

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

Frank Cole*

COLE0932@UMN.EDU

Yulong Lu*

YULONGLU@UMN.EDU

Riley O'Neill*

ONEIL571@UMN.EDU

Tianhao Zhang*

ZHAN7594@UMN.EDU

Abstract

Foundation models for natural language processing, powered by the transformer architecture, exhibit remarkable in-context learning (ICL) capabilities, allowing pre-trained models to adapt to downstream tasks using few-shot prompts without updating their weights. Recently, transformer-based foundation models have also emerged as versatile tools for solving scientific problems, particularly in the realm of partial differential equations (PDEs). However, the theoretical foundations of the ICL capabilities in these scientific models remain largely unexplored. This work develops a rigorous error analysis for transformer-based ICL applied to solution operators associated with a family of linear elliptic PDEs. We first demonstrate that a linear transformer, defined by a linear self-attention layer, can provably learn in-context to invert linear systems arising from the spatial discretization of PDEs. This is achieved by deriving theoretical scaling laws for the prediction risk of the proposed linear transformers in terms of spatial discretization size, the number of training tasks, and the lengths of prompts used during training and inference. These scaling laws also enable us to establish quantitative error bounds for learning PDE solutions. Furthermore, we quantify the adaptability of the pre-trained transformer on downstream PDE tasks that experience distribution shifts in both tasks (represented by PDE coefficients) and input covariates (represented by the source term). To analyze task distribution shifts, we introduce a novel concept of task diversity and characterize the transformer's prediction error in terms of the magnitude of task shift, assuming sufficient diversity in the pre-training tasks. We also establish sufficient conditions to ensure task diversity. Finally, we validate the ICL-capabilities of transformers through extensive numerical experiments.

Keywords: In-context learning, transformer, operator learning, task diversity, linear systems, partial differential equations.

1 Introduction

Foundation models (FMs) for natural language processing (NLP), exemplified by ChatGPT Achiam et al. (2023), have demonstrated unprecedented power in text generation tasks. From an architectural perspective, the main novelty of these models is the use of transformer-based neural networks Vaswani et al. (2017), which are distinguished from

*. 206 Church St. SE, School of Mathematics, University of Minnesota, Minneapolis, MN 55455

feedforward neural networks by their self-attention layers. Those transformer-based FMs pre-trained on a broad range of tasks with large amounts of data, exhibit remarkable adaptability to diverse downstream tasks with limited data Brown et al. (2020). The success of foundation models (FMs) for NLP has recently sparked a large amount of work on building FMs in domain-specific scientific fields Batatia et al. (2023); Celaj et al. (2023); Méndez-Lucio et al. (2022). Specifically, there is growing interest within the community of Scientific Machine Learning (SciML) in building scientific foundation models (SciFMs) to solve complex partial differential equations (PDEs) Subramanian et al. (2024); McCabe et al. (2023); Ye et al. (2024); Yang et al. (2023); Sun et al. (2024).

Traditional deep learning approaches for PDEs such as Physics-Informed Neural Networks Raissi et al. (2019) for learning solutions and neural operators Lu et al. (2019); Li et al. (2020) for learning solution operators need to be retrained from scratch for a different set of coefficients or different PDE systems. Instead, these SciFMs for PDEs, once pre-trained on large datasets of coefficients-solution pairs from multiple PDE systems, can be adapted to solving new PDE systems without training the model from scratch. Even more surprisingly, transformer-based FMs have demonstrated their **in-context learning (ICL)** capability in Achiam et al. (2023); Bubeck et al. (2023); Kirsch et al. (2022) and in SciML Yang et al. (2023); Chen et al. (2024b); Yang and Osher (2024): when given a prompt consisting of examples from a new learning task and a query, they are able to make correct predictions without updating their parameters. While the emergence of ICL has been deemed a paradigm shift in transformer-based FMs, its theoretical understandings remain underdeveloped.

The goal of this paper is to investigate the ICL capability of transformers for solving a class of linear elliptic PDEs. Our study is strongly motivated by the need of characterizing the **scaling and transferability/adaptability** of SciFMs Subramanian et al. (2024); McCabe et al. (2023). In fact, the computing and energy costs for training SciFMs are substantial and increase dramatically with the sizes of the model and data. In addition, the performance of FMs depends on several key parameters, including the size of training data, the model size, the amount of training time, etc. It is thus extremely critical to estimate those key parameters needed to achieve certain prediction accuracy given the allocated compute budget prior to training. This requests a **neural scaling law** that can quantify the prediction risk of an FM as a function of those parameters. While empirical scaling laws have been identified for Large Language Models (LLMs) Kaplan et al. (2020); Hoffmann et al. (2022) and more recently attempted in a SciFM for PDEs Subramanian et al. (2024); McCabe et al. (2023), a rigorous characterization of scaling laws remains open. Additionally, the generalization performance of SciFMs encounters significant hurdles due to prevalent **distribution shifts** between tasks and data used in pre-training and those in adaptation Subramanian et al. (2024); McCabe et al. (2023); Ye et al. (2024); Yang et al. (2023). For instance, the behavior and fine structures (e.g. multiple spatial-temporal scales) of a PDE solution can change dramatically in response to the changes in the physical parameters and conditions. Therefore a rigorous quantification of the generalization performance of SciFMs on downstream tasks due to the various domain shifts is an important step for understanding their capabilities and limitations.

1.1 Main contributions

We highlight our main contributions as follows:

- We formalize a framework for learning the solution operators of linear elliptic PDEs in-context. This is based on (1) reducing the infinite dimensional PDE problem into a problem of solving a finite dimensional linear system arising from spatial discretization of the PDE and (2) learning to invert the finite dimensional linear system in-context.
- We adopt transformers defined by single linear self-attention layers for ICL of the linear systems and establish a quantitative generalization error bound of ICL in terms of the discretization size, the number of pre-training tasks, and the lengths of prompts used in pre-training and downstream tasks; see Theorem 1. This bound further enables us to prove an H^1 -error bound for learning the solution of PDEs; see Theorem 2.
- We establish general prediction error bounds for the pre-trained transformer under distribution shifts with respect to tasks (represented by the coefficients of the PDE) and the data covariates (represented by the source term), in Theorem 3 and in Theorem 7 respectively. In the setting of task shifts, we introduce a novel concept of **task diversity** and show that pre-trained transformers can provably generalize even when the downstream task undergoes distribution shifts provided that the pre-training task distribution is sufficiently diverse; see Theorem 4.
- Additionally, we provide several sufficient conditions under which task diversity condition holds (see Theorem 5), and construct simple examples where the task diversity fails to hold (see Theorem 6).
- We demonstrate the ICL ability of linear transformers for learning both the PDE solutions and the associated linear systems through extensive numerical experiments.

1.2 Related work

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

Theory of ICL for linear regression and other statistical models. The work Garg et al. (2022) provides theoretical understanding of the ability of transformers in learning simple functions in-context. In the follow-up works Akyürek et al. (2022); Von Oswald et al. (2023), it is shown by explicit construction of attention matrices that linear transformers can implement a single step of gradient descent when given a new in-context linear regression task, and numerical experiments supported that trained transformer indeed implement

gradient descent on unseen tasks. Several recent works Mahankali et al. (2024); Zhang et al. (2023a); Ahn et al. (2024) extend the results of Von Oswald et al. (2023) by proving that one step of gradient descent is indeed optimal for learning linear models in-context. These works are further complemented by ICL guarantees for learning nonlinear functions Bai et al. (2024); Cheng et al. (2023); Kim et al. (2024) and for reinforcement learning problems Lin et al. (2023). Besides the explicit constructions of transformers to implement desired algorithms in various learning tasks, recent works Zhang et al. (2023a); Chen et al. (2024a) study the optimization landscape of single-layer linear transformers in learning of linear functions and characterize the convergence guarantees of gradient descent for training linear transformers.

Optimization analysis of transformers have also been studied in various settings, including for learning nonlinear functions Cheng et al. (2023); Kim et al. (2024), using nonlinear attentions Huang et al. (2023); Nichani et al. (2024) (e.g. softmax and ReLU), and multiple heads Chen et al. (2024a). Regarding the generalization error with respect to the number of tasks, the paper Wu et al. (2023) studies the behavior of in-context linear regression when the transformer parameters are trained by stochastic gradient descent. The work Mroueh (2023) proposes a general statistical learning framework for analyzing the generalization error of transformers for ICL.

Going beyond the i.i.d setting, several works also investigate the ICL capability of transformers when the data are defined by Markov chains Edelman et al. (2024) and dynamical systems Goel and Bartlett (2023); Li et al. (2023); Du et al. (2023). We also would like to mention several works that explain how transformers perform ICL from the Bayesian perspective Xie et al. (2021); Zhang et al. (2023b); An et al. (2020).

Among the aforementioned works, the settings of Zhang et al. (2023a); Ahn et al. (2024); Chen et al. (2024a) are closest to us. Our theoretical bound on the population risk extends the results of Zhang et al. (2023a); Ahn et al. (2024) for the linear regression tasks to the tasks of inverting linear systems that are associated to elliptic PDEs. Different from those works where the data and task distributions are assumed to be Gaussian, the task distributions considered here are non-Gaussian and are fully determined by the PDE structure. Our out-of-domain generalization error bounds substantially improve the earlier generalization bound Mroueh (2023) established for general ICL problems. In particular, we show that the error due to the distribution shift can be reduced by a factor of $1/m$, where m is the prompt-length of a downstream task. This rigorously justifies that the a key robust feature of pre-trained transformer model with respect to task distribution shifts, which has previously been empirically observed in ICL of PDEs (see e.g. Yang et al. (2023); Yang and Osher (2024)), but has only been rigorously studied by Zhang et al. (2023a) in the linear regression setting.

2 Problem set-up

2.1 In-context operator learning of linear elliptic PDEs

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

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

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 standard well-posedness of the elliptic PDE, the solution $u \in \mathcal{X}_u \subset H_0^1(\Omega)$. We are interested in learning the linear solution operator $\Psi : f \rightarrow u \in \mathcal{X}_u$ in context Yang et al. (2023). More specifically, at the training stage we are given a training dataset comprising N length- n prompts of source-solution 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 , 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 P_V are distributions on the coefficient a and V respectively. An ICL model, after pre-trained on the data above, is asked to predict the solution u for a new source term f conditioned on a new prompt $(f_i, u_i)_{i=1}^m$ which may or may not have the same distribution as the training prompts. Further, the prompt-length m in the downstream task may be different from the prompt-length n in the training.

While the ideal ICL problem above is stated for learning operators defined on infinite dimensional function spaces, a practical ICL model (e.g. a transformer) can only operate on finite dimensional data, which are typically observed in the form of finite dimensional projections or discrete evaluations. 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 set $\Phi(x) := [\phi_1, \dots, \phi_d(x)]$ for some $d < \infty$. An approximate solution \tilde{u} to problem (1) can be constructed in the framework of Galerkin method: we seek $\tilde{u}(x) = \langle \mathbf{u}, \Phi(x) \rangle$ where $\mathbf{u} \in \mathbb{R}^d$ 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 $\mathbf{f} = (f_i) \in \mathbb{R}^d$ are defined by

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

Note that the Galerkin method can be viewed as a special instance of the operator-network defined by principle component analysis or PCA-Net Bhattacharya et al. (2021). In fact, if we define the encoder $\mathcal{E}_f : \mathcal{X}_f \rightarrow \mathbb{R}^d$ on the space \mathcal{X}_f by $\mathcal{E}_f f := \mathbf{f}$, the decoder $\mathcal{D}_u : \mathbb{R}^d \rightarrow \mathcal{X}_u$ by $(\mathcal{D}_u \mathbf{u})(x) := \langle \mathbf{u}, \Phi(x) \rangle$ and $\mathcal{G} := \mathbb{R}^d \rightarrow \mathbb{R}^d$ by $\mathcal{G}\mathbf{f} = A^{-1}\mathbf{f}$, then $\tilde{\Psi} := \mathcal{D}_u \circ \mathcal{G} \circ \mathcal{E}_f$ defines an approximation to the solution operator Ψ . As quantitative discretization error bounds of PDEs are well established, e.g. for finite element methods Brenner and Scott (2007) and spectral methods Shen et al. (2011), this paper focuses on the error analysis of in-context learning of the finite dimensional linear systems defined by the matrix inversion A^{-1} , which will ultimately translate to estimation bounds for the solutions of PDEs.

2.2 ICL of linear systems

The consideration above reduces the original infinite dimensional in-context operator learning problem to the finite dimensional ICL problem of solving linear systems. To keep the framework more general, 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 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 new query input $\mathbf{y}_{n+1} \in \mathbb{R}^d$. Given multiple prompts, the model aims to predict \mathbf{x}_{n+1} corresponding to the new independent query input \mathbf{y}_{n+1} . Unlike in supervised learning, each prompt the model takes is drawn from a different data distribution. To be more precise, for $j = 1, \dots, N$, we assume 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 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 $A^{(j)} \stackrel{i.i.d.}{\sim} p_A$. Informed by task matrices derived from discretizations of PDEs as illuminated in (2), we make the following assumption on the task distribution p_A .

Assumption 1. *The task distribution p_A is supported on the set of symmetric positive definite 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 $A \in \text{supp}(p_A)$. The source term \mathbf{y} follows a Gaussian distribution $N(0, \Sigma)$.*

Observe that Assumption 1 on A is very mild and holds for instance whenever the coefficient a is strictly positive and V is non-negative and bounded. We will make repeated use of the bounds¹

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

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

2.3 Linear transformer architecture for linear systems

Inspired by the recent line of work on ICL of linear functions, we consider a linear transformer defined by a single-layer linear self-attention layer for our ICL model. Following the standard convention, we 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^{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. (2023a); Ahn et al. (2024),

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

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 $\theta = (P, Q)$,

$$\text{TF}_\theta(Z) = PA^{-1}Y_nQ\mathbf{y},$$

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 examples. We work with this simplified parameterization for the remainder of our theoretical analysis.

2.4 Generalization of ICL

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

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

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 have access to the distribution on tasks, P and Q are instead trained by minimizing the corresponding *empirical risk* functional defined on N tasks:

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

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 the in-context examples $\{\mathbf{y}_1^{(i)}, \dots, \mathbf{y}_n^{(i)}\}$ which are (jointly) independent from \mathbf{y}_i . Our empirical risk is closely related to the few-shot risk for in-context learning introduced in Mroueh (2023).

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

- **In-domain generalization:** The distributions of tasks and of prompt data in the pre-training are the same as these in the downstream task ($P_\mathbf{y} = P'_\mathbf{y}$ and $P_A = P'_A$). Thus in-domain generalization measures the testing performance on unseen samples in the downstream task that do not appear in the training samples. The in-domain generalization error is defined by

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

- **Out-of-domain (OOD) generalization:** The distributions of tasks or within-task data in the pre-training are different from those in the downstream task, i.e. $P_\mathbf{y} \neq P'_\mathbf{y}$ or $P_A \neq P'_A$. Specifically, the OOD-generalization error with respect to the task distribution shift is defined by

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

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

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

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

3 Theoretical results

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

Our first result studies the generalization ability of the transformer obtained by empirical risk minimization over a set of norm-constrained transformers, where the error is measured by the prediction risk \mathcal{R}_m . To state the result, it will be convenient to define the weighted trace of a matrix K with respect to the covariance $\Sigma = W \Lambda W^T$:

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

where $\{(\sigma_i^2, \varphi_i)\}_{i=1}^d = \{(\sigma_i^2, W e_i)\}_{i=1}^d$ are the eigenpairs of the covariance matrix Σ . Note that the Σ -weighted trace is independent of the choice of eigenbasis of Σ .

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

$$\begin{aligned} \mathcal{R}_m(\hat{\theta}) &\lesssim \frac{(c_A^2 + d) \text{Tr}(\Sigma)}{m} + \frac{c_A^2 C_A^4 \|\Sigma\|_{op}^2 \|\Sigma^{-1}\|_{op}^2 \left(1 + \text{Tr}_{\Sigma}(\mathbb{E}_{A \sim p_A}[A^{-2}])\right)^2 \text{Tr}(\Sigma)}{n^2} \\ &\quad + \frac{d^2 c_A^2 \|\Sigma\|_{op}^2 \max(1, \|\Sigma^{-1}\|_{op})^4}{\sqrt{N}} \\ &\quad + \max(1, \|\Sigma^{-1}\|_{op})^4 c_A^2 \max(\text{Tr}(\Sigma), \|\Sigma\|_{op}^2) \text{Tr}(\Sigma) \left| \frac{1}{n} - \frac{1}{m} \right|, \end{aligned} \quad (10)$$

where we have omitted factors which are polylog in N .

The first three terms on the right side of (10) are the generalization errors of the ICL model due to finite amount of training data and training tasks. The last term on the right side of (10), which we term the ‘‘context mismatch error’’, is likely due to an artifact of our proof strategy and can potentially be removed with a refined analysis. This term is not observed in our numerical experiments; see Figure 1. However, in the practical ² regime where the length of the testing prompts is less than that of the training prompts

2. The performance of GPTs is known to deteriorate when the test sequence length exceeds the train sequence length; Zhang et al. (2023a) conjectures this phenomenon to be the result of positional encoding.

(i.e. $m \leq n$), we have $|\frac{1}{n} - \frac{1}{m}| \leq \frac{1}{m}$, and hence the context-mismatch error is absorbed into the $O(\frac{1}{m})$ term, leading to the following overall generalization bound

$$\mathcal{R}_m(\hat{\theta}) \lesssim \frac{1}{m} + \frac{1}{n^2} + \frac{1}{\sqrt{N}}. \quad (11)$$

Notice also that the prompt lengths during training and testing contribute different rates to the overall sample complexity bound, with the sequence length n during training contributing an $O(n^{-2})$ rate while the sequence length m at inference contributing an $O(m^{-1})$ rate; a similar phenomenon was observed in (Zhang et al., 2023a, Theorem 4.2) for in-context linear regression.

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

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

Theorem 2. *Let Φ' be the stiffness matrix defined by $\Phi'_{ij} = (\phi'_i, \phi'_j)_{L^2(\Omega)}$ and let Φ 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,*

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

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

Theorem 2 bounds the in-domain generalization error of ICL for the PDE as a sum of the discretization error of the PDE solver and the statistical error of learning the linear system associated to the discretization of the PDE. It is worth-noting that there is a trade-off between the two terms; the first term decreases as the number of basis functions (or fineness of the mesh) increases, while the prefactor $\lambda_{\max}(\Phi^{-1/2}\Phi'\Phi^{-1/2})$ in the second term can grow as the number of basis functions tends to infinity. This suggests that the PDE recovery error in the H^1 -metric will be large if the dimension d is either too large or too small. We demonstrate this numerically in Figure 2 by plotting the H^1 -error between the learned solution and ground truth against the dimension d .

The abstract bound established in Theorem 2 is agnostic to the choice of PDE discretization. Below, we present how this result can be used to derive an explicit error estimate for the ICL in the context of a P^1 -finite element discretization of the PDE in one dimension.

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 Ω , 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

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

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

$$\mathbb{E}\left[\|u - \hat{u}_h\|_{H^1(\Omega)}^2\right] \lesssim h^2 + \frac{1}{h^{2m}} + \frac{1}{h^{10+4\alpha}n^2} + \frac{1}{h^{4+8\alpha}\sqrt{N}}, \quad (12)$$

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

$$\mathbb{E}\left[\|u - \hat{u}_h\|_{H^1(\Omega)}^2\right] \lesssim \frac{1}{d^2} + \frac{d^2}{m} + \frac{d^{10+4\alpha}}{n^2} + \frac{d^{4+8\alpha}}{\sqrt{N}}. \quad (13)$$

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

3.3 OOD-generalization under task distribution shift

Let $\hat{\theta}$ denote the minimizer of the empirical risk $\mathcal{R}_{n,N}$ over the bounded set $\{\|\theta\| \leq M\}$ for some $M > 0$, and recall that the training tasks (modeled by A) are drawn from a distribution p_A . Let p'_A denote the distribution of A in the downstream tasks, and let

$\mathcal{R}_m, \mathcal{R}'_m$ be the prediction risk functionals defined as in (8) where the expectations over tasks are taken with respect to p_A and p'_A respectively. We would like to bound the quantity $\mathcal{R}'_m(\hat{\theta})$, which represents the test error of the trained transformer under a shift on the task distribution. We say that a pre-trained model $\hat{\theta}$ **achieves OOD-generalization** if its population risk with respect to the downstream task distribution p'_A converges to zero in probability: $\lim_{(m,n,N) \rightarrow \infty} \mathcal{R}'_m(\hat{\theta}) \xrightarrow{P} 0$. In order to state our results on OOD-generalization, we first introduce the following “infinite-context” variant of the in-domain denoted by \mathcal{R}_∞ :

$$\mathcal{R}_\infty(\theta) = \mathbb{E}_{A \sim p_A} \|(PA^{-1}\Sigma Q - A^{-1})\Sigma^{1/2}\|_F^2. \quad (14)$$

We also define and OOD-generalization risk \mathcal{R}'_∞ similar to above with p_A replaced by p'_A . We denote by \mathcal{M}_∞ and \mathcal{M}'_∞ the sets of minimizers of \mathcal{R}_∞ and \mathcal{R}'_∞ respectively. We now present a rather general theorem that bounds the OOD-generalization error.

Theorem 3. *Let p_A and p'_A denote the pre-training and downstream task distributions respectively and assume both satisfy Assumption 1. Let $\mathcal{M}_\infty(p_A)$ and $\mathcal{M}_\infty(p'_A)$ denote the minimizers of \mathcal{R}_∞ and \mathcal{R}'_∞ respectively, and let $\hat{\theta} \in \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_{n,N}(\theta)$ denote the empirical risk minimizer. Then*

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

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

The precise definition of the discrepancy $d(p_A, p'_A)$ is technical and can be found in the statement of Lemma 2 in the appendix. Theorem 3 bounds OOD generalization error by a sum of three terms: the in-domain generalization error, the task-shift error, and the model error, the latter of which is captured by $\operatorname{dist}(\hat{\theta}_n, \mathcal{M}_\infty)$ and $\operatorname{dist}(\hat{\theta}, \mathcal{M}'_\infty)$. A salient feature of Theorem 4, compared to the prior ICL-generalization bound Mroueh (2023) under distribution shift, is that the task-shift error inherits a factor of m^{-1} , which elucidates the robustness of transformers under shifts in the task distribution. Theorem 4 also extends the prior OOD-generalization result of ICL for linear regression Zhang et al. (2023a) to learning linear systems. However, unlike in the linear regression setting, the set of minimizers of the population risk in the linear system setting can vary substantially when the task distribution changes. Because of this, the terms $\operatorname{dist}(\hat{\theta}, \mathcal{M}_\infty)$ and $\operatorname{dist}(\hat{\theta}, \mathcal{M}_\infty)'$ warrant a more careful examination. Since the empirical risk $\mathcal{R}_{n,N}$ recovers the infinite-context population risk \mathcal{R}_∞ as n and N tend to ∞ and $\hat{\theta}$ is a minimizer, we expect $\operatorname{dist}(\hat{\theta}, \mathcal{M}_\infty)$ to tend to zero as n and N tend to ∞ ; this is made precise in Lemma 3 in the appendix. Without further assumptions, we cannot expect the term $\operatorname{dist}(\hat{\theta}, \mathcal{M}'_\infty)$ to decay as the sample size increases, because $\hat{\theta}$ is not trained on data from p'_A . However, we note that if $\mathcal{M}_\infty \subseteq \mathcal{M}'_\infty$, then $\operatorname{dist}(\hat{\theta}, \mathcal{M}'_\infty) \leq \operatorname{dist}(\hat{\theta}, \mathcal{M}_\infty)$, and hence we can expect all terms in Theorem 3 to decay as the sample size increases. This motivates the following key notion of task diversity.

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

The importance of task diversity has been observed in the prior work Tripuraneni et al. (2020) for transfer learning. Our notion of diversity differs from the previous notion in that we compare the set of minimizers of population losses instead of the loss values. Theorem 4, which is a direct corollary of Theorem 3, shows that the task diversity, in the sense of Definition 1, is sufficient for the pre-trained transformer to achieve OOD-generalization.

Theorem 4. *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 $\hat{\theta} \in \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_{n,N}(\theta)$, we have*

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

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

Remark 1. *Theorem 4 can be combined with Theorem 2 to obtain bounds on the OOD-generalization error for learning the corresponding PDE solution. Specifically, if $\hat{\theta}$ is the transformer parameter obtained by empirical risk minimization on the in-domain risk, \hat{u} the corresponding PDE solution, and p'_A the downstream task distribution, then the proof of Theorem 2 immediately implies that*

$$\mathbb{E}_{u \sim p'_A} [\|u - \hat{u}\|_{H^1(\Omega)}^2] \lesssim \mathbb{E}_{u \sim p'_A} [\|u - u_d\|_{H^1(\Omega)}^2] + (1 + \lambda_{\max}(\Phi^{-1/2} \Phi' \Phi^{-1/2})) \mathcal{R}'_m(\hat{\theta}).$$

Notice that, under the task diversity assumption, Theorem 4 can be used to bound $\mathcal{R}'_m(\hat{\theta})$.

In Figure 3, we numerically demonstrate the robustness of pre-trained transformers under shifts on the PDE coefficients. The numerical results show that even under more extreme shifts, the OOD-generalization error is mitigated by increasing the prompt length, as predicted by Theorem 4. This suggests that the pre-training task distributions corresponding to elliptic PDE problems are sufficiently diverse in the sense of Definition 1, although proving this rigorously remains an open question.

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

Theorem 5. *Let p_A, p'_A be two distributions on the matrices A that satisfy Assumption 1. Then*

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

The first statement of Theorem 5 is natural: it says that the pre-training task distribution is diverse whenever the downstream task distribution is a “subset” of it, in the sense of supports. The second condition is particularly interesting because it implies OOD-generalization (by Theorem 4) regardless of the downstream task distribution. The second condition based on the centralizer of the set $\mathcal{S}(p_A)$ is less obvious, but heuristically it

enforces that the support of p_A must be large enough that the only matrices which can commute with all pairwise products in $\mathcal{S}(p_A)$ are scalars. Our empirical results suggest that the task distributions associated to elliptic PDE problems are diverse.

An inspection of the proof of Theorem 5 reveals that if $\text{supp}(p_A)$ satisfies the condition that the centralizer of $\{A_1 A_2^{-1} : A_1, A_2 \in \text{supp}(p_j)\}$ is trivial, then all minimizers of \mathcal{R}_∞ are of the form $\{(P, Q) = (c\mathbf{I}_d, c^{-1}\Sigma^{-1}) : c \neq 0\}$. In this case, it is worth noting that the discrepancy on task distributions $d(p_A, p'_A)$ defined in Theorem 4 admits a much simpler expression. We state this result as a corollary 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

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

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

To conclude this section, we investigate the diversity of task distributions whose support consists of simultaneously diagonalizable matrices. The simultaneous-diagonalizability of task matrices has been a key assumption in the existing theoretical analysis of in-context learning of linear systems (Chen et al. (2024a)) and in the in-context learning of linear dynamical systems (Sander et al. (2024)). It is therefore natural to ask whether a task distribution whose support consists of simultaneously diagonalizable matrices is diverse in the sense of Definition 1. This question has an important interpretation in the in-context learning of elliptic PDEs: if the diffusion coefficient $a(x)$ and potential $V(x)$ are both constant, $a(x) \equiv a_0$, $V(x) \equiv v_0$, then the solution operator of the corresponding elliptic PDE is given by $\left(-a_0\Delta + v_0\mathbf{I}\right)^{-1}$, whose diagonalization is independent of the constants a_0 and v_0 . Thus, in order to understand whether a transformer that is pre-trained to solve elliptic PDEs with constant coefficients can make accurate predictions on equations with non-constant coefficients, it is essential to understand the diversity of task distributions whose support consists of simultaneously diagonalizable matrices.

Theorem 6. *Let p_A and p'_A denote the pre-training and downstream task distributions, and suppose that the matrices in $\text{supp}(p_A)$ are simultaneously diagonalizable for a common orthogonal matrix U . Suppose additionally that there exist matrices $A_1, A_2 \in \text{supp}(p_A)$ and $A'_1 A'_2 \in \text{supp}(p'_A)$ such that $A_1 A_2^{-1}$ and $A'_1 (A'_2)^{-1}$ have no repeated eigenvalues.*

1. *If $\text{supp}(p'_A)$ is also simultaneously diagonalizable with respect to U , then p_A is diverse relative to p'_A .*
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 respect to U , then p_A is not diverse relative to p'_A .*

Theorem 6 reveals that a simultaneously-diagonalizable task distribution cannot achieve out-of-distribution generalization under arbitrary shifts in the downstream task distribution; in general, the downstream task distribution must also be simultaneously diagonalizable in the same basis. However, it also shows that, provided the pre-training and downstream task distributions are simultaneously diagonalizable, pre-trained transformers can generalize under arbitrary shifts on the distribution shifts on the eigenvalues of the task matrices. This provides a precise characterization of the diversity of a simultaneously diagonalizable task distribution. In Figure 4, we demonstrate the importance of task diversity by computing the OOD generalization error of a transformer pre-trained to solve Equation (1) with constant coefficients when the task distribution at inference corresponds to a PDE with non-constant coefficients. We find that the transformer is much more sensitive to task shifts in this case, but the OOD generalization error is improved if the transformer parameters are initialized around robust population risk minimizers.

3.4 OOD-generalization under covariate distribution shift

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

Theorem 7. *Let $\Sigma = W\Lambda W^T$ and $\tilde{\Sigma} = \tilde{W}\tilde{\Lambda}\tilde{W}^T$ be the covariance matrices of Gaussian covariates used in the training and testing respectively. Let (\hat{P}, \hat{Q}) be minimizers of the empirical risk associated to covariates sampled from $N(0, \Sigma)$ and take $M > 0$ such that $\max(\|\hat{P}\|_F, \|\hat{Q}\|_F) \leq M$. Then*

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

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

Theorem 7 states that the OOD-generalization error with respect to the covariate distribution shift is Lipschitz stable with respect to changes in the covariance matrix. However, unlike the case of task distribution shift, the covariate distribution shift error cannot be mitigated by increasing the prompt-length in the downstream task; see also Figure 3. A similar phenomenon was observed in Zhang et al. (2023a).

4 Numerical experiments

4.1 In-domain generalization

We investigate numerically the neural scaling law of the transformer model for solving the linear system associated to the Galerkin discretization of the elliptic PDE (1) in the setting of in-domain generalization. In more detail, 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 f is white noise, i.e. $f \sim N(0, \mathbb{I})$, where \mathbb{I} denotes the identity operator. Further, 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 takes values uniformly in $[1, 2]$. We discretize the PDE using Galerkin projection under d sine bases: $\phi_i(x) = \sin(\pi i x)$, $i = 1, \dots, d$. This leads to the linear system $A\mathbf{u} = \mathbf{f}$, where $\mathbf{f} \sim N(0, \mathbf{I}_d)$ and

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

The prompts used for pre-training are then built on observations of the form

$$((\mathbf{f}_1, A^{-1}\mathbf{f}_1), \dots, (\mathbf{f}_n, A^{-1}\mathbf{f}_n)).$$

In Figure 1: A-C, we demonstrate the empirical scaling law of the linear transformer for learning the discrete linear system by showing the log-log plots of the ℓ^2 -errors as functions of the number of pre-training tasks N , the sequence length n during training and the sequence length m at inference. These numerical results suggest that the decaying rates of the prediction errors are $O(N^{-\frac{1}{2}})$, $O(n^{-2})$ and $O(m^{-1})$ respectively, which almost agrees with the rates predicted in Theorem 1, except in the regime where $m > n$; in this case, the numerics suggest that the scaling of the population risk with respect to n is $O(n^{-2})$, whereas Theorem 1 predicts that the error is $O(n^{-1})$ in this regime. This gives further evidence towards the potential sub-optimality of our proof strategy. Figure 1:D shows that prediction error increases as d increases indicating that ICL of the linear system becomes harder as the d increases.

We also demonstrate the ICL-generalization error for learning the PDE solutions. Figure 2:B shows the H^1 -error curve between the numerical solution predicted by the ICL-model and the ground-truth as a function of the number of bases d , while fixing the prompt-lengths and the number of tasks. The U-shaped curve indicates the trade-off between the dimension of the discrete problem and the amount of data.

4.2 Out-of-domain generalization

We validate the ICL-capability of pre-trained transformers for learning the linear systems and PDEs under shifts in the distribution of both the PDE coefficients and the source term.

Task shifts on the PDE coefficients. For task shifts, we first test the behavior of pre-trained transformers when the smoothness of the PDE coefficients differs between training and inference, see Figure 3, Plots (A) and (B). Specifically, for the PDE (1) in one dimension, we consider the task distribution shifts in a and V exclusively. Throughout this experiment, the distribution of the source term f is fixed as the centered Gaussian measure with covariance operator $(-\Delta + \mathbb{I})^{-1}$.

We consider a as a random field sampled from a **log-normal** distribution, denoted by $p_a(\alpha, \tau)$, i.e. $a(x) = e^{b(x)}$ where $b(x)$ is a random field sampled from the centered Gaussian measure with covariance operator $(-\Delta + \tau \mathbb{I})^{-\alpha}$. For the potential V , we assume that V is uniform random field sampled from a distribution, denoted by $p_V(a, b)$, that is obtained by dividing the domain into $2d+1$ sub-intervals and in each cell independently, the potential takes values uniformly in $[a, b]$. In this experiment we fix the number of basis $d = 50$. The transformer model for solving the 50-dimensional linear system is trained over

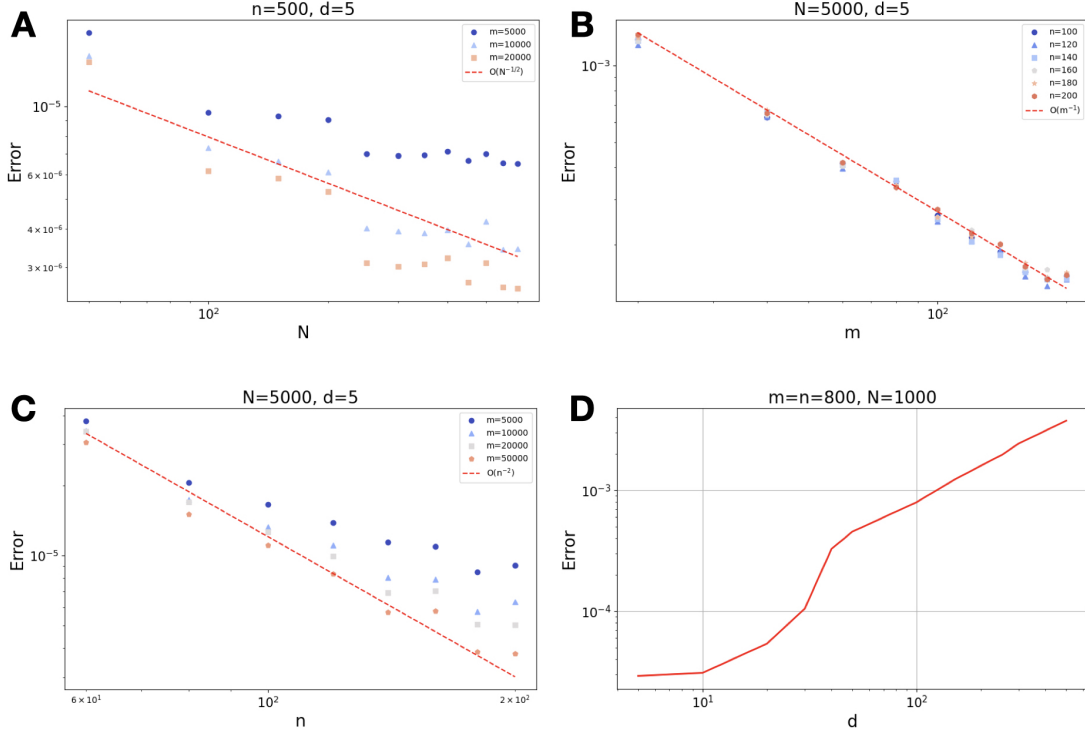


Figure 1: The figures A-D show the log-log plots for the ℓ^2 -error of learning the linear system associated to the PDE discretization with respect to the number of tasks N , the prompt length n during training, the prompt length m during inference, and the dimension d of the linear system.

$N = 5000$ tasks with $n = 300$ pairs of (\mathbf{f}, \mathbf{u}) . In Figure 3: A-B, we show the generalization errors of the pre-trained transformer under distribution shifts with respect to a with a fixed $V \sim p_V(1, 2)$ and with respect to V with a fixed $a \sim p_a(3, 5)$, respectively. As for error metrics, we consider the ℓ^2 -error for learning the coefficient vectors, as well as the relative L^2 and relative H^1 metrics for the PDE solutions. Figure 3: A shows that the pre-trained transformer can perform equally well on tasks on a variety of smoothness parameters α of the log-normal field a but perform slightly worse on tasks with less regular a . We refer to Figure 6 in the appendix for additional plots on the distribution shifts errors with a wide range of τ and α . Figure 3: B shows the OOD-generalization errors under distribution shifts on the range of the uniform field V . As shown in the figure, the errors increase as the distribution shift of the range becomes stronger, but they decrease as the context length at inference increases, as predicted by Theorem 4.

Covariate shifts. Next, we test the performance of the pre-trained transformer under covariate distribution shifts. Specifically, we train the model to solve the PDE (1), where the distribution on the coefficients is given by $a(x) \sim p_a(3, 5)$ and $V(x) \sim p_V(1, 2)$. We assume that the source term $f \sim N(0, C_{c,\beta})$ for $C_{c,\beta} = (-\Delta + c\mathbb{I})^{-\beta}$. During pre-training, we set $c = \beta = 1$. Then, at inference, we consider solving the PDE (1), with the same

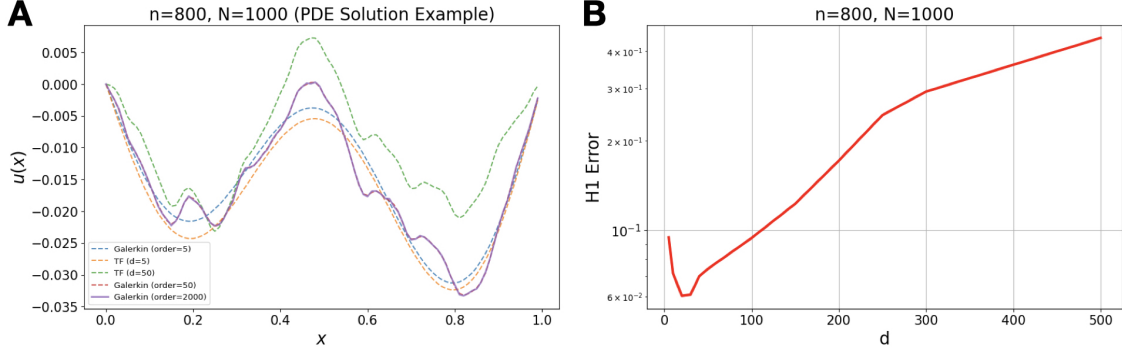


Figure 2: The left plot shows the PDE solution defined by the pre-trained transformer with the reference solution, obtained by Galerkin’s method with 2000 basis functions. The right plot shows the H^1 -error between the solution predicted by the transformer and reference solution with respect to the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

distribution on the coefficients $a(x), V(x)$, but under two types of covariate shifts on the source term f :

1. At inference, instead of sampling $f \sim N(0, C_{1,1})$ as in pre-training, we sample $f \sim N(0, 3 \cdot C_{1,1})$ and $f \sim N(0, 5 \cdot C_{1,1})$. In other words, we consider distribution shifts defined by scalar multiplication on the covariance matrix; the errors are shown in Figure 3:C.
2. At inference, we sample $f \sim N(0, C_{c,\beta})$ with $(c, \beta) \neq (1, 1)$. In other words, we shift the parameters of the covariance of the covariance matrix which control the smoothness of f ; the errors are shown in Figure 3:D.

Figure 3: C-D show that the pre-trained transformers are generally not robust to covariate distribution shifts. In addition, increasing the prompt-length m does not reduce the prediction errors. We refer to Figure 6 in the appendix for additional plots on the covariate shift errors for a wider range of parameters that specify the covariance operator for f .

Testing the task diversity assumption. As another task distribution shift, we consider pre-training a transformer to solve the equation $(-\Delta + V(x))u(x) = f$ on $[0, 1]$ with zero boundary conditions, where the potential function $V(x)$ is almost surely constant and $f \sim N(0, \mathbb{I})$ is white noise³. Specifically, during pre-training, we set $V(x) \equiv v_0$, where v_0 is sampled from the uniform distribution on $[1, 2]$. Then, at inference, we consider solving the same PDE with $f \sim N(0, \mathbb{I})$ but where the potential $V(x) \sim p_V(a, b)$ for various choices of a and b . In particular, the potential function is non-constant for the downstream tasks.

In this case, under a Galerkin discretization of the forward operator in the sine basis, Theorem 6 implies that the pre-training task distribution is not diverse in the sense of

3. We use white noise for f because we expect the attention matrices that minimize the empirical risk to be perturbations of the set $\{(P, \Sigma^{-1}P^{-1})\}$, and we want to ensure that the covariance matrix Σ of the Galerkin projection of f is well-conditioned.

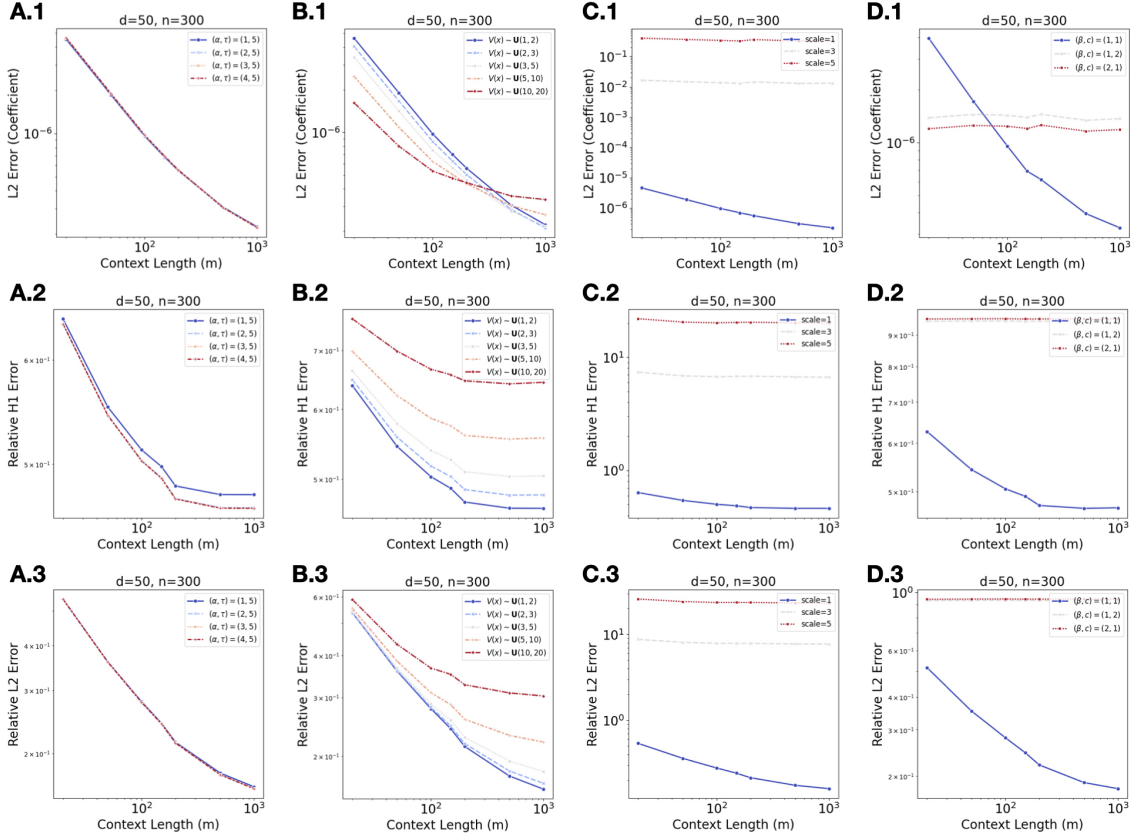


Figure 3: Columns A and B show the prediction error of the transformer under distribution shifts on the coefficients a and V respectively. Columns C and D show the error under the covariate shift in the source term f . The model is trained with $d = 50$, $n = 300$ and $N = 5000$.

Definition 1, since the pre-training task distribution consists of diagonal matrices. More specifically, any parameter of the form $(P, Q) = (P, P^{-1})$, where P is an invertible diagonal matrix, is a minimizer of the infinite-context risk for the pre-training task distribution. However, the test error of such minimizers is larger when P does not commute with the support of the task distribution, i.e., when P is not a multiple of the identity matrix. We therefore expect the OOD-generalization error of the pre-trained transformer to be quite sensitive to the initialization of the attention parameters.

In Figure 4, we plot the ℓ^2 -error for learning the linear system, the relative L^2 error for learning the PDE, and the relative H^1 -error for learning the PDE, of two pre-trained transformers with different initializations. We see that the pre-trained transformer achieves a lower OOD generalization error when the matrix P is initialized around the identity matrix, which agrees with Part (2) of Theorem 5 which shows that $P = \mathbf{I}_d$ is a sufficient condition on the minimizers of the infinite-context risk to achieve OOD generalization. In particular, Figure 4: A.3-B.3 demonstrate that when the attention matrix P is initialized near the identity matrix, the test error (in the relative L^2 metric) in solving the PDE

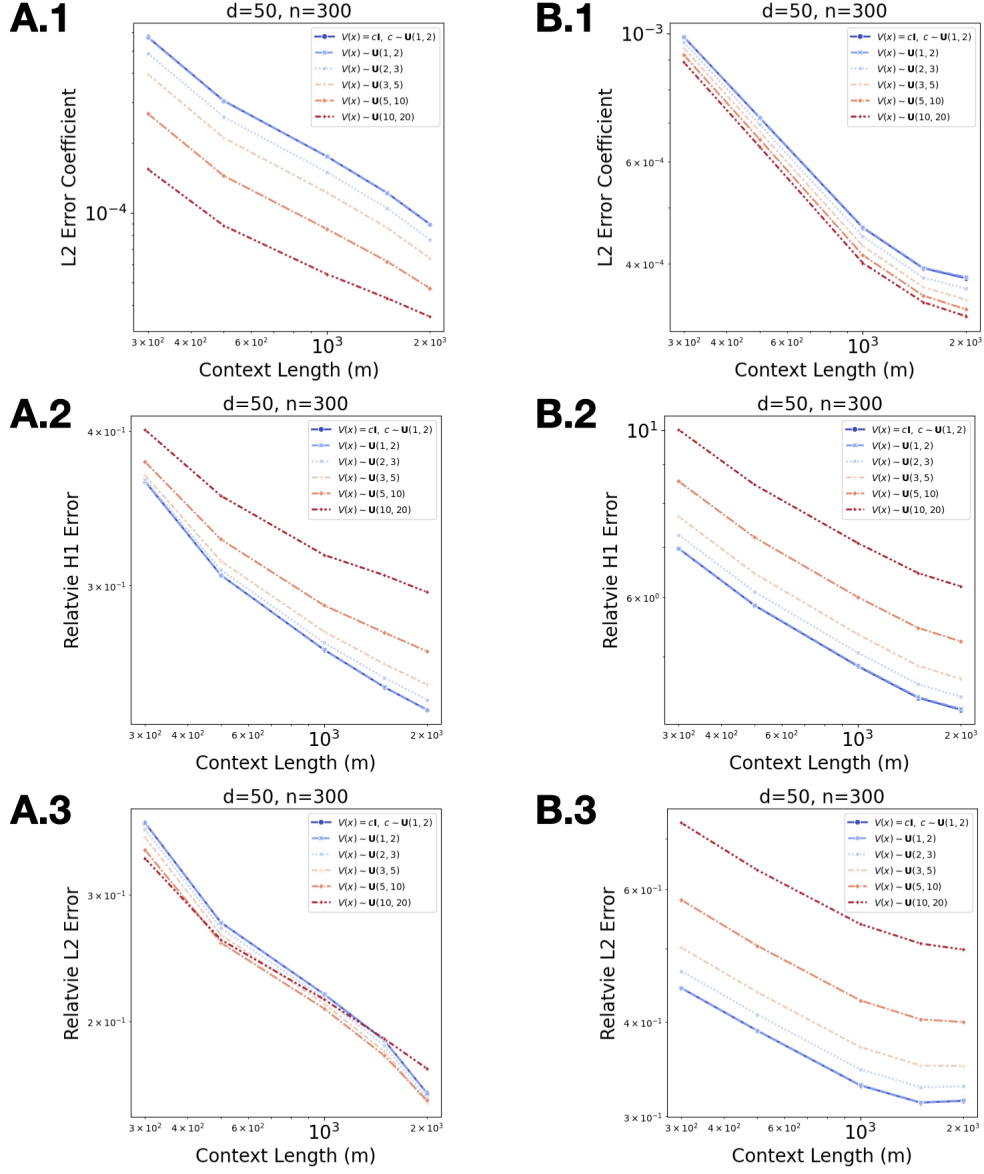


Figure 4: Column A shows the prediction errors of the transformer model with the attention parameters initialized around $(P, Q) = (\mathbf{I}_d, \mathbf{I}_d)$. Column B shows the prediction errors of the model with the attention parameters initialized around $(P, Q) = (D, D^{-1})$, where D is a diagonal matrix with diagonal term sampled from $\mathbf{U}(1, 5)$.

with non-constant coefficients is essentially the same as the error of solving the equation with constant coefficients; in contrast, when P is initialized away from the identity matrix, the test error becomes much larger when the PDE has non-constant coefficients. This demonstrates the claims of Theorem 6 and highlights the importance of task diversity in achieving OOD generalization.

5 Conclusion

In this work, we studied the ability of a transformer characterized by a single linear self-attention layer for learning the solution operator of a linear elliptic PDE in-context. We characterized the generalization error in learning the associated discrete linear systems and PDE solutions in terms of the number of pre-training task, the prompt length during pre-training and testing, the size of the discretization, and various distribution shifts on the PDE coefficients. We also conducted extensive numerical experiments to validate our theory. Several questions remain open for future investigations. First, it remains to check the validity of the centralizer condition in Theorem 5 for random matrices that arise from the discretization of PDEs with random coefficients. Our numerical experiments empirically confirmed the validity of such condition for elliptic PDEs with a wide range of random coefficients, but a rigorous proof is still lacking. Second, it would be interesting to establish analogous theory for nonlinear and time-dependent PDE problems. In these more complex settings, it is crucial to characterize the role that depth and nonlinearity play in the ability of transformers to approximate the PDE solution. Moreover, it is also important to quantify the prediction error under distribution shifts in the tasks and covariates. We plan to explore these directions and report the results in the future.

Acknowledgements

F. Cole and Y. Lu acknowledge the support from NSF (award DMS-2343135). The computations presented here were conducted on the Minnesota Supercomputing Institute at the University of Minnesota Twin Cities.

References

- J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Al-tenschmidt, S. Altman, S. Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- K. Ahn, X. Cheng, H. Daneshmand, and S. Sra. Transformers learn to implement pre-conditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- E. Akyürek, D. Schuurmans, J. Andreas, T. Ma, and D. Zhou. What learning algorithm is in-context learning? investigations with linear models. In *The Eleventh International Conference on Learning Representations*, 2022.
- B. An, J. Lyu, Z. Wang, C. Li, C. Hu, F. Tan, R. Zhang, Y. Hu, and C. Chen. Repulsive attention: Rethinking multi-head attention as bayesian inference. *arXiv preprint arXiv:2009.09364*, 2020.
- Y. Bai, F. Chen, H. Wang, C. Xiong, and S. Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *Advances in neural information processing systems*, 36, 2024.

- I. Batatia, P. Benner, Y. Chiang, A. M. Elena, D. P. Kovács, J. Riebesell, X. R. Advincula, M. Asta, W. J. Baldwin, N. Bernstein, et al. A foundation model for atomistic materials chemistry. *arXiv preprint arXiv:2401.00096*, 2023.
- K. Bhattacharya, B. Hosseini, N. B. Kovachki, and A. M. Stuart. Model reduction and neural networks for parametric pdes. *The SMAI journal of computational mathematics*, 7:121–157, 2021.
- D. Boffi. Finite element approximation of eigenvalue problems. *Acta numerica*, 19:1–120, 2010.
- S. Brenner and R. Scott. *The Mathematical Theory of Finite Element Methods*, volume 15. Springer Science & Business Media, 2007.
- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, H. Nori, H. Palangi, M. T. Ribeiro, and Y. Zhang. Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv:2303.12712*, 2023.
- A. Celaj, A. J. Gao, T. T. Lau, E. M. Holgersen, A. Lo, V. Lodaya, C. B. Cole, R. E. Denroche, C. Spickett, O. Wagih, et al. An rna foundation model enables discovery of disease mechanisms and candidate therapeutics. *bioRxiv*, pages 2023–09, 2023.
- S. Chen, H. Sheen, T. Wang, and Z. Yang. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv preprint arXiv:2402.19442*, 2024a.
- W. Chen, J. Song, P. Ren, S. Subramanian, D. Morozov, and M. W. Mahoney. Data-efficient operator learning via unsupervised pretraining and in-context learning. *arXiv preprint arXiv:2402.15734*, 2024b.
- X. Cheng, Y. Chen, and S. Sra. Transformers implement functional gradient descent to learn non-linear functions in context. *arXiv preprint arXiv:2312.06528*, 2023.
- F. Cole and Y. Lu. Score-based generative models break the curse of dimensionality in learning a family of sub-gaussian probability distributions. *arXiv preprint arXiv:2402.08082*, 2024.
- Z. Du, H. Balim, S. Oymak, and N. Ozay. Can transformers learn optimal filtering for unknown systems? *IEEE Control Systems Letters*, 7:3525–3530, 2023.
- R. M. Dudley. The sizes of compact subsets of hilbert space and continuity of gaussian processes. *Journal of Functional Analysis*, 1(3):290–330, 1967.
- B. L. Edelman, E. Edelman, S. Goel, E. Malach, and N. Tsilivis. The evolution of statistical induction heads: In-context learning markov chains. *arXiv preprint arXiv:2402.11004*, 2024.

- A. Ern and J.-L. Guermond. *Theory and practice of finite elements*, volume 159. Springer, 2004.
- S. Garg, D. Tsipras, P. S. Liang, and G. Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- G. Goel and P. Bartlett. Can a transformer represent a kalman filter? *arXiv preprint arXiv:2312.06937*, 2023.
- J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Y. Huang, Y. Cheng, and Y. Liang. In-context convergence of transformers. *arXiv preprint arXiv:2310.05249*, 2023.
- J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- J. Kim, T. Nakamaki, and T. Suzuki. Transformers are minimax optimal nonparametric in-context learners. *arXiv preprint arXiv:2408.12186*, 2024.
- L. Kirsch, J. Harrison, J. Sohl-Dickstein, and L. Metz. General-purpose in-context learning by meta-learning transformers. *arXiv preprint arXiv:2212.04458*, 2022.
- Y. Li, M. E. Ildiz, D. Papailiopoulos, and S. Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *International Conference on Machine Learning*, pages 19565–19594. PMLR, 2023.
- Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar. Fourier neural operator for parametric partial differential equations. *arXiv preprint arXiv:2010.08895*, 2020.
- L. Lin, Y. Bai, and S. Mei. Transformers as decision makers: Provable in-context reinforcement learning via supervised pretraining. *arXiv preprint arXiv:2310.08566*, 2023.
- L. Lu, P. Jin, and G. E. Karniadakis. Deeponet: Learning nonlinear operators for identifying differential equations based on the universal approximation theorem of operators. *arXiv preprint arXiv:1910.03193*, 2019.
- A. V. Mahankali, T. Hashimoto, and T. Ma. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. In *The Twelfth International Conference on Learning Representations*, 2024.
- M. McCabe, B. R.-S. Blancard, L. H. Parker, R. Ohana, M. Cranmer, A. Bietti, M. Eickenberg, S. Golkar, G. Krawezik, F. Lanusse, et al. Multiple physics pretraining for physical surrogate models. *arXiv preprint arXiv:2310.02994*, 2023.

- O. Méndez-Lucio, C. Nicolaou, and B. Earnshaw. Mole: a molecular foundation model for drug discovery. *arXiv preprint arXiv:2211.02657*, 2022.
- Y. Mroueh. Towards a statistical theory of learning to learn in-context with transformers. In *NeurIPS 2023 Workshop Optimal Transport and Machine Learning*, 2023.
- E. Nichani, A. Damian, and J. D. Lee. How transformers learn causal structure with gradient descent. *arXiv preprint arXiv:2402.14735*, 2024.
- J. Park, I. Pelakh, and S. Wojtowytsch. Minimum norm interpolation by perceptrons: Explicit regularization and implicit bias. *Advances in Neural Information Processing Systems*, 36, 2023.
- M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics*, 378:686–707, 2019.
- M. Rudelson and R. Vershynin. Hanson-wright inequality and sub-gaussian concentration. 2013.
- M. E. Sander, R. Giryes, T. Suzuki, M. Blondel, and G. Peyré. How do transformers perform in-context autoregressive learning? *arXiv preprint arXiv:2402.05787*, 2024.
- S. Shalev-Shwartz and S. Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- J. Shen, T. Tang, and L.-L. Wang. *Spectral methods: algorithms, analysis and applications*, volume 41. Springer Science & Business Media, 2011.
- G. Strang. *Introduction to linear algebra*. SIAM, 2022.
- S. Subramanian, P. Harrington, K. Keutzer, W. Bhimji, D. Morozov, M. W. Mahoney, and A. Gholami. Towards foundation models for scientific machine learning: Characterizing scaling and transfer behavior. *Advances in Neural Information Processing Systems*, 36, 2024.
- J. Sun, Y. Liu, Z. Zhang, and H. Schaeffer. Towards a foundation model for partial differential equation: Multi-operator learning and extrapolation. *arXiv preprint arXiv:2404.12355*, 2024.
- N. Tripuraneni, M. Jordan, and C. Jin. On the theory of transfer learning: The importance of task diversity. *Advances in neural information processing systems*, 33:7852–7862, 2020.
- A. W. Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

- J. Von Oswald, E. Niklasson, E. Randazzo, J. Sacramento, A. Mordvintsev, A. Zhmoginov, and M. Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023.
- M. J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- J. Wu, D. Zou, Z. Chen, V. Braverman, Q. Gu, and P. L. Bartlett. How many pre-training tasks are needed for in-context learning of linear regression? *arXiv preprint arXiv:2310.08391*, 2023.
- S. M. Xie, A. Raghunathan, P. Liang, and T. Ma. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- L. Yang and S. J. Osher. Pde generalization of in-context operator networks: A study on 1d scalar nonlinear conservation laws. *arXiv preprint arXiv:2401.07364*, 2024.
- L. Yang, S. Liu, T. Meng, and S. J. Osher. In-context operator learning with data prompts for differential equation problems. *Proceedings of the National Academy of Sciences*, 120(39):e2310142120, 2023.
- Z. Ye, X. Huang, L. Chen, H. Liu, Z. Wang, and B. Dong. Pdeformer: Towards a foundation model for one-dimensional partial differential equations. *arXiv preprint arXiv:2402.12652*, 2024.
- R. Zhang, S. Frei, and P. L. Bartlett. Trained transformers learn linear models in-context. *arXiv preprint arXiv:2306.09927*, 2023a.
- Y. Zhang, F. Zhang, Z. Yang, and Z. Wang. What and how does in-context learning learn? bayesian model averaging, parameterization, and generalization. *arXiv preprint arXiv:2305.19420*, 2023b.

Part I

Appendix

Table of Contents

A	Notation	25
B	Proofs for Subsection 3.1	26
C	Proofs and additional results for Subsection 3.2	34
D	Proofs and additional results for Subsection 3.3	35
E	Proofs for Subsection 3.4	39
F	Discussion on dependence of constants on dimension	41
G	Auxiliary lemmas	41
H	Additional numerical results	48

Appendix A. Notation

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

- Physical dimension of PDE problem: d_0
- Dimension of task matrix for ICL: d
- Task matrix for ICL: A
- Covariates for ICL: $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$
- Prompt matrix for ICL: Z
- Empirical covariance matrix of $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$: Y_n
- Distribution on tasks: p_A
- Upper bound on largest eigenvalue of A^{-1} over $\text{supp}(p_A)$: c_A
- Covariance operator of the distribution on $L^2(\Omega)$ -valued covariates: Σ_f
- Covariance matrix of the distribution on \mathbb{R}^d -valued covariates: Σ
- Parameters of transformer: $\theta = (P, Q)$
- Prediction of the transformer with parameters θ : $\text{TF}_\theta(Z)$
- Population risk for training: \mathcal{R}_n
- Population risk for inference: \mathcal{R}_m
- Empirical risk: $\mathcal{R}_{n,N}$

- "Infinite-context" population risk: \mathcal{R}_∞
- Number of context examples per prompt during training: n
- Number of context examples per prompt during inference: m
- Number of pre-training tasks: N

Appendix B. Proofs for Subsection 3.1

In this section we prove Theorem 1, which controls the (in-distribution) generalization error for in-context learning of linear systems in terms of the context length during training, the context length during inference, and the number of pre-training tasks.

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

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

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

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

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

$$\mathcal{R}_m(\hat{\theta}) = \left(\mathcal{R}_m(\hat{\theta}) - \mathcal{R}_m^t(\hat{\theta}) \right) + \left(\mathcal{R}_m^t(\hat{\theta}) - \mathcal{R}_{m, N}^t(\hat{\theta}) \right) + \left(\mathcal{R}_{m, N}^t(\hat{\theta}) - \mathcal{R}_{m, N}^t(\theta^*) \right) \quad (15)$$

$$+ \left(\mathcal{R}_{m, N}^t(\theta^*) - \mathcal{R}_m^t(\theta^*) \right) + \left(\mathcal{R}_m^t(\theta^*) - \mathcal{R}_m(\theta^*) \right) + \mathcal{R}_m(\theta^*) \quad (16)$$

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

$$+ \left(\mathcal{R}_{m, N}^t(\hat{\theta}) - \mathcal{R}_{m, N}^t(\theta^*) \right) + \inf_{\|\theta^*\| \leq M} \mathcal{R}(\theta^*). \quad (18)$$

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

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

$$+ \sup_{\|\theta\| \leq M} \left(\mathcal{R}_m(\theta) - \mathcal{R}_m^t(\theta) \right) + \sup_{\|\theta\| \leq M} \left(\mathcal{R}_n(\theta) - \mathcal{R}_n^t(\theta) \right) \quad (20)$$

$$+ 2 \sup_{\|\theta\| \leq M} \left| \mathcal{R}_m(\theta) - \mathcal{R}_n(\theta) \right| + \left(\mathcal{R}_{n,N}^t(\widehat{\theta}) + \mathcal{R}_{n,N}^t(\theta^*) \right). \quad (21)$$

Plugging the estimate 19 into the bound from (15) gives the final bound

$$\mathcal{R}_m(\widehat{\theta}) \leq 2 \underbrace{\sup_{\|\theta\| \leq M} \left(\mathcal{R}_m - \mathcal{R}_m^t \right)(\theta) + \sup_{\|\theta\| \leq M} \left(\mathcal{R}_n - \mathcal{R}_n^t \right)(\theta)}_{\text{data truncation error}} \quad (22)$$

$$+ 4 \underbrace{\sup_{\|\theta\| \leq M} \left(\mathcal{R}_m^t - \mathcal{R}_{m,N}^t \right)(\theta) + 2 \sup_{\|\theta\| \leq M} \left(\mathcal{R}_n^t - \mathcal{R}_{n,N}^t \right)(\theta)}_{\text{statistical error}} \quad (23)$$

$$+ 2 \underbrace{\sup_{\|\theta\| \leq M} \left| \mathcal{R}_m(\theta) - \mathcal{R}_n(\theta) \right|}_{\text{context mismatch error}} + \underbrace{\left(\mathcal{R}_{n,N}^t(\widehat{\theta}) - \mathcal{R}_{n,N}^t(\theta^*) \right)}_{\leq 0 \text{ w.h.p.}} + \underbrace{\mathcal{R}_m(\theta^*)}_{\text{approx. error}} \quad (24)$$

$$= I + II + III + IV + V. \quad (25)$$

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

Step 2 - bounding the truncation error: By Lemma 7 and Example 6.2 in Wainwright (2019), when $\mathbf{y} \sim N(0, \Sigma)$ 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

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

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

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

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

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

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

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

where $\ell_\theta(A, Y_n, \mathbf{y}) = \|(PA^{-1}Y_nQ - A^{-1})\mathbf{y}\|^2$. By Theorem 26.5 in Shalev-Shwartz and Ben-David (2014), we have with probability at least $1 - \delta$,

$$\sup_{\|\theta\| \leq M} \left| \mathbb{E}_{A,Y_n,\mathbf{y}} \ell_\theta(A, Y_n, \mathbf{y}) - \sum_{i=1}^N \ell_\theta(A_i, Y_n^{(i)}, \mathbf{y}_i) \right| \leq \operatorname{Rad}_N(\{\ell_\theta : \|\theta\| \leq M\}) \quad (26)$$

$$+ \sup_{\|\theta\| \leq M} \|\ell_\theta\|_\infty \cdot \sqrt{\frac{2 \log(1/\delta)}{N}}, \quad (27)$$

where the $\operatorname{Rad}_N(\{\ell_\theta : \|\theta\| \leq M\})$ is the Rademacher complexity defined by

$$\operatorname{Rad}_N(\{\ell_\theta : \|\theta\| \leq M\}) = \mathbb{E}_{\epsilon_i \sim \text{unif}(\{\pm 1\})} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \ell_\theta(A_i, Y_n^{(i)}, \mathbf{y}_i),$$

and the expectations over Y_n and \mathbf{y} are taken over the truncated versions of the original distributions. Note that on the event $\mathcal{B}_{N,t}$, we have

$$\|(PA_i^{-1}Y_n^{(i)}Q - A_i^{-1})\mathbf{y}_i\|^2 \leq M^4 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}} \right)^2 \left(\sqrt{\operatorname{Tr}(\Sigma)} + t \right)^2,$$

and hence the second term in Inequality (26) is bounded from above by

$$M^4 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}}\right)^2 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^2 \cdot \sqrt{\frac{2 \log(1/\delta)}{N}}.$$

It remains to bound the Rademacher complexity of $\{\ell_\theta : \|\theta\| \leq M\}$. Notice that by the triangle inequality, we have

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

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

$$\begin{aligned} \Theta_1(M) &= \{(A, Y_n, \mathbf{y}) \mapsto \|PA^{-1} Y_n Q \mathbf{y}\|^2 : \|\theta\| \leq M\}, \\ \Theta_2(M) &= \{(A, Y_n, \mathbf{y}) \mapsto \langle PA^{-1} Y_n Q \mathbf{y}, A^{-1} \mathbf{y} \rangle : \|\theta\| \leq M\}. \end{aligned}$$

By Dudley's integral theorem Dudley (1967), it holds that

$$\mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \|PA_i^{-1} Y_n^{(i)} Q \mathbf{y}_i\|^2 \leq \inf_{\epsilon > 0} \frac{12\sqrt{2}}{\sqrt{N}} \int_{\epsilon}^{D_1(M)} \sqrt{\log \mathcal{N}(\Theta_1(M), \|\cdot\|_N, \tau)} d\tau, \quad (28)$$

where $\mathcal{N}(\Theta_1(M), \|\cdot\|_N, \tau)$ is the τ -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

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

Note the bound

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

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 $\theta_1 = (P_1, Q_1), \theta_2 = (P_2, Q_2)$, 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

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

Optimizing over the choice of ϵ in Equation (28), this proves that

$$\mathbb{E}_{\epsilon_i} \sup_{\|\theta\| \leq M} \frac{1}{N} \sum_{i=1}^N \epsilon_i \|PA_i^{-1}Y_n^{(i)}Q\mathbf{y}_i\|^2 \quad (29)$$

$$= O\left(\frac{d^2 M^4 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}}\right)^2 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^2}{\sqrt{N}}\right), \quad (30)$$

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

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

which in turn bounds the Rademacher complexity $\text{Rad}_N(\{\ell_\theta : \|\theta\| \leq M\})$ by the right-hand side of Equation (29). Combining this Rademacher complexity estimate with the overall statistical error bound (26), we conclude that

$$\begin{aligned} &\sup_{\|\theta\| \leq M} \left| \mathcal{R}_n^t(\theta) - \mathcal{R}_{n,N}^t(\theta) \right| \\ &= O\left(\frac{(d^2 + \sqrt{2\log(1/\delta)})M^4 c_A^2 \|\Sigma\|_{\text{op}}^2 \left(1 + t + \sqrt{\frac{d}{n}}\right)^2 \left(\sqrt{\text{Tr}(\Sigma)} + t\right)^2}{\sqrt{N}}\right) \end{aligned}$$

holds with probability at least

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

by a union bound. The same argument proves in analogous bound on the statistical error term

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

where n is replaced by m in the bound of Equation (29).

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

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

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

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

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

Step 7 - Balancing error terms: Putting everything together and applying the error decomposition from step 1, we have shown that ⁴

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

with probability at least

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

We choose t and δ such that

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

4. 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 .

It is clear that for this to be satisfied, we can take both t and δ to be logarithmic in N . For such t and δ , we have, omitting universal constants and $\log(N)$ factors, that

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

We omit the third term from the final bound, since, asymptotically, it is dominated by the fourth term. \blacksquare

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 first observe that for $\theta = (P, Q)$, the output of the transformer TF_{θ} 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

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

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 $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$. Note that if we set $P = Id$ and $Q = \Sigma^{-1}$ to be the inverse of the data covariance matrix, then for sufficiently large m we have $\text{TF}_{\theta}(Z) \approx A^{-1} \mathbf{y}$. This suggests that the transformer can learn to solve linear systems in a way that is extremely robust to shifts in the distribution on the task matrices. We note that similar choices of attention matrices have been studied in the linear regression setting (Ahn et al. (2024), Zhang et al. (2023a)). Our result essentially employs the parameterization $P = \mathbf{I}_d$ and $Q = \Sigma^{-1}$, but with an additional bias term to account for the fact that the sequence 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

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

where $\sigma_1^2, \dots, \sigma_d^2$ are the eigenvalues of Σ and $\varphi_i = W e_i$ are the eigenvectors. Note that the weighted 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

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

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

$$\|K\|_{\text{op}} \leq \|\Sigma^{-1}\|_{\text{op}} \|\Sigma\|_{\text{op}} \left(1 + \text{Tr}_{\Sigma}(B)\right) C_A^2. \quad (31)$$

It follows that

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

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

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

where we used the fact that $\text{Tr}_{\Sigma}(\Sigma^{-1}) = d$. Using the bound on the norm of K stated in Equation (31), 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

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

■

To justify our ansatz for upper bounding the approximation error (i.e., how the matrix Q_n in Proposition 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*

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

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

Proof Let us recall the definition of the population risk functional

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

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 A and $\{\mathbf{y}_i\}_{i=1}^n$, $A^{-1} (Y_n Q - I) y$ is a centered Gaussian random vector with covariance $A^{-1} (Y_n Q - I) \Sigma (Q Y_n - I) A^{-1}$. In addition, since the task and data distributions are independent, we can replace 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

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

This equation is solved by the matrix Q_n defined in the statement of the Lemma. \blacksquare

Appendix C. Proofs and additional results for Subsection 3.2

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

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

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

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

and therefore that

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

■

Appendix D. Proofs and additional results for Subsection 3.3

Proof of Theorem 3. Recall that $\hat{\theta} \in \operatorname{argmin}_{\|\theta\| \leq M} \mathcal{R}_{n,N}(\theta)$ is the ERM. Let $\theta_* = (P_*, Q_*)$ denote a projection of $\hat{\theta}$ onto the set \mathcal{M}_∞ and let $\theta'_* = (P'_*, Q'_*)$ denote a projection of $\hat{\theta}$ onto \mathcal{M}_∞ . Let $\epsilon_1 = \|\hat{\theta} - \theta_*\|$ and $\epsilon_2 = \|\hat{\theta} - \theta'_*\|$. Then we have the error decomposition

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

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

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

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

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

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

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

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

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

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 distribution 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 $\theta = (P, Q)$,

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

and

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

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 its minimum value, we have

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

It follows that

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

where, again, the infimum is taken over all $\theta_* \in \operatorname{argmin}_{\theta \in \mathcal{M}_{\infty}(p_A)} \|\theta - \hat{\theta}\|^2$ and $\theta'_* \in \operatorname{argmin}_{\theta' \in \mathcal{M}_{\infty}(p'_A)} \|\theta' - \hat{\theta}\|^2$. Combining the estimates for each individual term in the error decomposition, we obtain the final bound in the statement of Theorem 3. The fact that the bound we have obtained tends to zero as the sample size $(m, n, N) \rightarrow \infty$ follows from examination of each term in the estimate: the in-domain generalization error $\mathcal{R}_m(\hat{\theta})$ tends to zero in probability by 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 to zero as N and n tend to infinity, respectively, by Proposition 3. \blacksquare

The discrepancy $d(p_A, p'_A)$ defined in the proof of Theorem 4 may not be a metric, but, crucially, it satisfies $d(p_A, p_A) = 0$. This ensures that the error term due to distribution shift in Theorem 4 vanishes when the pre-training and downstream tasks coincide. We give a simple proof of this fact 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 θ_* and θ'_* of $\hat{\theta}$ onto the sets $\mathcal{M}_{\infty}(p_A)$ and $\mathcal{M}_{\infty}(p'_A)$ respectively, and

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

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 θ_* of $\widehat{\theta}$ onto $\mathcal{M}_\infty(p_A)$. Clearly we have

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

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

The next proposition gives a characterization of the minimizers of the functionals \mathcal{R}_∞ and \mathcal{R}'_∞ . Apart from being interesting in its own right, it is a key tool to prove Theorem 5.

Proposition 2. *Fix a task distribution p_A satisfying Assumption 1. Then $\theta = (P, Q)$ is a minimizer of \mathcal{R}_∞ 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 is given by $Q = \Sigma^{-1} A_0 P^{-1} A_0^{-1}$ for any $A_0 \in \text{supp}(p_A)$.*

Proof Recall that

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

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}_∞ 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 ϵ 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

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

for each $\epsilon > 0$. Since $p_{A, \epsilon}(A_0)$ converges weakly to the Dirac measure centered at A_0 , we have 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 concludes the first part of the proof.

Now, suppose $\theta = (P, Q)$ is a minimizer of \mathcal{R}_∞ . By the previous argument, this is equivalent to the system of equations $PA^{-1}\Sigma Q = A^{-1}$ holding simultaneously for all $A \in \text{supp}(p_A)$. In particular, 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 $Q = \Sigma^{-1} A_0 P^{-1} A_0^{-1}$. Since the matrix Q is constant, this implies that the function $A \mapsto AP^{-1}A^{-1}$ is a constant on the support of p_A . We have therefore shown that the minimizers of \mathcal{R}_∞ can be 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 $\text{supp}(p_A)$. In addition, the fact that the function $A \mapsto AP^{-1}A^{-1}$ is constant on the support of p_A implies that P commutes with all products of the form $\{A_1 A_2^{-1} : A_1, A_2 \in \text{supp}(p_A)\}$. \blacksquare

We now give a proof of Theorem 5.

Proof of Theorem 5. 1) This is a direct corollary of Proposition 2.

2) Let $\theta_* = (P_*, Q_*)$ be a minimizer of \mathcal{R}_∞ . Then Proposition 2 implies that $P_* \in \mathcal{C}(\mathcal{S}(p_A))$. Since the centralizer of $\mathcal{S}(p_A)$ is trivial by assumption, this implies that $P_* = c\mathbf{I}_d$

for some $c \in \mathbb{R} \setminus \{0\}$. Using the characterization of minimizers of \mathcal{R}_∞ derived in Proposition 2, we have that Q_* solves the equation $cA^{-1}\Sigma Q_* = A^{-1}$ for all $A \in \text{supp}(p_A)$, and therefore $Q = c^{-1}\Sigma^{-1}$. \blacksquare

Proof of Corollary 1. By combining Theorems 4 and 5, 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 θ_* 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

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

Before proving Theorem 6, 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*

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

Proof By Proposition 2, a parameter (P, Q) belongs to $\mathcal{M}_\infty(p_A)$ if and only if P commutes 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 $Q = \Sigma^{-1}A_0 P^{-1}A_0^{-1}$ 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 of the lemma. Since P and $A_1 A_2^{-1}$ are commuting diagonalizing matrices and $A_1 A_2^{-1}$ has no repeated eigenvalues (Strang (2022)), they must be simultaneously diagonalizable. This implies that P is diagonal in the basis U , and hence Q is given by $Q = \Sigma^{-1}A_0 P^{-1}A_0^{-1} = \Sigma^{-1}P^{-1}$. \blacksquare

Proof of Theorem 6. For 1), if the support of p'_A is also simultaneously diagonalizable with respect to U , then Lemma 3 implies that $\mathcal{M}_\infty(p_A) = \mathcal{M}_\infty(p'_A) = \Theta_{U, \Sigma}$, where $\Theta_{U, \Sigma}$ is as defined in the statement of Lemma 3. This proves that if the support of p'_A is also simultaneously diagonalizable with respect to U , then p_A is diverse.

For 2), we must find a minimizer of \mathcal{R}_∞ which is not a minimizer of \mathcal{R}'_∞ . Consider the parameter $\theta = (P, \Sigma^{-1}P^{-1})$, where $P = UDU^T$ for D an invertible diagonal matrix

with no repeated entries. By Lemma 3, θ is a minimizer of \mathcal{R}_∞ . Let $A'_3, A'_4 \in \text{supp}(p'_A)$ be such that $A'_3(A'_4)^{-1}$ is not diagonalizable with respect to U . Since $A'_3(A'_4)^{-1}$ and P are not simultaneously diagonalizable and P has no repeated eigenvalues (Strang (2022)), P does not commute with $A'_3(A'_4)^{-1}$. By Proposition 2, θ is therefore not a minimizer of \mathcal{R}'_∞ , completing the proof. \blacksquare

Appendix E. Proofs for Subsection 3.4

We begin by stating a more formal version of Theorem 7 where the constants are more explicit.

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 minimizers of the empirical risk when the in-context examples follow the distribution $N(0, \Sigma)$ and 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 7 then follows from Theorem 8 by bounding $\|\Lambda - \tilde{\Lambda}\|_1 \lesssim \|\Sigma - \tilde{\Sigma}\|_{\text{op}}$, merging the term

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

into the second term, and omitting the constant factors.

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

$$\mathcal{R}_m^{\tilde{\Sigma}}(\hat{P}, \hat{Q}) \leq \mathcal{R}_m^{\Sigma}(\hat{P}, \hat{Q}) + \sup_{\|P\|_{\text{op}}, \|Q\|_{\text{op}} \leq M} \left| \mathcal{R}_m^{\tilde{\Sigma}}(P, Q) - \mathcal{R}_m^{\Sigma}(P, Q) \right|. \quad (32)$$

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

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

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

Similarly, we have

$$\mathcal{R}_m^{\tilde{\Sigma}}(P, Q) = \mathbb{E}_A \left[\frac{m+1}{m} \text{Tr}(PA^{-1}\tilde{\Sigma} Q\tilde{\Sigma} Q^T \tilde{\Sigma} A^{-1}P^T) + \frac{\text{Tr}_{\tilde{\Sigma}}(Q\tilde{\Sigma} Q^T)}{m} \text{Tr}(PA^{-1}\tilde{\Sigma} A^{-1}P^T) \right] \quad (35)$$

$$+ \mathbb{E}_A \left[\text{Tr}(A^{-1}\tilde{\Sigma} A^{-1}) - \text{Tr}(PA^{-1}\tilde{\Sigma} Q\tilde{\Sigma} A^{-1}) - \text{Tr}(A^{-1}\tilde{\Sigma} Q^T \tilde{\Sigma} A^{-1}P^T) \right]. \quad (36)$$

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 term appearing in the expressions for \mathcal{R}_m^Σ and $\mathcal{R}_m^{\tilde{\Sigma}}$. By a simple applications of Hölder's inequality 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

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

and

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

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

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

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

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

Similarly, using Lemma 11, the first term satisfies

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

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

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

■

Appendix F. Discussion on dependence of constants on dimension

It is important to consider the dependence of the constants appearing in Theorem 1 on the dimension of the linear system. Recall that in the PDE setting, the dimension d corresponds to the number of basis functions used in Galerkin's method, and hence the true PDE solution is only recovered in the 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 Σ are dimension-independent, since we always have

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

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

Appendix G. Auxiliary lemmas

We make frequent use of the following lemma to compute expectations of products of empirical 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*

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

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 Σ .

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

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

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

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

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

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

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 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 $\{W^T y_1, \dots, W^T y_n\}$. Noting that $W^T y \sim N(0, \Lambda)$ for $y \sim N(0, \Sigma)$, we can apply the above result to $W^T K W$:

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

■

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

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

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

Proof This follows from a direct computation of the expectation with respect to Y_n :

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

where we used Lemma 4 to compute the expectation over Y_n in the second-to-last line. ■

It will also be useful to derive a simpler expression for the population risk $\mathcal{R}_m(\theta)$ when θ belongs to 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*

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

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

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

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

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*

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

where $C > 0$ is a constant independent of Σ and d .

We use the following result to control the error between Q_n and Σ^{-1} .

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

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

Then we can write

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

where \mathcal{E}_1 satisfies

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

Proof Using some algebra, we find

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

By Lemma 9, we have

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

where

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

This gives the final bound

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

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

The following result, used to bound the inverse of a perturbed matrix, is a standard application of 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

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

It follows that

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

In turn, this gives the bound

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

■

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

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

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

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,

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

■

In order to prove Theorem 7, we also need the following stability bound of $\text{Tr}_{\Sigma}(K)$ with respect to perturbations of both Σ 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 Σ and $\tilde{\Sigma}$ and let $\{\varphi_i\}_{i=1}^d$ and $\{\tilde{\varphi}_i\}_{i=1}^d$ be the respective eigenvectors. Then*

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

Proof We have

$$\text{Tr}_{\Sigma}(K) - \text{Tr}_{\tilde{\Sigma}}(\tilde{K}) \leq \left| \text{Tr}_{\Sigma}(K) - \text{Tr}_{\tilde{\Sigma}}(K) \right| + \left| \text{Tr}_{\tilde{\Sigma}}(K - \tilde{K}) \right|. \quad (37)$$

The second term in (37) can be bounded by an application of Lemma 10, which yields

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

To bound the first term in (37), we first use the estimate

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

The first term above can be bounded by

$$\left| \sum_{i=1}^d \left(\sigma_i^2 - \tilde{\sigma}_i^2 \right) \langle K\varphi_i, \varphi_i \rangle \right| \leq \|K\|_{\text{op}} \cdot \sum_{i=1}^d \left| \sigma_i^2 - \tilde{\sigma}_i^2 \right| = \|K\|_{\text{op}} \cdot \|\Lambda - \tilde{\Lambda}\|_1. \quad (38)$$

To bound the second term in (38), note that for any $1 \leq i \leq d$, we have

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

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

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

Combining all terms yields the final estimate

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

■

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

Lemma 12. *The bound*

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

holds.

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

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

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

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

■

The following lemma is an adaptation of Wald's consistency theorem of M-estimators (Van der Vaart, 2000, Theorem 5.14). We use it to prove the convergence in probability of empirical risk minimizers 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 semi-continuous in θ . 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*

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

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

Proposition 3. *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}_{∞} are polynomials in θ such that the coefficients of \mathcal{R}_n converge to the coefficients of \mathcal{R}_{∞} as $n \rightarrow \infty$, we have by the triangle inequality that

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

■

Appendix H. Additional numerical results

In this section, we present some additional numerics. The plots in Figure 5: A.1-C.1 are identical 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 function of the sample size. This makes it easier to compare the empirical scaling laws with those derived in Theorem 1. Figure 6: A.1-A.3 shows the heat maps for the PDE error with respect to the parameters α and τ that define the log-normal random field $a(x)$, while Figure 6: B.1-B.3 shows the heat maps for the PDE error with respect to the parameters c and β of the covariance of the source term f . These plots confirm that pre-trained transformers are more robust under task shifts than they are under covariate shifts. They also suggest that the pre-trained transformer is better at tolerating shifts on the parameter c of the covariance operator compared to shifts on β .

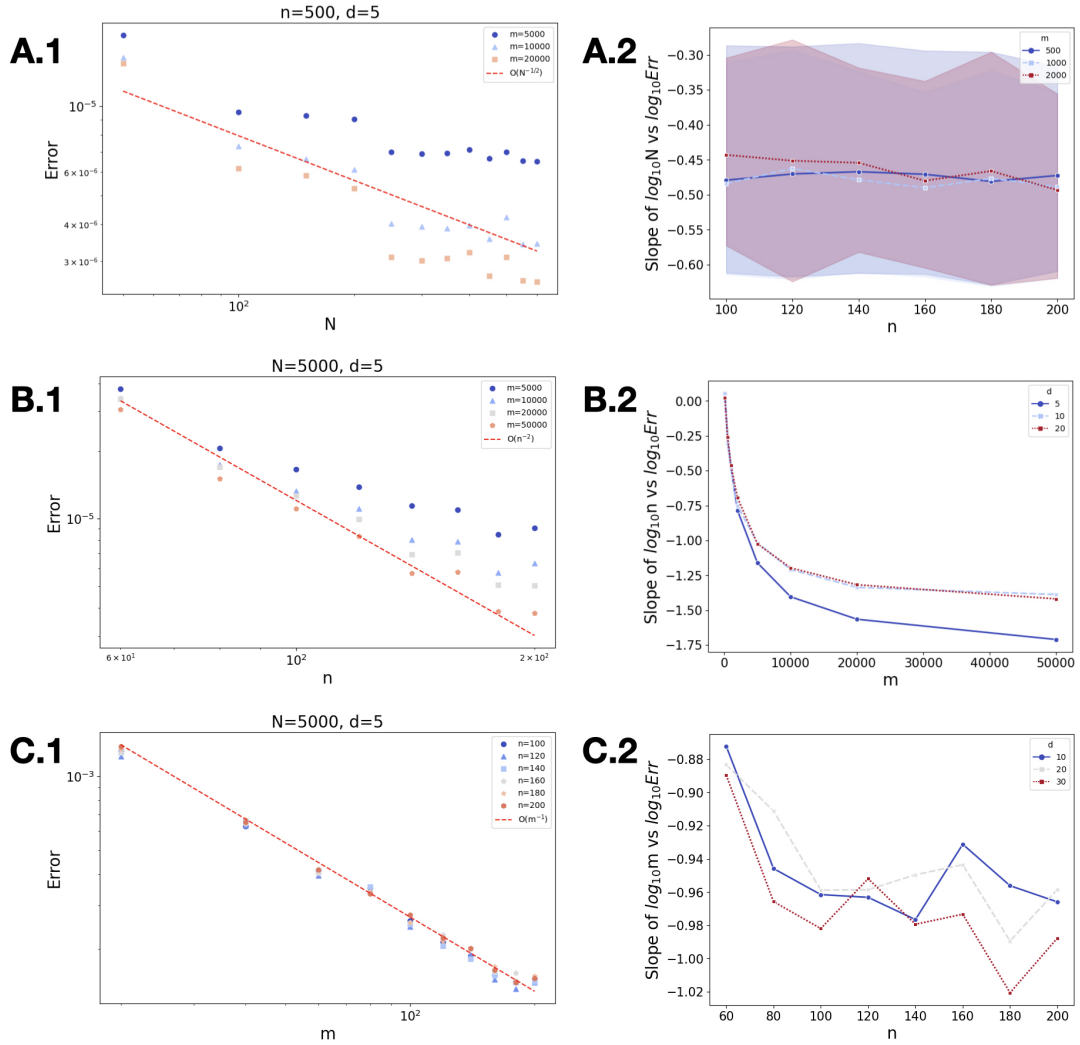


Figure 5: Plots A.1-C.1 are identical to those shown in Figure 1. Plots A.2-C.2 show the slopes of the error curves in the left column as functions of various sample sizes.

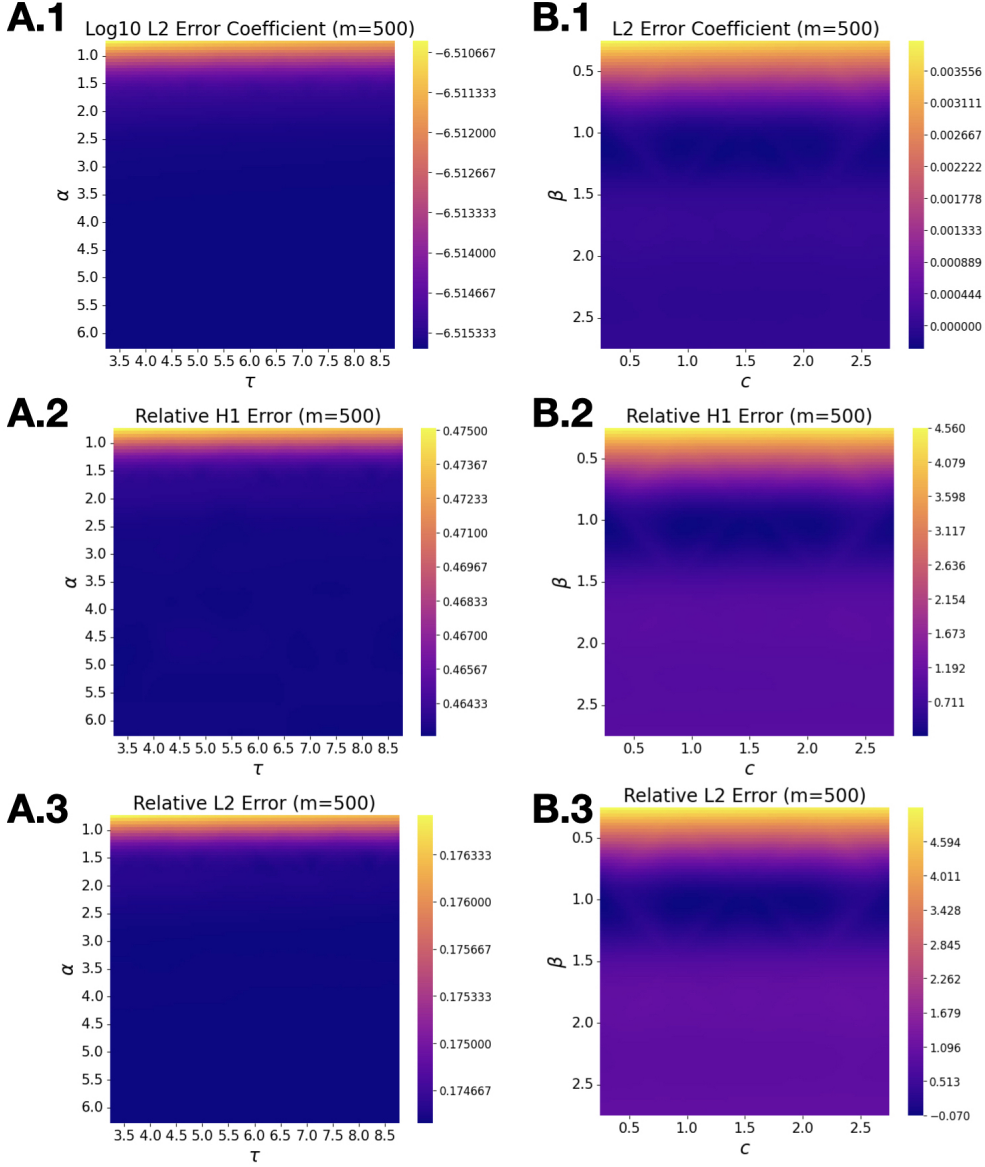


Figure 6: Column A shows the heat map for the error with respect to the parameters α and τ on the distribution of $a(x)$ (the training distribution is $a(x) = e^{b(x)}$ with $b(x) \sim N(0, (-\Delta + \tau \mathbb{I})^{-\alpha})$, $\alpha = 3$ and $\tau = 5$). Column B plots the heat map for the error with respect to the parameters β and c on the distribution of the data $f(x)$ (the training distribution is $f(x) = b(x)$ with $b(x) \sim N(0, (-\Delta + c \mathbb{I})^{-\beta})$, $\beta = c = 1$). The transformer model is trained with $n = 300$ and $d = 50$, and error is computed on a new task with $m = 500$ prompts.