{\text{op}} \\ &= 2c_A^2 M^2 \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}}) \|\Sigma - \tilde{\Sigma}\|_{\text{op}}. \end{aligned}$$

Analogous arguments can be used to prove the bounds

$$\mathbb{E}_A \text{Tr}(A^{-1}(\Sigma Q^T \Sigma - \tilde{\Sigma} Q^T \tilde{\Sigma}) A^{-1} P^T) \leq 2c_A^2 M^2 \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}}) \|\Sigma - \tilde{\Sigma}\|_{\text{op}},$$

$$\mathbb{E}_A \text{Tr}(A^{-1}(\Sigma - \tilde{\Sigma}) A^{-1}) \leq c_A^2 \|\Sigma - \tilde{\Sigma}\|_{\text{op}}$$

676 and

$$\mathbb{E}_A \text{Tr}(P A^{-1}(\Sigma Q \Sigma Q^T \Sigma - \tilde{\Sigma} Q \tilde{\Sigma} Q^T \tilde{\Sigma}) A^{-1} P^T) \leq c_A^2 M^4 \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}})^2 \|\Sigma - \tilde{\Sigma}\|_{\text{op}}.$$

677 Notice that the term above dominates each of the preceding three terms. For the final term, we have

$$\begin{aligned} & \text{Tr}_{\Sigma}(Q \Sigma Q^T) \text{Tr}(P A^{-1} \Sigma A^{-1} P^T) - \text{Tr}_{\tilde{\Sigma}}(Q \tilde{\Sigma} Q^T) \text{Tr}(P A^{-1} \tilde{\Sigma} A^{-1} P^T) \\ & \leq \left| \text{Tr}_{\Sigma}(Q \Sigma Q^T) - \text{Tr}_{\tilde{\Sigma}}(Q \tilde{\Sigma} Q^T) \right| \left| \text{Tr}(P A^{-1} \Sigma A^{-1} P^T) \right| \\ & \quad + \left| \text{Tr}_{\tilde{\Sigma}}(Q \tilde{\Sigma} Q^T) \right| \left| \text{Tr}(P A^{-1}(\Sigma - \tilde{\Sigma}) A^{-1} P^T) \right|. \end{aligned}$$

678 By Lemma 10 and Holder's inequality, the second term satisfies

$$\left| \text{Tr}_{\tilde{\Sigma}}(Q \tilde{\Sigma} Q^T) \right| \left| \text{Tr}(P A^{-1}(\Sigma - \tilde{\Sigma}) A^{-1} P^T) \right| \leq c_A^2 M^4 \|\tilde{\Sigma}\|_{\text{op}} \text{Tr}(\tilde{\Sigma}) \cdot \|\Sigma - \tilde{\Sigma}\|_{\text{op}}.$$

679 Similarly, using Lemma 11, the first term satisfies

$$\begin{aligned} & \left| \text{Tr}_{\Sigma}(Q \Sigma Q^T) - \text{Tr}_{\tilde{\Sigma}}(Q \tilde{\Sigma} Q^T) \right| \left| \text{Tr}(P A^{-1} \Sigma A^{-1} P^T) \right| \\ & \leq c_A^2 M^4 \|\Sigma\|_{\text{op}} \left( \text{Tr}(\tilde{\Sigma}) \|\Sigma - \tilde{\Sigma}\|_{\text{op}} + \|\Sigma\|_{\text{op}} \left( \|\Lambda - \tilde{\Lambda}\|_1 + 2\text{Tr}(\tilde{\Sigma}) \|W - \tilde{W}\|_{\text{op}} \right) \right) \end{aligned}$$

680 Combining the estimates for each individual term and taking the supremum over the all  $P, Q$  with  
681 Frobenius norm bounded by  $M$  yields the final bound

$$\begin{aligned} \mathcal{R}_m^{\tilde{\Sigma}}(\hat{P}, \hat{Q}) & \lesssim \mathcal{R}_m^{\Sigma}(\hat{P}, \hat{Q}) + c_A^2 M^4 \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}})^2 \|\Sigma - \tilde{\Sigma}\|_{\text{op}} \\ & \quad + \frac{1}{m} \cdot c_A^2 M^4 \max(\|\Sigma\|_{\text{op}}, \|\tilde{\Sigma}\|_{\text{op}})^2 \text{Tr}(\tilde{\Sigma}) \left( \|\Sigma - \tilde{\Sigma}\|_{\text{op}} + \|\Lambda - \tilde{\Lambda}\|_1 + \|W - \tilde{W}\|_{\text{op}} \right). \end{aligned}$$

682 □

## 683 F Discussion on dependence of constants on dimension

684 It is important to consider the dependence of the constants appearing in Theorem 1 on the dimension  
685 of the linear system. Recall that in the PDE setting, the dimension  $d$  corresponds to the number of  
686 basis functions used in Galerkin's method, and hence the true PDE solution is only recovered in the  
687 limit  $d \rightarrow \infty$ .

Since the solution operator of the PDE is a bounded operator on  $L^2(\Omega)$ , the norm of the inverse  $A^{-1}$  is uniformly bounded in  $d$ , and hence the constant  $c_A = \sup_{A \in \text{supp}(p_A)} \|A^{-1}\|_{\text{op}}$  is dimension-independent. Similarly, constants involving the norm of the covariance  $\Sigma$  are dimension-independent, since we always have

$$\|\Sigma\|_{\text{op}} \leq \|\Sigma_f\|_{\text{op}}, \quad \text{Tr}(\Sigma) \leq \text{Tr}(\Sigma_f),$$

688 where  $\Sigma_f$  is the covariance of the source  $f$  on the infinite-dimensional space. However, the constant  
689  $C_A = \sup_{A \in \text{supp}(p_A)} \|A\|_{\text{op}}$  is unbounded as  $d \rightarrow \infty$ , because the limiting forward operator is  
690 unbounded on  $L^2(\Omega)$ . Similarly, the constant  $\|\Sigma^{-1}\|_{\text{op}}$  is unbounded as  $d \rightarrow \infty$ . The precise growth  
691 of these constants depends on the distributions on the coefficients of the PDE; as a prototypical  
692 example, we have  $\|A\|_{\text{op}} = O(d^2)$  for the Laplace operator under FEM discretization in 1D. It is thus  
693 important to consider the trade-offs between discretization and generalization error with respect to  
694 the dimension  $d$ ; this is explored in Example 1 for the specific case of FEM discretization.

## 695 G Auxiliary lemmas

696 We make frequent use of the following lemma to compute expectations of products of empirical  
697 covariance matrices.

**Lemma 4.** Let  $\{y_1, \dots, y_n\} \subseteq \mathbb{R}^d$  be iid samples from  $N(0, \Sigma)$  and assume that  $\Sigma = W\Lambda W^T$ , where  $\Lambda = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ . Let  $Y_n = \frac{1}{n} \sum_{k=1}^n y_k y_k^T$  associated to  $\{y_1, \dots, y_n\}$  and let  $K \in \mathbb{R}^{d \times d}$  denote a deterministic symmetric matrix. Then

$$\mathbb{E}[Y_n K Y_n] = \frac{n+1}{n} \Sigma K \Sigma + \frac{\text{Tr}_\Sigma(K)}{n} \Sigma,$$

698 where  $\text{Tr}_\Sigma(K) := \sum_{\ell=1}^d \sigma_\ell^2 \langle K \varphi_\ell, \varphi_\ell \rangle$  and  $\varphi_\ell := W e_\ell$  denote the eigenvectors of  $\Sigma$ .

699 *Proof.* Let us first consider the case that  $W = \mathbf{I}_d$ , so that the covariance is diagonal with entries  
700  $\sigma_1^2, \dots, \sigma_d^2$ . Observe that

$$\begin{aligned} \mathbb{E}[(Y_n K Y_n)_{ij}] &= \mathbb{E} \left[ \sum_{\ell, \ell'=1}^d \frac{1}{n^2} \left( \sum_{k \neq k'} \langle y_k, e_i \rangle \langle y_{k'}, e_j \rangle \langle y_k, e_\ell \rangle \langle y_{k'}, e_{\ell'} \rangle K_{\ell, \ell'} \right. \right. \\ &\quad \left. \left. + \sum_{k=1}^n \langle e_i, y_k \rangle \langle e_j, y_k \rangle \langle e_\ell, y_k \rangle \langle e_{\ell'}, y_k \rangle K_{\ell, \ell'} \right) \right]. \end{aligned}$$

701 When  $i \neq j$ , we compute that

$$\sum_{\ell, \ell'=1}^d \mathbb{E} \left[ \langle y_k, e_i \rangle \langle y_{k'}, e_j \rangle \langle y_k, e_\ell \rangle \langle y_{k'}, e_{\ell'} \rangle K_{\ell, \ell'} \right] = \sigma_i^2 \sigma_j^2 K_{i, j}$$

and

$$\sum_{\ell, \ell'=1}^d \mathbb{E} \left[ \langle y_k, e_i \rangle \langle y_k, e_j \rangle \langle y_k, e_\ell \rangle \langle y_k, e_{\ell'} \rangle K_{\ell, \ell'} \right] = 2\sigma_i^2 \sigma_j^2 K_{i, j}.$$

On the other hand, for  $i = j$ , we have

$$\sum_{\ell, \ell'=1}^d \mathbb{E} \left[ \langle y_k, e_i \rangle \langle y_{k'}, e_i \rangle \langle y_k, e_\ell \rangle \langle y_{k'}, e_{\ell'} \rangle K_{\ell, \ell'} \right] = \sigma_i^4 K_{i, i}$$

702 and

$$\sum_{\ell, \ell'=1}^d \mathbb{E} \left[ \langle y_k, e_i \rangle^2 \langle y_k, e_\ell \rangle \langle y_k, e_{\ell'} \rangle K_{\ell, \ell'} \right] = 2\sigma_i^4 K_{i, i} + \sigma_i^2 \sum_{\ell=1}^d \sigma_\ell^2 K_{\ell, \ell}.$$

Putting everything together, we have shown that

$$\mathbb{E}(Y_n K Y_n)_{i, j} = \frac{n+1}{n} \sigma_i^2 \sigma_j^2 K_{i, j} + \delta_{ij} \cdot \frac{\text{Tr}_\Sigma(K)}{n} \sigma_i^2.$$

703 The result then follows since  $(\Sigma K \Sigma)_{i, j} = \sigma_i^2 \sigma_j^2 K_{i, j}$ . For general covariance  $\Sigma = W\Lambda W^T$ , we  
704 have  $Y_n K Y_n = W(Z_n W^T K W Z_n) W^T$ , where  $Z_n$  is the empirical covariance matrix associated to  
705  $\{W^T y_1, \dots, W^T y_n\}$ . Noting that  $W^T y \sim N(0, \Lambda)$  for  $y \sim N(0, \Sigma)$ , we can apply the above result  
706 to  $W^T K W$ :

$$\begin{aligned} \mathbb{E}[Y_n K Y_n] &= W \mathbb{E}[Z_n (W^T K W) Z_n] W^T \\ &= W \left( \frac{n+1}{n} \Lambda W^T K W \Lambda + \frac{\text{Tr}_\Sigma(K)}{n} \Lambda \right) W^T \\ &= \frac{n+1}{n} \Sigma K \Sigma + \frac{\text{Tr}_\Sigma(K)}{n} \Sigma. \end{aligned}$$

707 □

708 We quickly put Lemma 4 to work to give a tractable expression for the population risk.

709 **Lemma 5.** For  $\theta = (P, Q)$ , we have

$$\begin{aligned} \mathcal{R}_n(\theta) &:= \mathbb{E}_{A, Y_n} [\| (PA^{-1} Y_n Q - A^{-1}) \Sigma^{1/2} \|_F^2] = \mathbb{E}_A [\| (PA^{-1} \Sigma Q - A^{-1}) \Sigma^{1/2} \|_F^2] \\ &\quad + \frac{1}{n} \mathbb{E}_A \left[ \text{Tr}(PA^{-1} \Sigma Q \Sigma Q^T \Sigma A^{-1} P^T) + \text{Tr}_\Sigma(Q \Sigma Q^T) \text{Tr}(PA^{-1} \Sigma A^{-1} P^T) \right]. \end{aligned}$$

710 *Proof.* This follows from a direct computation of the expectation with respect to  $Y_n$  :

$$\begin{aligned}
& \mathbb{E}_{A, Y_n} [\| (PA^{-1}Y_nQ - A^{-1})\Sigma^{1/2} \|_F^2] = \mathbb{E}_{A, Y_n} [\text{Tr}((PA^{-1}Y_nQ - A^{-1})\Sigma(Q^TY_nA^{-1}P^T - A^{-1}))] \\
& = \mathbb{E}_{A, Y_n} [\text{Tr}(A^{-1}\Sigma A^{-1} + PA^{-1}Y_nQ\Sigma Q^TY_nA^{-1}P^T - PA^{-1}Y_nQ\Sigma A^{-1} - A^{-1}\Sigma Q^TY_nA^{-1}P^T)] \\
& = \mathbb{E}_A [\text{Tr}(A^{-1}\Sigma A^{-1} - PA^{-1}\Sigma Q\Sigma A^{-1} - A^{-1}\Sigma Q^T\Sigma A^{-1}P^T) \\
& \quad + \mathbb{E}_{A, Y_n} [\text{Tr}(PA^{-1}Y_nQ\Sigma Q^TY_nA^{-1}P^T)]] \\
& = \mathbb{E}_A [\text{Tr}(A^{-1}\Sigma A^{-1} - PA^{-1}\Sigma Q\Sigma A^{-1} - A^{-1}\Sigma Q^T\Sigma A^{-1}P^T) \\
& \quad + \frac{n+1}{n} \mathbb{E}_A [\text{Tr}(PA^{-1}\Sigma Q\Sigma Q^T\Sigma A^{-1}P^T)]] + \frac{1}{n} \mathbb{E}_A [\text{Tr}_\Sigma(Q\Sigma Q^T)\text{Tr}(PA^{-1}\Sigma A^{-1}P^T)] \\
& = \mathbb{E}_A [\| (PA^{-1}\Sigma Q - A^{-1})\Sigma^{1/2} \|_F^2] \\
& \quad + \frac{1}{n} \mathbb{E}_A [\text{Tr}(PA^{-1}\Sigma Q\Sigma Q^T\Sigma A^{-1}P^T) + \text{Tr}_\Sigma(Q\Sigma Q^T)\text{Tr}(PA^{-1}\Sigma A^{-1}P^T)],
\end{aligned}$$

711 where we used Lemma 4 to compute the expectation over  $Y_n$  in the second-to-last line.  $\square$

712 It will also be useful to derive a simpler expression for the population risk  $\mathcal{R}_m(\theta)$  when  $\theta$  belongs to  
713 the set  $\Theta_\Sigma = \{(c\mathbf{I}_d, c^{-1}\Sigma^{-1}) : c \in \mathbb{R} \setminus \{0\}\}$ .

**Lemma 6.** *Let  $P = c\mathbf{I}_d$ ,  $Q = c^{-1}\Sigma^{-1}$  for  $c \in \mathbb{R} \setminus \{0\}$ . Then*

$$\mathcal{R}_m(\theta) = \frac{d+1}{n} \mathbb{E}_A [\text{Tr}(A^{-1}\Sigma A^{-1})].$$

714 *Proof.* Using Lemma 4 to compute the expectations defining  $\mathcal{R}_m$ , we have

$$\begin{aligned}
\mathcal{R}_m(\theta) & = \mathbb{E}_A [\text{Tr}(A^{-1}\Sigma A^{-1} - PA^{-1}\Sigma Q\Sigma A^{-1} - A^{-1}\Sigma Q^T\Sigma A^{-1}P^T) \\
& \quad + \frac{n+1}{n} \mathbb{E}_A [\text{Tr}(PA^{-1}\Sigma Q\Sigma Q^T\Sigma A^{-1}P^T)]] + \frac{1}{n} \mathbb{E}_A [\text{Tr}_\Sigma(Q\Sigma Q^T)\text{Tr}(PA^{-1}\Sigma A^{-1}P^T)].
\end{aligned}$$

Since  $P = c\mathbf{I}_d$  and  $Q = c^{-1}\Sigma^{-1}$ , we have that  $PA^{-1}\Sigma Q\Sigma A^{-1}$ ,  $A^{-1}\Sigma Q^T\Sigma A^{-1}P^T$ , and  $PA^{-1}\Sigma Q\Sigma Q^T\Sigma A^{-1}P^T$  are all equal to  $A^{-1}\Sigma A^{-1}$ , and

$$\mathbb{E}_A \text{Tr}_\Sigma(Q\Sigma Q^T)\text{Tr}(PA^{-1}\Sigma A^{-1}P^T) = \mathbb{E}_A \text{Tr}_\Sigma(\Sigma^{-1})\text{Tr}(A^{-1}\Sigma A^{-1}).$$

Therefore, after some algebra, the population risk simplifies to

$$\mathcal{R}_m(\theta) = \frac{1 + \text{Tr}_\Sigma(\Sigma^{-1})}{n} \mathbb{E}_A [\text{Tr}(A^{-1}\Sigma A^{-1})].$$

715 Noting that  $\text{Tr}_\Sigma(\Sigma^{-1}) = d$ , we conclude the expression for  $\mathcal{R}_m(\theta)$  as stated in the lemma.  $\square$

716 We quote the following result from Theorem 2.1 of Rudelson and Vershynin [2013].

**Lemma 7.** *[Gaussian concentration bound] Let  $y \sim N(0, \Sigma)$ . Then*

$$\mathbb{P} \left\{ \|y\| \geq \sqrt{\text{Tr}(\Sigma)} + t \right\} \leq 2 \exp \left( - \frac{t^2}{C \|\Sigma\|_{op}} \right),$$

717 where  $C > 0$  is a constant independent of  $\Sigma$  and  $d$ .

718 We use the following result to control the error between  $Q_n$  and  $\Sigma^{-1}$ .

719 **Lemma 8.**

Let  $Q_n = B \left( \frac{n+1}{n} B\Sigma + \frac{\text{Tr}_\Sigma(B)}{n} \Sigma \right)^{-1}$  be as defined in Lemma 1. Assume that  $n$  satisfies

$$\frac{\|\Sigma^{-1}\|_{op} \left\| \Sigma \left( \mathbf{I}_d + \text{Tr}_\Sigma(B)B^{-1} \right) \right\|_{op}}{n} \leq \frac{1}{2}.$$

Then we can write

$$Q_n = \Sigma^{-1} + \frac{1}{n} \mathcal{E}_1,$$

where  $\mathcal{E}_1$  satisfies

$$\|\mathcal{E}_1\| \lesssim \|\Sigma^{-1}\|_{op} \|\Sigma\|_{op} \left( 1 + \text{Tr}_\Sigma(B) \right) C_A^2.$$

720 *Proof.* Using some algebra, we find

$$\begin{aligned} Q_n &= B \left( \frac{n+1}{n} B \Sigma + \frac{\text{Tr}_\Sigma(B)}{n} \Sigma \right)^{-1} \\ &= \left( \frac{n+1}{n} \Sigma + \frac{\text{Tr}_\Sigma(B)}{n} \Sigma B^{-1} \right)^{-1} \\ &= \left( \Sigma + \frac{1}{n} \Sigma (\mathbf{I}_d + \text{Tr}_\Sigma(B) B^{-1}) \right)^{-1}. \end{aligned}$$

By Lemma 9, we have

$$\|Q_n - \Sigma^{-1}\|_{\text{op}} \leq \|\Sigma^{-1}\|_{\text{op}} \cdot \frac{\epsilon^*}{1 - \epsilon^*},$$

where

$$\epsilon^* = \frac{\|\Sigma^{-1}\|_{\text{op}} \left\| \Sigma (\mathbf{I}_d + \text{Tr}_\Sigma(B) B^{-1}) \right\|_{\text{op}}}{n}.$$

This gives the final bound

$$\|Q_n - \Sigma^{-1}\|_{\text{op}} \lesssim \frac{\|\Sigma^{-1}\|_{\text{op}} \left\| \Sigma (\mathbf{I}_d + \text{Tr}_\Sigma(B) B^{-1}) \right\|_{\text{op}}}{n} \leq \frac{\|\Sigma^{-1}\|_{\text{op}} \|\Sigma\|_{\text{op}} (1 + \text{Tr}_\Sigma(B) \|B^{-1}\|_{\text{op}})}{n},$$

721 Here, we used the bound  $\frac{\epsilon}{1-\epsilon} \lesssim \epsilon$  which holds for  $\epsilon$  sufficiently small; in particular, for  $\epsilon \in (0, 1/2)$ ,  
722 we have  $\frac{\epsilon}{1-\epsilon} \leq 2\epsilon$ .  $\square$

723 The following result, used to bound the inverse of a perturbed matrix, is a standard application of  
724 matrix power series.

**Lemma 9.** *Suppose that  $A$  is an invertible  $d \times d$  matrix and  $D \in \mathbb{R}^{d \times d}$  satisfies  $\|D\|_{\text{op}} \leq \frac{\epsilon}{\|A^{-1}\|_{\text{op}}}$  for some  $\epsilon < 1$ . Then*

$$\|(A + D)^{-1} - A^{-1}\|_{\text{op}} \leq \|A^{-1}\|_{\text{op}} \cdot \frac{\epsilon}{1 - \epsilon}.$$

*Proof.* Note that  $A + D = (\mathbf{I}_d + DA^{-1})A$ . Under our assumption on  $D$ , we have  $\|DA^{-1}\|_{\text{op}} \leq \|D\|_{\text{op}} \|A^{-1}\|_{\text{op}} < 1$ , which implies the series expansion

$$(I + DA^{-1})^{-1} = \sum_{k=0}^{\infty} (-DA^{-1})^k.$$

725 It follows that

$$\begin{aligned} (A + D)^{-1} &= \left( (I + DA^{-1})A \right)^{-1} \\ &= A^{-1} (I + DA^{-1})^{-1} \\ &= A^{-1} \sum_{k=0}^{\infty} (-DA^{-1})^k. \end{aligned}$$

726 In turn, this gives the bound

$$\begin{aligned} \|(A + D)^{-1} - A^{-1}\|_{\text{op}} &= \left\| A^{-1} \sum_{k=1}^{\infty} (-DA^{-1})^k \right\|_{\text{op}} \\ &\leq \|A^{-1}\|_{\text{op}} \sum_{k=1}^{\infty} \|DA^{-1}\|_{\text{op}}^k \\ &\leq \|A^{-1}\|_{\text{op}} \sum_{k=1}^{\infty} \epsilon^k \\ &= \|A^{-1}\|_{\text{op}} \frac{\epsilon}{1 - \epsilon}. \end{aligned}$$

727  $\square$

Recall that for a positive definite matrix  $\Sigma = W\Lambda W^T$  and a symmetric matrix  $K$ ,

$$\text{Tr}_\Sigma(K) = \sum_{i=1}^d \sigma_i^2 \langle K\varphi_i, \varphi_i \rangle,$$

728 where  $\sigma_1^2, \dots, \sigma_d^2$  are the eigenvalues of  $\Sigma$  and  $\varphi_i = We_i$  are the eigenvectors of  $\Sigma$ .

**Lemma 10.** *For any symmetric matrix  $K$ , we have*

$$\text{Tr}_\Sigma(K) \leq \|K\|_{\text{op}} \text{Tr}(\Sigma).$$

*Proof.* For each  $1 \leq i \leq d$ , we have  $\langle K\varphi_i, \varphi_i \rangle \leq \|K\varphi_i\| \|\varphi_i\| \leq \|K\|_{\text{op}}$ . Therefore,

$$\text{Tr}_\Sigma(K) = \sum_{i=1}^d \sigma_i^2 \langle K\varphi_i, \varphi_i \rangle \leq \|K\|_{\text{op}} \sum_{i=1}^d \sigma_i^2 = \|K\|_{\text{op}} \text{Tr}(\Sigma).$$

729

□

730 In order to prove Theorem 5, we also need the following stability bound of  $\text{Tr}_\Sigma(K)$  with respect to  
731 perturbations of both  $\Sigma$  and  $K$ .

**Lemma 11.** *Let  $\Sigma = W\Lambda W^T$  and  $\tilde{\Sigma} = \tilde{W}\tilde{\Lambda}\tilde{W}^T$  be two symmetric positive definite matrices and  $K, \tilde{K}$  two symmetric matrices, let  $\{\sigma_i^2\}_{i=1}^d$  and  $\{\tilde{\sigma}_i^2\}_{i=1}^d$  be the respective eigenvalues of  $\Sigma$  and  $\tilde{\Sigma}$  and let  $\{\varphi_i\}_{i=1}^d$  and  $\{\tilde{\varphi}_i\}_{i=1}^d$  be the respective eigenvectors. Then*

$$\left| \text{Tr}_\Sigma(K) - \text{Tr}_{\tilde{\Sigma}}\tilde{K} \right| \leq \text{Tr}(\tilde{\Sigma})\|K - \tilde{K}\|_{\text{op}} + \|K\|_{\text{op}} \left( \|\Lambda - \tilde{\Lambda}\|_1 + 2\text{Tr}(\tilde{\Sigma})\|W - \tilde{W}\|_{\text{op}} \right).$$

732 *Proof.* We have

$$\text{Tr}_\Sigma(K) - \text{Tr}_{\tilde{\Sigma}}(\tilde{K}) \leq \left| \text{Tr}_\Sigma(K) - \text{Tr}_{\tilde{\Sigma}}(K) \right| + \left| \text{Tr}_{\tilde{\Sigma}}(K - \tilde{K}) \right|. \quad (35)$$

733

□

The second term in 35 can be bounded by an application of Lemma 10, which yields

$$\left| \text{Tr}_{\tilde{\Sigma}}(K - \tilde{K}) \right| \leq \text{Tr}(\tilde{\Sigma})\|K - \tilde{K}\|_{\text{op}}.$$

734 To bound the first term in 35, we first use the estimate

$$\left| \text{Tr}_\Sigma(K) - \text{Tr}_{\tilde{\Sigma}}(K) \right| \leq \left| \sum_{i=1}^d (\sigma_i^2 - \tilde{\sigma}_i^2) \langle K\varphi_i, \varphi_i \rangle \right| + \left| \sum_{i=1}^d \tilde{\sigma}_i^2 \left( \langle K(\varphi_i - \tilde{\varphi}_i), \varphi_i \rangle + \langle K\tilde{\varphi}_i, \varphi_i - \tilde{\varphi}_i \rangle \right) \right|.$$

735 The first term above can be bounded by

$$\left| \sum_{i=1}^d (\sigma_i^2 - \tilde{\sigma}_i^2) \langle K\varphi_i, \varphi_i \rangle \right| \leq \|K\|_{\text{op}} \cdot \sum_{i=1}^d |\sigma_i^2 - \tilde{\sigma}_i^2| = \|K\|_{\text{op}} \cdot \|\Lambda - \tilde{\Lambda}\|_1. \quad (36)$$

To bound the second term in 36, note that for any  $1 \leq i \leq d$ , we have

$$\langle K(\varphi_i - \tilde{\varphi}_i), \varphi_i \rangle \leq \|K\|_{\text{op}} \|\varphi_i - \tilde{\varphi}_i\| \leq \|K\|_{\text{op}} \|W - \tilde{W}\|_{\text{op}},$$

736 and similarly  $\langle K\tilde{\varphi}_i, \varphi_i - \tilde{\varphi}_i \rangle \leq \|K\|_{\text{op}} \|W - \tilde{W}\|_{\text{op}}$ . It therefore holds that

$$\left| \sum_{i=1}^d \tilde{\sigma}_i^2 \left( \langle K(\varphi_i - \tilde{\varphi}_i), \varphi_i \rangle + \langle K\tilde{\varphi}_i, \varphi_i - \tilde{\varphi}_i \rangle \right) \right| \leq 2\|K\|_{\text{op}} \text{Tr}(\tilde{\Sigma}) \|W - \tilde{W}\|_{\text{op}}.$$

737 Combining all terms yields the final estimate

$$\left| \text{Tr}_\Sigma(K) - \text{Tr}_{\tilde{\Sigma}}\tilde{K} \right| \leq \text{Tr}(\tilde{\Sigma})\|K - \tilde{K}\|_{\text{op}} + \|K\|_{\text{op}} \left( \|\Lambda - \tilde{\Lambda}\|_1 + 2\text{Tr}(\tilde{\Sigma})\|W - \tilde{W}\|_{\text{op}} \right).$$

738 The following lemma bounds the 'context mismatch error', which arises in the proof of Theorem 1.

**Lemma 12.** *The bound*

$$\sup_{\|\theta\| \leq M} \left| \mathcal{R}_m(\theta) - \mathcal{R}_n(\theta) \right| \leq 2M^4 c_A^2 \max(\text{Tr}(\Sigma), \|\Sigma\|_{\text{op}}^2) \text{Tr}(\Sigma) \left| \frac{1}{n} - \frac{1}{m} \right|$$

739 *holds.*

740 *Proof.* Denote  $\theta = (P, Q)$ . Recall that, as a direct consequence of Lemma 5, we have

$$\begin{aligned} \mathcal{R}_n(\theta) &= \mathbb{E}_A \left[ \text{Tr}(A^{-1} \Sigma A^{-1}) - \text{Tr}(P A^{-1} \Sigma Q \Sigma A^{-1}) - \text{Tr}(A^{-1} \Sigma Q^T \Sigma A^{-1} P^T) \right. \\ &\quad \left. + \frac{n+1}{n} \text{Tr}(P A^{-1} \Sigma Q \Sigma Q^T \Sigma A^{-1} P^T) + \frac{\text{Tr}_\Sigma(Q \Sigma Q^T)}{n} \text{Tr}(P A^{-1} \Sigma A^{-1} P^T) \right], \end{aligned}$$

741 An analogous expression holds for  $\mathcal{R}_m(\theta)$ . Therefore, for  $\theta$  satisfying  $\|\theta\| = \max(\|P\|_{\text{op}}, \|Q\|_{\text{op}}) \leq$   
742  $M$ , we have the bound

$$\begin{aligned} \left| \mathcal{R}_m(\theta) - \mathcal{R}_n(\theta) \right| &= \left| \frac{1}{n} - \frac{1}{m} \right| \left| \mathbb{E}_A \left[ \text{Tr}(P A^{-1} \Sigma Q \Sigma Q^T \Sigma A^{-1} P^T) + \text{Tr}_\Sigma(Q \Sigma Q^T) \text{Tr}(P A^{-1} \Sigma A^{-1} P^T) \right] \right| \\ &\leq \left| \frac{1}{n} - \frac{1}{m} \right| \cdot 2M^4 c_A^2 \max(\text{Tr}(\Sigma), \|\Sigma\|_{\text{op}}^2) \text{Tr}(\Sigma). \end{aligned}$$

743

□

744 The following lemma is an adaptation of Wald's consistency theorem of M-estimators [Van der Vaart,  
745 2000, Theorem 5.14]. We use it to prove the convergence in probability of empirical risk minimizers  
746 to population risk minimizers.

**Lemma 13.** *Let  $\theta \in \mathbb{R}^m$ ,  $x \in \mathbb{R}^d$ , and suppose  $\ell(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^m \rightarrow [0, \infty)$  is lower semicontinuous in  $\theta$ . Let  $m_0 = \min_{\theta} \mathbb{E}[\ell(x, \theta)]$  for some fixed distribution on  $x$ , and let  $\Theta_0 = \text{argmin}_{\theta} \mathbb{E}[\ell(x, \theta)]$ . Let  $\{\theta_N\}_{N \in \mathbb{N}}$  be a collection of estimators such that  $\sup_N \|\theta_N\| < \infty$  and*

$$m_0 - \mathbb{E}_N[\ell(x, \theta_0)] = o_P(1)$$

747 *Then  $\text{dist}(\theta_N, \Theta_0) \xrightarrow{P} 0$ .*

**Proposition 4.** *For any sequence  $\{\hat{\theta}_{n,N}\}_{n,N \in \mathbb{N}}$  of minimizers of the empirical risk  $\mathcal{R}_{n,N}$  with  $\sup_N \|\hat{\theta}_{n,N}\| < \infty$  for all  $n$ , we have*

$$\lim_{n \rightarrow \infty} \lim_{N \rightarrow \infty} \text{dist}(\hat{\theta}_{n,N}, \mathcal{M}_{\infty}) = 0, \text{ in probability.}$$

*Proof.* For each fixed  $n \in \mathbb{N}$ , we can apply Lemma 13 to the empirical risk minimizer  $\hat{\theta}_{n,N}$ . In this context, the condition of the lemma amounts to the condition that  $\mathcal{R}_n(\theta_*) - \mathcal{R}_{n,N}(\hat{\theta}_{n,N}) = o_P(1)$ , for any  $\theta_* \in \text{argmin}_{\theta} \mathcal{R}_n$ , which is satisfied since

$$\mathcal{R}_n(\theta_*) - \mathcal{R}_{n,N}(\hat{\theta}_{n,N}) = \left( \mathcal{R}_n(\theta_*) - \mathcal{R}_{n,N}(\theta_*) \right) + \left( \mathcal{R}_{n,N}(\theta_*) - \mathcal{R}_{n,N}(\hat{\theta}_{n,N}) \right).$$

The first term tends to zero in probability by the law of large numbers, and the second term is non-negative by the minimality of  $\hat{\theta}_{n,N}$ . This proves that

$$\lim_{N \rightarrow \infty} \text{dist}(\hat{\theta}_{n,N}, \mathcal{M}_n) = 0, \text{ in probability,}$$

where  $\mathcal{M}_n = \text{argmin}_{\theta} \mathcal{R}_n(\theta)$ . Consequently, since  $\mathcal{R}_n$  and  $\mathcal{R}_{\infty}$  are polynomials in  $\theta$  such that the coefficients of  $\mathcal{R}_n$  converge to the coefficients of  $\mathcal{R}_{\infty}$  as  $n \rightarrow \infty$ , we have by the triangle inequality that

$$\lim_{n \rightarrow \infty} \lim_{N \rightarrow \infty} \text{dist}(\hat{\theta}_{n,N}, \mathcal{M}_{\infty}) = 0, \text{ in probability.}$$

748

□

## 749 H Experimental setup

### 750 H.1 In-domain generalization

We recapitulate the experimental set-up described in Subsection 4.1 for our in-domain experiments. We consider the one dimensional elliptic PDE  $(-\Delta + V(x))u(x) = f(x)$  on  $\Omega = [0, 1]$  with Dirichlet boundary condition. We assume that the source term is a Gaussian white noise, i.e.  $f = N(0, \mathbb{I})$ , where  $\mathbb{I}$  denotes the identity operator. We discretize the PDE using Galerkin projection onto the sine basis  $\phi_k(x) = \sin(k\pi x)$ ,  $k \in \{1, \dots, d\}$ . Furthermore, we assume that the potential  $V$  is uniform random field that is obtained by dividing the domain into  $2d + 1$  sub-intervals and in each cell independently, the potential  $V$  takes values uniformly in  $[1, 2]$ . This leads to the linear system  $A\mathbf{u} = \mathbf{f}$ , where  $\mathbf{f} \sim N(0, \mathbf{I}_d)$  and

$$A_{ij} = k^2\pi^2\delta_{ij} + \langle \phi_i, V\phi_j \rangle_{L^2}.$$

751 The prompts used for pre-training are then built on observations of the form  
752  $((\mathbf{f}_1, A^{-1}\mathbf{f}_1), \dots, (\mathbf{f}_n, A^{-1}\mathbf{f}_n))$ .

### 753 H.2 Out-of-domain generalization

754 For out-of-domain generalization, we consider the PDE defined by  $-\nabla \cdot (a(x)\nabla u(x)) + V(x)u(x) =$   
755  $f(x)$  on  $[0, 1]$  with Dirichlet boundary conditions.

756 **Task shifts:** During both training and inference, we assume that  $f$  is a centered Gaussian with  
757 covariance operator defined by  $(-\Delta + c\mathbb{I})^{-\beta}$  for some fixed  $c, \beta > 0$ . We parameterize  $a(x)$  as a *log-*  
758 *normal random field*, i.e., we write  $a(x) = e^{b(x)}$ , where  $b(x)$  is sampled from an infinite-dimensional  
759 Gaussian measure  $N(0, C_{\alpha, \tau})$ , where  $C_{\alpha} = (-\Delta + \tau\mathbb{I})^{-\alpha}$ . The parameter  $\alpha$  governs the smoothness  
760 of the field. During training, we set  $\alpha = 3, \tau = 5$ , and during inference we use  $\alpha = 1, 2, 4$ . For  $V$ ,  
761 we assume during training that  $V$  is piecewise constant, and the constant values are iid according to  
762 the uniform distribution  $U(1, 2)$ . During inference, we shift the distribution on the pieces of  $V$  to  
763  $U(3, 4), U(5, 10)$ , and  $U(10, 20)$ .

764 **Covariate shifts:** We train the model to solve the PDE (1), where the source term is defined by a  
765 Gaussian measure  $N(0, C)$  for  $C = (-\Delta + c\mathbb{I})^{-\beta}$ , where  $c = \beta = 1$ . Then, at inference, we consider  
766 solving the same PDE, but where the source term is defined by  $N(0, 3C)$  or  $N(0, 5C)$ ; see Figure 3:  
767 C. We also consider covariate shifts defined by changing the parameters  $c$  and  $\beta$  in the covariance;  
768 see Figure 5 in Appendix I.

## 769 I Additional numerical results

770 In this section, we present some additional numerics. The plots in Figure 4: A.1-C.1 are identical  
771 to those in Figure 1: A-C, but Figure 1: A.2 - C.2 also show the slopes of the log-log plots as a  
772 function of the sample size. This makes it easier to compare the empirical scaling laws with those  
773 derived in Theorem 1. Figure 5 depicts the relative  $H^1$ -error of the pre-trained transformer under  
774 covariate shifts with respect to a set of parameters in the covariance operator that are different from  
775 the one discussed in Section 4.2. More precisely, we recall that the source term  $f$  is sampled from a  
776 centered Gaussian measure on  $L^2([0, 1])$  with covariance operator given by  $(-\Delta + c\mathbb{I})^{-\beta}$ . During  
777 training, we set the parameters of the covariance as  $\beta = c = 1$ . We then shift the parameters of  
778 the covariance during inference, as defined by the legend of Figure 5: A. Figures 5: B shows the  
779 heat map of the relative  $H^1$ -error with respect to the parameters  $\alpha$  and  $\tau$ . Note that the shift on the  
780 covariance operator of  $f$  defined in Figure 5 differs from the shift defined in Figure 3: C, where the  
781 shift on the covariance operator was defined by constant multiplication. Both cases validate Theorem  
782 5 and provide further evidence that pre-trained transformers are not robust with respect to covariate  
783 shifts. In particular, the prediction errors are more sensitive to the shifts in the amplitude of field and  
784 the smoothness parameter  $\beta$  than the shift in the shift parameter  $c$ . Figure 6 complements Figure 3  
785 with an additional heat map of the relative  $H^1$ -error under tasks shift in the diffusion coefficient  $a$   
786 with respect to the parameters  $\alpha$  and  $\tau$ . Figure 6 shows that the prediction errors under task shifts  
787 remain decently small in a wide range of parameter shifts.



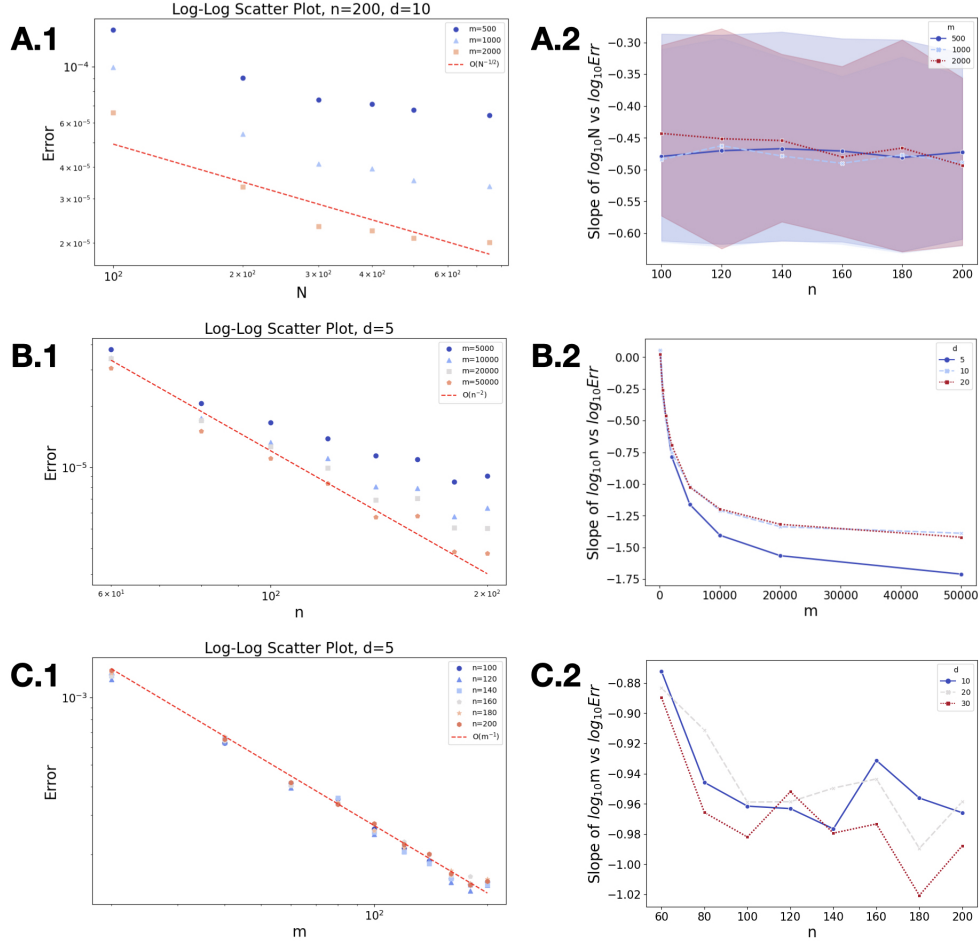


Figure 4: Plots A.1-C.1 are identical to those shown in Figure 1. Plots A.2-C.2 show the slopes of the error curves in the left column as functions of various sample sizes.

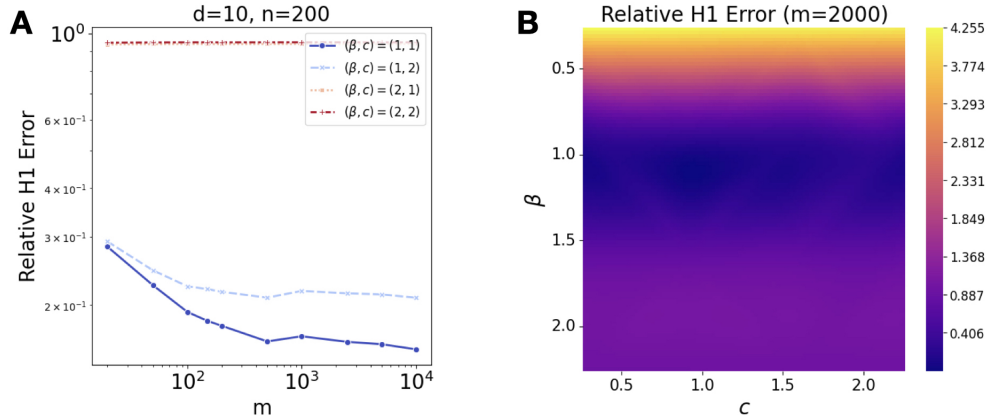


Figure 5: The figures show the relative  $H^1$ -error of learning the linear systems under covariate shifts in the covariance operator  $C = (-\Delta + c\mathbb{I})^{-\beta}$  with respect to the parameters  $c$  and  $\beta$ . During training, we set  $c = \beta = 1$ . Figure A plots the error curves corresponding to four parameter pairs  $(\beta, c)$  as a function of the testing prompt length. Figure B plots the errors for the data corresponding to a wide range of parameter pairs.

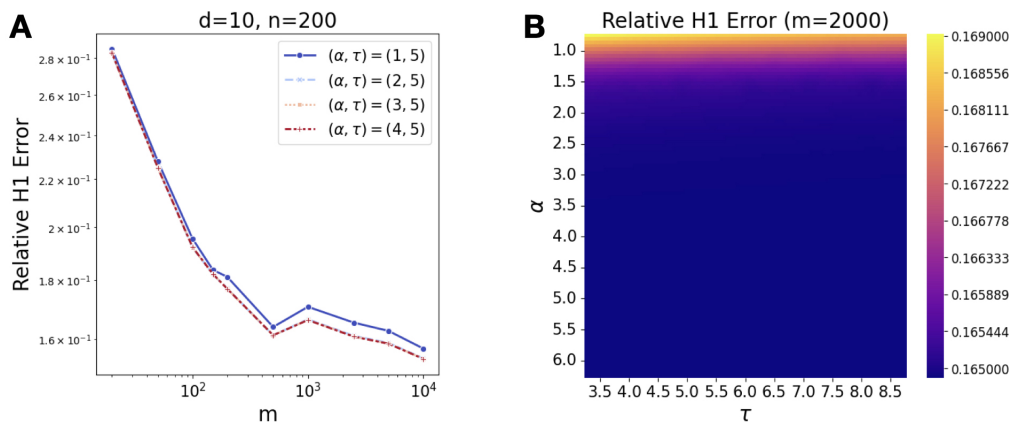


Figure 6: Figure A shows the relative  $H^1$  error as a function of the prompt length under shifts on the distribution of  $a(x)$  (the training distribution is  $a(x) = e^{b(x)}$  with  $b(x) \sim N(0, (-\Delta + \tau \mathbb{I})^{-\alpha})$ ,  $\alpha = 3$  and  $\tau = 5$ ). Figure B shows the corresponding heat map for the relative  $H^1$  error with respect to the parameters  $\alpha$  and  $\tau$ .