

Trained Transformers Learn Linear Models In-Context

Ruiqi Zhang

*Department of Statistics
University of California, Berkeley
367 Evans Hall, Berkeley, CA 94720-3860, USA*

RQZHANG@BERKELEY.EDU

Spencer Frei

*Department of Statistics
University of California, Davis
4118 Mathematical Sciences Building
399 Crocker Ave., Davis, CA 95616, USA*

SFREI@UCDAVIS.EDU

Peter L. Bartlett

*Department of Statistics and Department of Electrical Engineering and Computer Sciences
University of California, Berkeley
367 Evans Hall, Berkeley, CA 94720-3860, USA
Google DeepMind
1600 Amphitheatre Parkway
Mountain View, CA 94040, USA*

PETER@BERKELEY.EDU

Editor: Daniel Hsu

Abstract

Attention-based neural networks such as transformers have demonstrated a remarkable ability to exhibit in-context learning (ICL): Given a short prompt sequence of tokens from an unseen task, they can formulate relevant per-token and next-token predictions without any parameter updates. By embedding a sequence of labeled training data and unlabeled test data as a prompt, this allows for transformers to behave like supervised learning algorithms. Indeed, recent work has shown that when training transformer architectures over random instances of linear regression problems, these models' predictions mimic those of ordinary least squares.

Towards understanding the mechanisms underlying this phenomenon, we investigate the dynamics of ICL in transformers with a single linear self-attention layer trained by gradient flow on linear regression tasks. We show that despite non-convexity, gradient flow with a suitable random initialization finds a global minimum of the objective function. At this global minimum, when given a test prompt of labeled examples from a new prediction task, the transformer achieves prediction error competitive with the best linear predictor over the test prompt distribution. We additionally characterize the robustness of the trained transformer to a variety of distribution shifts and show that although a number of shifts are tolerated, shifts in the covariate distribution of the prompts are not. Motivated by this, we consider a generalized ICL setting where the covariate distributions can vary across prompts. We show that although gradient flow succeeds at finding a global minimum in this setting, the trained transformer is still brittle under mild covariate shifts. We complement this finding with experiments on large, nonlinear transformer architectures which we show are more robust under covariate shifts.

Keywords: in-context learning, transformers, neural networks, self-attention, generalization

1. Introduction

Transformer-based neural networks have quickly become the default machine learning model for problems in natural language processing, forming the basis of chatbots like ChatGPT (OpenAI, 2023), and are increasingly popular in computer vision (Dosovitskiy et al., 2021). These models can take as input sequences of tokens and return relevant next-token predictions. When trained on sufficiently large and diverse datasets, these models are often able to perform *in-context learning* (ICL): when given a short sequence of input-output pairs (called a *prompt*) from a particular task as input, the model can formulate predictions on test examples without having to make any updates to the parameters in the model.

Recently, Garg et al. (2022) initiated the investigation of ICL from the perspective of learning particular function classes. At a high-level, this refers to when the model has access to instances of prompts of the form $(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})$ where x_i, x_{query} are sampled i.i.d. from a distribution \mathcal{D}_x and h is sampled independently from a distribution over functions in a function class \mathcal{H} . The transformer succeeds at in-context learning if when given a new prompt $(x'_1, h'(x'_1), \dots, x'_N, h'(x'_N), x'_{\text{query}})$ corresponding to an independently sampled h' it is able to formulate a prediction for x'_{query} that is close to $h'(x'_{\text{query}})$ given a sufficiently large number of examples N . The authors showed that when transformer models are trained on prompts corresponding to instances of training data from a particular function class (e.g., linear models, neural networks, or decision trees), they succeed at in-context learning, and moreover the behavior of the trained transformers can mimic those of familiar learning algorithms like ordinary least squares.

Following this, a number of follow-up works provided constructions of transformer-based neural network architectures which are capable of achieving small prediction error for query examples when the prompt takes the form $(x_1, \langle w, x_1 \rangle, \dots, x_N, \langle w, x_N \rangle, x_{\text{query}})$ where $x_i, x_{\text{query}}, w \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$ (von Oswald et al., 2022; Akyürek et al., 2022). However, this leaves open the question of how it is that *gradient-based optimization algorithms* over transformer architectures produce models which are capable of in-context learning.¹

In this work, we investigate the learning dynamics of gradient flow in a simplified transformer architecture when the training prompts consists of random instances of linear regression datasets. Our main contributions are as follows.

- We establish that for a class of transformers with a single layer and with a linear self-attention module (LSAs), gradient flow on the population loss with a suitable random initialization converges to a global minimum of the population objective, despite the non-convexity of the underlying objective function.
- We characterize the learning algorithm that is encoded by the transformer at convergence, as well as the prediction error achieved when the model is given a test prompt corresponding to a new (and possibly nonlinear) prediction task.

1. We note a concurrent work also explores the optimization question we consider here (Ahn et al., 2023); we shall provide a more detailed comparison to this work in Section 2.

- We use this to conclude that transformers trained by gradient flow indeed in-context learn the class of linear models. Moreover, we characterize the robustness of the trained transformer to a variety of distribution shifts. We show that although a number of shifts can be tolerated, shifts in the covariate distribution of the features x_i cannot.
- Motivated by this failure under covariate shift, we consider a generalized setting of in-context learning where the covariate distribution can vary across prompts. We provide global convergence guarantees for LSAs trained by gradient flow in this setting and show that even when trained on a variety of covariate distributions, LSAs still fail under covariate shift.
- We then empirically investigate the behavior of large, nonlinear transformers when trained on linear regression prompts. We find that these more complex models are able to generalize better under covariate shift, especially when trained on prompts with varying covariate distributions.

2. Additional Related Work

The literature on transformers and non-convex optimization in machine learning is vast. In this section, we will focus on those works most closely related to theoretical understanding of in-context learning of function classes.

As mentioned previously, Garg et al. (2022) empirically investigated the ability for transformer architectures to in-context learn a variety of function classes. They showed that when trained on random instances of linear regression, the models’ predictions are very similar to those of ordinary least squares. Additionally, they showed that transformers can in-context learn two-layer ReLU networks and decision trees, showing that by training on differently-structured data, the transformers learn to implement distinct learning algorithms. A number of works further investigated the types of algorithms implemented by transformers trained on in-context examples of linear models (Ahuja et al., 2023; Ahuja and Lopez-Paz, 2023).

Akyürek et al. (2022) and von Oswald et al. (2022) examined the behavior of transformers when trained on random instances of linear regression, as we do in this work. They considered the setting of isotropic Gaussian data with isotropic Gaussian weight vectors, and showed that the trained transformer’s predictions mimic those of a single step of gradient descent. They also provided a construction of transformers which implement this single step of gradient descent. By contrast, we explicitly show that gradient flow provably converges to transformers which learn linear models in-context. Moreover, our analysis holds when the covariates are anisotropic Gaussians, for which a single step of vanilla gradient descent is unable to achieve small prediction error.²

Let us briefly mention a number of other works on understanding in-context learning in transformers and other sequence-based models. Han et al. (2023) suggests that Bayesian inference on prompts can be asymptotically interpreted as kernel regression. Dai et al. (2022)

2. To see this, suppose (x_i, y_i) are i.i.d. with $x \sim \mathbf{N}(0, \Lambda)$ and $y = \langle w, x \rangle$. A single step of gradient descent under the squared loss from a zero initialization yields the predictor $x \mapsto x^\top \left(\frac{1}{n} \sum_{i=1}^n y_i x_i \right) = x^\top \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right) w \approx x^\top \Lambda w$. Clearly, this can differ from $x^\top w$ when $\Lambda \neq I_d$.

interprets ICL as implicit fine-tuning, viewing large language models as meta-optimizers performing gradient-based optimization. Xie et al. (2021) regards ICL as implicit Bayesian inference, with transformers learning a shared latent concept between prompts and test data, and they prove the ICL property when the training distribution is a mixture of HMMs. Similarly, Wang et al. (2023) perceives ICL as a Bayesian selection process, implicitly inferring information pertinent to the designated tasks. Li et al. (2023a) explores the functional resemblance between a single layer of self-attention and gradient descent on a softmax regression problem, offering upper bounds on their difference. Min et al. (2022) notes that the alteration of label parts in prompts does not drastically impair the ICL ability. They contend that ICL is invoked when prompts reveal information about the label space, input distribution, and sequence structure.

Another collection of works have sought to understand transformers from an approximation theoretic perspective. Yun et al. (2019, 2020) established that transformers can universally approximate any sequence-to-sequence function under some assumptions. Investigations by Edelman et al. (2022); Likhoshesterov et al. (2021) indicate that a single-layer self-attention can learn sparse functions of the input sequence, where sample complexity and hidden size are only logarithmic relative to the sequence length. Further studies by Pérez et al. (2019); Deghani et al. (2019); Bhattamishra et al. (2020) indicate that the vanilla transformer and its variants exhibit Turing completeness. Liu et al. (2023) showed that transformers can approximate finite-state automata with few layers. Bai et al. (2023) showed that transformers can implement a variety of statistical machine learning algorithms as well as model selection procedures. Abernethy et al. (2023) showed that a pretrained transformer can be used to define a transformer that segments a prompt into examples and labels and learns to solve a sparse retrieval task. Zhang et al. (2023) interpreted in-context learning via a Bayesian model averaging process.

A handful of recent works have developed provable guarantees for transformers trained with gradient-based optimization. Jelassi et al. (2022) analyzed the dynamics of gradient descent in vision transformers for data with spatial structure. Li et al. (2023c) demonstrated that a single-layer transformer trained by a gradient method could learn a topic model, treating learning semantic structure as detecting co-occurrence between words and theoretically analyzing the two-stage dynamics during the training process.

Finally, we note a concurrent work by Ahn et al. (2023) on the optimization landscape of single layer transformers with linear self-attention layers. They show that there exist global minima of the population objective of the transformer that can achieve small prediction error with anisotropic Gaussian data, and they characterize some critical points of deep linear self-attention networks. In this work, we show that despite nonconvexity, gradient flow with a suitable random initialization converges to a global minimum that achieves small prediction error for anisotropic Gaussian data. We also characterize the prediction error when test prompts come from a new (possibly nonlinear) task, when there is distribution shift, and when transformers are trained on prompts with possibly different covariate distributions across prompts.

3. Preliminaries

Notation We first describe the notation we use in the paper. We write $[n] = \{1, 2, \dots, n\}$. We use \otimes to denote the Kronecker product, and Vec the vectorization operator in column-wise order. For example, $\text{Vec}\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} = (1, 3, 2, 4)^\top$. We write the inner product of two matrices $A, B \in \mathbb{R}^{m \times n}$ as $\langle A, B \rangle = \text{tr}(AB^\top)$. We use 0_n and $0_{m \times n}$ to denote the zero vector and zero matrix of size n and $m \times n$, respectively. For a general matrix A , $A_{k\cdot}$ and $A_{\cdot k}$ denote the k -th row and k -th column, respectively. We denote the matrix operator norm and Frobenius norm as $\|\cdot\|_{op}$ and $\|\cdot\|_F$. We use I_d to denote the d -dimensional identity matrix and sometimes we also use I when the dimension is clear from the context. For a positive semi-definite matrix A , we write $\|x\|_A^2 := x^\top Ax$. Unless otherwise defined, we use lower case letters for scalars and vectors, and use upper case letters for matrices.

3.1 In-context learning

We begin by describing a framework for in-context learning of function classes, as initiated by Garg et al. (2022). In-context learning refers to the behavior of models that operate on sequences, called *prompts*, of input-output pairs $(x_1, y_1, \dots, x_N, y_N, x_{\text{query}})$, where $y_i = h(x_i)$ for some (unknown) function h and examples x_i and query x_{query} . The goal for an in-context learner is to use the prompt to form a prediction $\hat{y}(x_{\text{query}})$ for the query such that $\hat{y}(x_{\text{query}}) \approx h(x_{\text{query}})$.

From this high-level description, one can see that at a surface level, the behavior of in-context learning is no different than that of a standard learning algorithm: the learner takes as input a training dataset and returns predictions on test examples. For instance, one can view ordinary least squares as an ‘in-context learner’ for linear models. However, the rather unique feature of in-context learners is that these learning algorithms can be the solutions to stochastic optimization problems defined over a distribution of prompts. We formalize this notion in the following definition.

Definition 1 (Trained on in-context examples) Let \mathcal{D}_x be a distribution over an input space \mathcal{X} , $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$ a set of functions $\mathcal{X} \rightarrow \mathcal{Y}$, and $\mathcal{D}_{\mathcal{H}}$ a distribution over functions in \mathcal{H} . Let $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. Let $\mathcal{S} = \cup_{n \in \mathbb{N}} \{(x_1, y_1, \dots, x_n, y_n) : x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$ be the set of finite-length sequences of (x, y) pairs and let

$$\mathcal{F}_{\Theta} = \{f_{\theta} : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}, \theta \in \Theta\}$$

be a class of functions parameterized by θ in some set Θ . For $N > 0$, we say that a model $f : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}$ is trained on in-context examples of functions in \mathcal{H} under loss ℓ w.r.t. $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$ if $f = f_{\theta^*}$ where $\theta^* \in \Theta$ satisfies

$$\theta^* \in \underset{\theta \in \Theta}{\text{argmin}} \mathbb{E}_{P=(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})} [\ell(f_{\theta}(P), h(x_{\text{query}}))], \quad (1)$$

where $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$ and $h \sim \mathcal{D}_{\mathcal{H}}$ are independent. We call N the length of the prompts seen during training.

As mentioned above, this definition naturally leads to a method for *learning a learning algorithm from data*: Sample independent prompts by sampling a random function $h \sim \mathcal{D}_{\mathcal{H}}$

and feature vectors $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$, and then minimize the objective function appearing in (1) using stochastic gradient descent or other stochastic optimization algorithms. This procedure returns a model that is learned from in-context examples and can form predictions for test (query) examples given a sequence of training data. This leads to the following natural definition that quantifies how well such a model performs on in-context examples corresponding to a particular hypothesis class.

Definition 2 (In-context learning of a hypothesis class) *Let \mathcal{D}_x be a distribution on an input space \mathcal{X} , $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$ a class of functions $\mathcal{X} \rightarrow \mathcal{Y}$, and $\mathcal{D}_{\mathcal{H}}$ a distribution on functions in \mathcal{H} . Let $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. Let $\mathcal{S} = \cup_{n \in \mathbb{N}} \{(x_1, y_1, \dots, x_n, y_n) : x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$ be the set of finite-length sequences of (x, y) pairs. We say that a model $f : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}$ defined on prompts of the form $P = (x_1, h(x_1), \dots, x_M, h(x_M), x_{\text{query}})$ in-context learns a hypothesis class \mathcal{H} under loss ℓ with respect to $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$ up to error $\eta \in \mathbb{R}$ if there exists a function $M_{\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x}(\varepsilon) : (0, 1) \rightarrow \mathbb{N}$ such that for every $\varepsilon \in (0, 1)$, and for every prompt P of length $M \geq M_{\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x}(\varepsilon)$,*

$$\mathbb{E}_{P=(x_1, h(x_1), \dots, x_M, h(x_M), x_{\text{query}})} \left[\ell \left(f(P), h(x_{\text{query}}) \right) \right] \leq \eta + \varepsilon, \quad (2)$$

where the expectation is over the randomness in $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$ and $h \sim \mathcal{D}_{\mathcal{H}}$.

The additive error term η in Definition 2 above allows for the possibility that the model does not achieve arbitrarily small error. This error could come from using a model which is not complex enough to learn functions in \mathcal{H} or from considering a non-realizable setting where it is not possible to achieve arbitrarily small error.

With these two definitions in hand, we can formulate the following questions: suppose a function class \mathcal{F}_{Θ} is given and $\mathcal{D}_{\mathcal{H}}$ corresponds to random instances of hypotheses in a hypothesis class \mathcal{H} . Can a model from \mathcal{F}_{Θ} that is trained on in-context examples of functions in \mathcal{H} w.r.t. $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$ in-context learn the hypothesis class \mathcal{H} w.r.t. $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$ with small prediction error? Do standard gradient-based optimization algorithms suffice for training the model from in-context examples? How long must the contexts be during training and at test time to achieve small prediction error? In the remaining sections, we shall answer these questions for the case of one-layer transformers with linear self-attention modules when the hypothesis class is linear models, the loss of interest is the squared loss, and the marginals are (possibly anisotropic) Gaussian marginals.

3.2 Linear self-attention networks

Before describing the particular transformer models we analyze in this work, we first recall the definition of the softmax-based single-head self-attention module (Vaswani et al., 2017). Let $E \in \mathbb{R}^{d_e \times d_N}$ be an embedding matrix formed using a prompt $(x_1, y_1, \dots, x_N, y_N, x_{\text{query}})$ of length N . The user has the freedom to determine how this embedding matrix is formed from the prompt. One natural way to form E is to stack $(x_i, y_i)^{\top} \in \mathbb{R}^{d+1}$ as the first N columns of E and to let the final column be $(x_{\text{query}}, 0)^{\top}$; if $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$, we would then have $d_e = d + 1$ and $d_N = N + 1$. Let $W^K, W^Q \in \mathbb{R}^{d_k \times d_e}$ and $W^V \in \mathbb{R}^{d_v \times d_e}$ be the key, query, and value weight matrices, $W^P \in \mathbb{R}^{d_e \times d_v}$ the projection matrix, and $\rho > 0$

a normalization factor. The softmax self-attention module takes as input an embedding matrix E of width d_N and outputs a matrix of the same size,

$$f_{\text{Attn}}(E; W^K, W^Q, W^V, W^P) = E + W^P W^V E \cdot \text{softmax} \left(\frac{(W^K E)^\top W^Q E}{\rho} \right),$$

where softmax is applied column-wise and, given a vector input of v , the i -th entry of $\text{softmax}(v)$ is given by $\exp(v_i) / \sum_s \exp(v_s)$. The $d_N \times d_N$ matrix appearing inside the softmax is referred to as the *self-attention matrix*. Note that f_{Attn} can take as its input a sequence of arbitrary length.

In this work, we consider a simplified version of the single-layer self-attention module, which is more amenable to theoretical analysis and yet is still capable of in-context learning linear models. In particular, we consider a single-layer linear self-attention (LSA) model, which is a modified version of f_{Attn} where we remove the softmax nonlinearity, merge the projection and value matrices into a single matrix $W^{PV} \in \mathbb{R}^{d_e \times d_e}$, and merge the query and key matrices into a single matrix $W^{KQ} \in \mathbb{R}^{d_e \times d_e}$. We concatenate these matrices into $\theta = (W^{KQ}, W^{PV})$ and denote

$$f_{\text{LSA}}(E; \theta) = E + W^{PV} E \cdot \frac{E^\top W^{KQ} E}{\rho}. \quad (3)$$

We note that recent theoretical works on understanding transformers looked at identical models (von Oswald et al., 2022; Li et al., 2023b; Ahn et al., 2023). It is noteworthy that recent empirical work has shown that state-of-the-art trained vision transformers with standard softmax-based attention modules are such that $(W^K)^\top W^Q$ and $W^P W^V$ are nearly multiples of the identity matrix (Trockman and Kolter, 2023), which can be represented under the parameterization we consider.

The user has the flexibility to determine the method for constructing the embedding matrix from a prompt $P = (x_1, y_1, \dots, x_N, y_N, x_{\text{query}})$. In this work, for a prompt of length N , we shall use the following embedding, which stacks $(x_i, y_i)^\top \in \mathbb{R}^{d+1}$ into the first N columns with $(x_{\text{query}}, 0)^\top \in \mathbb{R}^{d+1}$ as the last column:

$$E = E(P) = \begin{pmatrix} x_1 & x_2 & \cdots & x_N & x_{\text{query}} \\ y_1 & y_2 & \cdots & y_N & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (N+1)}. \quad (4)$$

We take the normalization factor ρ to be the width of embedding matrix E minus one, i.e., $\rho = d_N - 1$, since each element in $E \cdot E^\top$ is a inner product of two vectors of length d_N . Under the above token embedding, we take $\rho = N$. We note that there are alternative ways to form the embedding matrix with this data, e.g. by padding all inputs and labels into vectors of equal length and arranging them into a matrix (Akyürek et al., 2022), or by stacking columns that are linear transformations of the concatenation (x_i, y_i) (Garg et al., 2022), although the dynamics of in-context learning will differ under alternative parameterizations.

The network’s prediction for the token x_{query} will be the bottom-right entry of matrix output by f_{LSA} , namely,

$$\hat{y}_{\text{query}} = \hat{y}_{\text{query}}(E; \theta) = [f_{\text{LSA}}(E; \theta)]_{(d+1), (N+1)}.$$

Here and after, we may occasionally suppress dependence on θ and write $\hat{y}_{\text{query}}(E; \theta)$ as \hat{y}_{query} . Since the prediction takes only the right-bottom entry of the token matrix output

by the LSA layer, actually only part of W^{PV} and W^{KQ} affect the prediction. To see how, let us denote

$$W^{PV} = \begin{pmatrix} W_{11}^{PV} & w_{12}^{PV} \\ (w_{21}^{PV})^\top & w_{22}^{PV} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}, \quad W^{KQ} = \begin{pmatrix} W_{11}^{KQ} & w_{12}^{KQ} \\ (w_{21}^{KQ})^\top & w_{22}^{KQ} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)},$$

where $W_{11}^{PV} \in \mathbb{R}^{d \times d}$; $w_{12}^{PV}, w_{21}^{PV} \in \mathbb{R}^d$; $w_{22}^{PV} \in \mathbb{R}$; and $W_{11}^{KQ} \in \mathbb{R}^{d \times d}$; $w_{12}^{KQ}, w_{21}^{KQ} \in \mathbb{R}^d$; $w_{22}^{KQ} \in \mathbb{R}$. Then, the prediction \hat{y}_{query} is

$$\hat{y}_{\text{query}} = \begin{pmatrix} (w_{21}^{PV})^\top & w_{22}^{PV} \end{pmatrix} \cdot \left(\frac{EE^\top}{N} \right) \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} x_{\text{query}}, \quad (6)$$

since only the last row of W^{PV} and the first d columns of W^{KQ} affects the prediction, which means we can simply take all other entries zero in the following sections.

3.3 Training procedure

In this work, we will consider the task of in-context learning linear predictors. We will assume training prompts are sampled as follows. Let Λ be a positive definite covariance matrix. Each training prompt, indexed by $\tau \in \mathbb{N}$, takes the form

$$P_\tau = (x_{\tau,1}, h_\tau(x_{\tau,1}), \dots, x_{\tau,N}, h_\tau(x_{\tau,N}), x_{\tau,\text{query}}),$$

where task weights $w_\tau \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$, inputs $x_{\tau,i}, x_{\tau,\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Lambda)$, and labels $h_\tau(x) = \langle w_\tau, x \rangle$.

Each prompt corresponds to an embedding matrix E_τ , formed using the transformation (4):

$$E_\tau := \begin{pmatrix} x_{\tau,1} & x_{\tau,2} & \cdots & x_{\tau,N} & x_{\tau,\text{query}} \\ \langle w_\tau, x_{\tau,1} \rangle & \langle w_\tau, x_{\tau,2} \rangle & \cdots & \langle w_\tau, x_{\tau,N} \rangle & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (N+1)}.$$

We denote the prediction of the LSA model on the query label in the task τ as $\hat{y}_{\tau,\text{query}}$, which is the bottom-right element of $f_{\text{LSA}}(E_\tau)$, where f_{LSA} is the linear self-attention model defined in (3). The empirical risk over B independent prompts is defined as

$$\hat{L}(\theta) = \frac{1}{2B} \sum_{\tau=1}^B \left(\hat{y}_{\tau,\text{query}} - \langle w_\tau, x_{\tau,\text{query}} \rangle \right)^2. \quad (7)$$

We shall consider the behavior of gradient flow-trained networks over the population loss induced by the limit of infinite training tasks/prompts $B \rightarrow \infty$:

$$L(\theta) = \lim_{B \rightarrow \infty} \hat{L}(\theta) = \frac{1}{2} \mathbb{E}_{w_\tau, x_{\tau,1}, \dots, x_{\tau,N}, x_{\tau,\text{query}}} \left[(\hat{y}_{\tau,\text{query}} - \langle w_\tau, x_{\tau,\text{query}} \rangle)^2 \right] \quad (8)$$

Above, the expectation is taken w.r.t. the covariates $\{x_{\tau,i}\}_{i=1}^N \cup \{x_{\text{query}}\}$ in the prompt and the weight vector w_τ , i.e. over $x_{\tau,i}, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Lambda)$ and $w_\tau \sim \mathcal{N}(0, I_d)$. Gradient

flow captures the behavior of gradient descent with infinitesimal step size and has dynamics given by the following differential equation:

$$\frac{d}{dt}\theta = -\nabla L(\theta). \quad (9)$$

We will consider gradient flow with an initialization that satisfies the following.

Assumption 3 (Initialization) *Let $\sigma > 0$ be a parameter, and let $\Theta \in \mathbb{R}^{d \times d}$ be any matrix satisfying $\|\Theta\Theta^\top\|_F = 1$ and $\Theta\Lambda \neq 0_{d \times d}$. We assume*

$$W^{PV}(0) = \sigma \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}, \quad W^{KQ}(0) = \sigma \begin{pmatrix} \Theta\Theta^\top & 0_d \\ 0_d^\top & 0 \end{pmatrix}. \quad (10)$$

This initialization is satisfied for a particular class of random initialization schemes: if M has i.i.d. entries from a continuous distribution, then by setting $\Theta\Theta^\top = MM^\top / \|MM^\top\|_F$, the assumption is satisfied almost surely. The reason we use this particular initialization scheme will be made more clear in Section 5 when we describe the proof, but at a high-level this is due to the fact that the predictions (6) can be viewed as the output of a two-layer linear network, and initializations satisfying Assumption 3 allow for the layers to be ‘balanced’ throughout the gradient flow trajectory. Random initializations that induce this balancedness condition have been utilized in a number of theoretical works on deep linear networks (Du et al., 2018; Arora et al., 2018, 2019; Azulay et al., 2021). We leave the question of convergence under alternative random initialization schemes for future work.

4. Main results

In this section, we present the main results of this paper. First, in Section 4.1, we prove the gradient flow on the population loss will converge to a specific global optimum. We characterize the prediction error of the trained transformer at this global minimum when given a prompt from a new prediction task. Our characterization allows for the possibility that this new prompt comes from a nonlinear prediction task. We then instantiate our results for well-specified linear regression prompts and characterize the number of samples needed to achieve small prediction error, showing that transformers can in-context learn linear models when trained on in-context examples of linear models.

Next, in Section 4.2, we analyze the behavior of the trained transformer under a variety of distribution shifts. We show the transformer is robust to a number of distribution shifts, including task shift (when the labels in the prompt are not deterministic linear functions of their input) and query shift (when the query example x_{query} has a possibly different distribution than the test prompt). On the other hand, we show that the transformer suffers from covariate distribution shifts, i.e. when the training prompt covariate distribution differs from the test prompt covariate distribution.

Finally, motivated by the failure of the trained transformer under covariate distribution shift, we consider in Section 4.3 the setting of training on in-context examples with varying covariate distributions across prompts. We prove that transformers with a single linear self-attention layer trained by gradient flow converge to a global minimum of the population objective, but that the trained transformer still fails to perform well on new prompts. We complement our proof in the linear self-attention case with experiments on large, nonlinear transformer architectures which we show are more robust under covariate shifts.

4.1 Convergence of gradient flow and prediction error for new tasks

First, we prove that under suitable initialization, gradient flow will converge to a global optimum.

Theorem 4 (Convergence and limits) *Consider gradient flow of a linear self-attention network f_{LSA} defined in (3) over the population loss (8). Suppose the initialization satisfies Assumption 3 with initialization scale $\sigma > 0$ satisfying $\sigma^2 \|\Gamma\|_{\text{op}} \sqrt{d} < 2$ where we have defined*

$$\Gamma := \left(1 + \frac{1}{N}\right) \Lambda + \frac{1}{N} \text{tr}(\Lambda) I_d \in \mathbb{R}^{d \times d}.$$

Then gradient flow converges to a global minimum of the population loss (8). Moreover, W^{PV} and W^{KQ} converge to W_^{PV} and W_*^{KQ} respectively, where*

$$W_*^{KQ} = [\text{tr}(\Gamma^{-2})]^{-\frac{1}{4}} \begin{pmatrix} \Gamma^{-1} & 0_d \\ 0_d^\top & 0 \end{pmatrix}, \quad W_*^{PV} = [\text{tr}(\Gamma^{-2})]^{\frac{1}{4}} \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}. \quad (11)$$

The full proof of this theorem appears in Appendix A. We note that if we restrict our setting to $\Lambda = I_d$, then the limiting solution described found by gradient flow is quite similar to the construction of von Oswald et al. (2022). Since the prediction of the transformer is the same if we multiply W^{PV} by a constant $c \neq 0$ and simultaneously multiply W^{KQ} by c^{-1} , the only difference (up to scaling) is that the top-left entry of their W^{KQ} matrix is I_d rather than the $(1 + (d+1)/N)^{-1} I_d$ that we find for the case $\Lambda = I_d$.

Next, we would like to characterize the prediction error of the trained network described above when the network is given a new prompt. Let us consider a prompt of the form $(x_1, \langle w, x_1 \rangle, \dots, x_M, \langle w, x_M \rangle, x_{\text{query}})$ where $w \in \mathbb{R}^d$ and $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Lambda)$. A simple calculation shows that the prediction \hat{y}_{query} at the global optimum with parameters W_*^{KQ} and W_*^{PV} is given by

$$\begin{aligned} \hat{y}_{\text{query}} &= \begin{pmatrix} 0_d^\top & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{M} \sum_{i=1}^M x_i x_i^\top + \frac{1}{M} x_{\text{query}} x_{\text{query}}^\top & \frac{1}{M} \sum_{i=1}^M x_i x_i^\top w \\ \frac{1}{M} \sum_{i=1}^M w^\top x_i x_i^\top & \frac{1}{M} \sum_{i=1}^M w^\top x_i x_i^\top w \end{pmatrix} \begin{pmatrix} \Gamma^{-1} & 0_d \\ 0_d^\top & 0 \end{pmatrix} \begin{pmatrix} x_{\text{query}} \\ 0 \end{pmatrix} \\ &= x_{\text{query}}^\top \Gamma^{-1} \left(\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right) w. \end{aligned} \quad (12)$$

When the length of prompts seen during training N is large, $\Gamma^{-1} \approx \Lambda^{-1}$, and when the test prompt length M is large, $\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \approx \Lambda$, so that $\hat{y}_{\text{query}} \approx x_{\text{query}}^\top w$. Thus, for sufficiently large prompt lengths, *the trained transformer indeed in-context learns the class of linear predictors.*

In fact, we can generalize the above calculation for test prompts which could take a significantly different form than the training prompts. Consider prompts that are of the form $(x_1, y_1, \dots, x_n, y_n, x_{\text{query}})$ where, for some joint distribution \mathcal{D} over (x, y) pairs with marginal

distribution $x \sim \mathbf{N}(0, \Lambda)$, we have $(x_i, y_i) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$ and $x_{\text{query}} \sim \mathbf{N}(0, \Lambda)$ independently. Note that this allows for a label y_i to be a nonlinear function of the input x_i . The prediction of the trained transformer for this prompt is then

$$\begin{aligned} \hat{y}_{\text{query}} &= \begin{pmatrix} 0_d^\top & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{M} \sum_{i=1}^M x_i x_i^\top + \frac{1}{M} x_{\text{query}} x_{\text{query}}^\top & \frac{1}{M} \sum_{i=1}^M x_i y_i \\ \frac{1}{M} \sum_{i=1}^M x_i^\top y_i & \frac{1}{M} \sum_{i=1}^M y_i^2 \end{pmatrix} \begin{pmatrix} \Gamma^{-1} & 0_d \\ 0_d^\top & 0 \end{pmatrix} \begin{pmatrix} x_{\text{query}} \\ 0 \end{pmatrix} \\ &= x_{\text{query}}^\top \Gamma^{-1} \begin{pmatrix} \frac{1}{M} \sum_{i=1}^M y_i x_i \end{pmatrix}. \end{aligned} \quad (13)$$

Just as before, when N is large we have $\Gamma^{-1} \approx \Lambda^{-1}$, and so when M is large as well this implies

$$\hat{y}_{\text{query}} \approx x_{\text{query}}^\top \Lambda^{-1} \mathbb{E}_{(x,y) \sim \mathcal{D}}[yx] = x_{\text{query}}^\top \left(\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}_{(x,y) \sim \mathcal{D}}[(y - \langle w, x \rangle)^2] \right). \quad (14)$$

This suggests that trained transformers in-context learn the *best linear predictor* over a distribution when the test prompt consists of i.i.d. samples from a joint distribution over feature-response pairs. In the following theorem, we formalize the above and characterize the prediction error when prompts take this form.

Theorem 5 *Let \mathcal{D} be a distribution over $(x, y) \in \mathbb{R}^d \times \mathbb{R}$, whose marginal distribution on x is $\mathcal{D}_x = \mathbf{N}(0, \Lambda)$. Assume $\mathbb{E}_{\mathcal{D}}[y], \mathbb{E}_{\mathcal{D}}[xy], \mathbb{E}_{\mathcal{D}}[y^2 x x^\top]$ exist and are finite. Assume the test prompt is of the form $P = (x_1, y_1, \dots, x_M, y_M, x_{\text{query}})$, where $(x_i, y_i), (x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$. Let f_{LSA}^* be the LSA model with parameters W_*^{PV} and W_*^{KQ} in (11), and \hat{y}_{query} is the prediction for x_{query} given the prompt. If we define*

$$a := \Lambda^{-1} \mathbb{E}_{(x,y) \sim \mathcal{D}}[xy], \quad \Sigma := \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[(xy - \mathbb{E}(xy))(xy - \mathbb{E}(xy))^\top \right], \quad (15)$$

then, for $\Gamma = \Lambda + \frac{1}{N} \Lambda + \frac{1}{N} \operatorname{tr}(\Lambda) I_d$. we have,

$$\begin{aligned} \mathbb{E}(\hat{y}_{\text{query}} - y_{\text{query}})^2 &= \underbrace{\min_{w \in \mathbb{R}^d} \mathbb{E}(\langle w, x_{\text{query}} \rangle - y_{\text{query}})^2}_{\text{Error of best linear predictor}} \\ &\quad + \frac{1}{M} \operatorname{tr}[\Sigma \Gamma^{-2} \Lambda] + \frac{1}{N^2} \left[\|a\|_{\Gamma^{-2} \Lambda^3}^2 + 2 \operatorname{tr}(\Lambda) \|a\|_{\Gamma^{-2} \Lambda^2}^2 + \operatorname{tr}(\Lambda)^2 \|a\|_{\Gamma^{-2} \Lambda}^2 \right], \end{aligned} \quad (16)$$

where the expectation is over $(x_i, y_i), (x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$.

The full proof is deferred to Appendix B. Let us now make a few remarks on the above theorem before considering particular instances of \mathcal{D} where we may provide more explicit bounds on the prediction error.

First, this theorem shows that, provided the length of prompts seen during training (N) and the length of the test prompt (M) is large enough, a transformer trained by

gradient flow from in-context examples achieves prediction error competitive with the best linear model. Next, our bound shows that the length of prompts seen during training and the length of prompts seen at test-time have different effects on the expected prediction error: ignoring dimension and covariance-dependent factors, the prediction error is at most $O(1/M + 1/N^2)$, decreasing more rapidly as a function of the training prompt length N compared to the test prompt length M .

Let us now consider when \mathcal{D} corresponds to noiseless linear models, so that for some $w \in \mathbb{R}^d$, we have $(x, y) = (x, \langle w, x \rangle)$, in which case the prediction of the trained transformer is given by (12). Moreover, a simple calculation shows that the Σ from Theorem 5 takes the form $\Sigma = \|w\|_\Lambda^2 \Lambda + \Lambda w w^\top \Lambda$. Hence Theorem 5 implies the prediction error for the prompt $P = (x_1, \langle w, x_1 \rangle, \dots, x_M, \langle w, x_M \rangle, x_{\text{query}})$ is

$$\begin{aligned} & \mathbb{E}_{x_1, \dots, x_M, x_{\text{query}}} (\hat{y}_{\text{query}} - \langle w, x_{\text{query}} \rangle)^2 \\ &= \frac{1}{M} \left\{ \|w\|_{\Gamma^{-2}\Lambda^3}^2 + \text{tr}(\Gamma^{-2}\Lambda^2) \|w\|_\Lambda^2 \right\} \\ & \quad + \frac{1}{N^2} \left\{ \|w\|_{\Gamma^{-2}\Lambda^3}^2 + 2 \|w\|_{\Gamma^{-2}\Lambda^2}^2 \text{tr}(\Lambda) + \|w\|_{\Gamma^{-2}\Lambda}^2 \text{tr}(\Lambda)^2 \right\} \\ & \leq \frac{d+1}{M} \|w\|_\Lambda^2 + \frac{1}{N^2} \left[\|w\|_\Lambda^2 + 2 \|w\|_2^2 \text{tr}(\Lambda) + \|w\|_{\Lambda^{-1}}^2 \text{tr}(\Lambda)^2 \right]. \end{aligned}$$

The inequality above uses that $\Gamma \succ \Lambda$. Finally, if we assume that $w \sim \mathbf{N}(0, I_d)$ and denote κ as the condition number of Λ , then by taking expectations over w we get the following:

$$\begin{aligned} & \mathbb{E}_{x_1, \dots, x_M, x_{\text{query}}, w} (\hat{y}_{\text{query}} - \langle w, x_{\text{query}} \rangle)^2 \\ & \leq \frac{(d+1) \text{tr}(\Lambda)}{M} + \frac{1}{N^2} [\text{tr}(\Lambda) + 2d \text{tr}(\Lambda) + \text{tr}(\Lambda^{-1}) \text{tr}(\Lambda)^2] \\ & \leq \frac{(d+1) \text{tr}(\Lambda)}{M} + \frac{(1+2d+d^2\kappa) \text{tr}(\Lambda)}{N^2}, \end{aligned}$$

From the upper bound above, we can see the rate w.r.t M and N are still at most $O(1/M)$ and $O(1/N^2)$ respectively. Moreover, the generalization risk also scales with dimension d , $\text{tr}(\Lambda)$ and the condition number κ . This suggests that for in-context examples involving covariates of greater variance, or a more ill-conditioned covariance matrix, the generalization risk will be higher for the same lengths of training and testing prompts. Putting the above together with Theorem 5, Definition 1 and Definition 2, we get the following corollary.

Corollary 6 *The transformer f_{LSA} trained on length- N prompts of in-context examples of functions in $\{x \mapsto \langle w, x \rangle\}$ w.r.t. $w \sim \mathbf{N}(0, I_d)$ and $\mathcal{D}_x = \mathbf{N}(0, \Lambda)$ by gradient flow on the population loss (8) for initializations satisfying Assumption 3 converges to the model $f_{\text{LSA}}(\cdot; W_*^{KQ}, W_*^{PV})$. This model takes a prompt $P = (x_1, y_1, \dots, x_M, y_M, x_{\text{query}})$ and returns a prediction \hat{y}_{query} for x_{query} given by*

$$\hat{y}_{\text{query}} = [f_{\text{LSA}}(P; W_*^{KQ}, W_*^{PV})]_{d+1, M+1} = x_{\text{query}}^\top \left(\Lambda + \frac{1}{N} \Lambda + \frac{\text{tr}(\Lambda)}{N} I_d \right)^{-1} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i \right).$$

This model in-context learns the class of linear models $\{x \mapsto \langle w, x \rangle\}$ with respect to $w \sim \mathbf{N}(0, I_d)$ and $\mathcal{D}_x = \mathbf{N}(0, \Lambda)$ up to error $\eta := (1+2d+d^2\kappa) \text{tr}(\Lambda)/N^2$ (where κ is the condition

number of Λ): provided $M \geq (d+1) \text{tr}(\Lambda)\varepsilon^{-1}$, the model achieves prediction error at most $\eta + \varepsilon$.

It is worth emphasizing that the transformer $f_{\text{LSA}}(\cdot; W_*^{KQ}, W_*^{PV})$ only learns the function class up to error $\eta = O(1/N^2)$ in the sense of Definition 2. In particular, training on finite-length prompts leads to prediction error bounded away from zero.

4.2 Behavior of trained transformer under distribution shifts

Using the identity (13), it is straightforward to characterize the behavior of the trained transformer under a variety of distribution shifts. In this section, we shall examine a number of shifts that were first explored empirically for transformer architectures by Garg et al. (2022). Although their experiments were for transformers trained by gradient descent, we find that (in the case of linear models) many of the behaviors of the trained transformers under distribution shift are identical to those predicted by our theoretical characterizations of the performance of transformers with a single linear self-attention layer trained by gradient flow on the population.

Following Garg et al. (2022), for prompts of the form $(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})$, let us assume for training prompts that $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x^{\text{train}}$ and $h \sim \mathcal{D}_h^{\text{train}}$, while for test prompts $x_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x^{\text{test}}$, $x_{\text{query}} \sim \mathcal{D}_{\text{query}}^{\text{test}}$, and $h \sim \mathcal{D}_h^{\text{test}}$. We will consider the following distinct categories of shifts:

- Task shifts: $\mathcal{D}_h^{\text{train}} \neq \mathcal{D}_h^{\text{test}}$.
- Query shifts: $\mathcal{D}_{\text{query}}^{\text{test}} \neq \mathcal{D}_x^{\text{test}}$.
- Covariate shifts: $\mathcal{D}_x^{\text{train}} \neq \mathcal{D}_x^{\text{test}}$.

In the following, we shall fix $\mathcal{D}_x^{\text{train}} = \mathbf{N}(0, \Lambda)$ and vary the other distributions. Recall from (13) that the prediction for a test prompt $(x_1, y_1, \dots, x_N, y_N, x_{\text{query}})$ is given by (for N large),

$$\hat{y}_{\text{query}} = x_{\text{query}}^\top \Gamma^{-1} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i \right) \approx x_{\text{query}}^\top \Lambda^{-1} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i \right). \quad (17)$$

Task shifts. These shifts are tolerated easily by the trained transformer. As Theorem 5 shows, the trained transformer is competitive with the best linear model provided the prompt length during training and at test time is large enough. In particular, even if the prompt is such that the labels y_i are not given by $\langle w, x_i \rangle$ for some $w \sim \mathbf{N}(0, I_d)$, the trained transformer will compute a prediction which has error competitive with the best linear model that fits the test prompt.

For example, consider a prompt corresponding to a noisy linear model, so that the prompt consists of a sequence of (x_i, y_i) pairs where $y_i = \langle w, x_i \rangle + \varepsilon_i$ for some arbitrary vector $w \in \mathbb{R}^d$ and independent sub-Gaussian noise ε_i . Then from (17), the prediction of the transformer on query examples is

$$\hat{y}_{\text{query}} \approx x_{\text{query}}^\top \Lambda^{-1} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i \right) = x_{\text{query}}^\top \Lambda^{-1} \left(\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right) w + x_{\text{query}}^\top \Lambda^{-1} \left(\frac{1}{M} \sum_{i=1}^M \varepsilon_i x_i \right).$$

Since ε_i is mean zero and independent of x_i , this is approximately $x_{\text{query}}^\top w$ when M is large. And note that this calculation holds for an *arbitrary* vector w , not just those which are sampled from an isotropic Gaussian or those with a particular norm. This behavior coincides with that of the trained transformers observed by Garg et al. (2022).

Query shifts. Continuing from (17), since $y_i = \langle w, x_i \rangle$,

$$\hat{y}_{\text{query}} \approx x_{\text{query}}^\top \Lambda^{-1} \left(\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right) w.$$

From this we see that whether query shifts can be tolerated hinges upon the distribution of the x_i 's. Since $\mathcal{D}_x^{\text{train}} = \mathcal{D}_x^{\text{test}}$, if M is large then

$$\hat{y}_{\text{query}} \approx x_{\text{query}}^\top \Lambda^{-1} \Lambda w = x_{\text{query}}^\top w. \tag{18}$$

Thus, very general shifts in the query distribution can be tolerated. On the other hand, very different behavior can be expected if M is not large and the query example depends on the training data. For example, if the query example is orthogonal to the subspace spanned by the x_i 's, the prediction will be zero, as was observed with transformer architectures by Garg et al. (2022).

Covariate shifts. In contrast to task and query shifts, covariate shifts cannot be fully tolerated in the transformer. This can be easily seen due to the identity (13): when $\mathcal{D}_x^{\text{train}} \neq \mathcal{D}_x^{\text{test}}$, then the approximation in (18) does not hold as $\frac{1}{M} \sum_{i=1}^M x_i x_i^\top$ will not cancel Γ^{-1} when M and N are large. For instance, if we consider test prompts where the covariates are scaled by a constant $c \neq 1$, then

$$\hat{y}_{\text{query}} \approx x_{\text{query}}^\top \Lambda^{-1} \left(\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right) w \approx x_{\text{query}}^\top \Lambda^{-1} c^2 \Lambda w = c^2 x_{\text{query}}^\top w \neq x_{\text{query}}^\top w.$$

This failure mode of the trained transformer with linear self-attention was also observed in the trained transformer architectures by Garg et al. (2022). This suggests that although the predictions of the transformer may look similar to those of ordinary least squares in some settings, the algorithm implemented by the transformer is not the same since ordinary least squares is robust to scaling of the features by a constant.

It may seem surprising that a transformer trained on linear regression tasks fails in settings where ordinary least squares performs well. However, both the linear self-attention transformer we consider and the transformers considered by Garg et al. (2022) were trained on instances of linear regression when the covariate distribution \mathcal{D}_x over the features was fixed across instances. This leads to the natural question of what happens if the transformers instead are trained on prompts where the covariate distribution varies across instances, which we explore in the following section.

4.3 Transformers trained on prompts with random covariate distributions

In this section, we will consider a variant of training on in-context examples (in the sense of Definition 1) where the distribution \mathcal{D}_x is itself sampled randomly from a distribution, and training prompts are of the form $(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})$ where $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$ and $h \sim \mathcal{D}_H$. More formally, we can generalize Definition 1 as follows.

Definition 7 (In-context training with random covariate distributions) Let Δ be a distribution over distributions \mathcal{D}_x defined on an input space \mathcal{X} , $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$ a set of functions $\mathcal{X} \rightarrow \mathcal{Y}$, and $\mathcal{D}_{\mathcal{H}}$ a distribution over functions in \mathcal{H} . Let $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a loss function. Let $\mathcal{S} = \cup_{n \in \mathbb{N}} \{(x_1, y_1, \dots, x_n, y_n) : x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$ be the set of finite-length sequences of (x, y) pairs and let

$$\mathcal{F}_{\Theta} = \{f_{\theta} : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}, \theta \in \Theta\}$$

be a class of functions parameterized by some set Θ . We say that a model $f : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}$ is trained on in-context examples of functions in \mathcal{H} under loss ℓ w.r.t. $\mathcal{D}_{\mathcal{H}}$ and distribution over covariate distributions Δ if $f = f_{\theta^*}$ where $\theta^* \in \Theta$ satisfies

$$\theta^* \in \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}_{P=(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})} [\ell(f_{\theta}(P), h(x_{\text{query}}))], \quad (19)$$

where $\mathcal{D}_x \sim \Delta$, $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$ and $h \sim \mathcal{D}_{\mathcal{H}}$.

We recover the previous definition of training on in-context examples by taking Δ to be concentrated on a singleton, $\operatorname{supp}(\Delta) = \{\mathcal{D}_x\}$. The natural question is then, if a model f is trained on in-context examples from a function class \mathcal{H} w.r.t. $\mathcal{D}_{\mathcal{H}}$ and a distribution Δ over covariate distributions, and if one then samples some covariate distribution $\mathcal{D}_x \sim \Delta$, does f in-context learn \mathcal{H} w.r.t. $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$ for that \mathcal{D}_x (cf. Definition 2)? Since \mathcal{D}_x is random, we can hope that this may hold in expectation or with high probability over the sampling of the covariate distribution. In the remainder of this section, we will explore this question for transformers with a linear self-attention layer trained by gradient flow on the population loss.

We shall again consider the case where the covariates have Gaussian marginals, $x_i \sim \mathbf{N}(0, \Lambda)$, but we shall now assume that within each prompt we first sample a random covariance matrix Λ . For simplicity, we will restrict our attention to the case where Λ is diagonal. More formally, we shall assume training prompts are sampled as follows. For each independent task indexed by $\tau \in [B]$, we first sample $w_{\tau} \sim \mathbf{N}(0, I_d)$. Then, for each task τ and coordinate $i \in [d]$, we sample $\lambda_{\tau, i}$ independently such that the distribution of each $\lambda_{\tau, i}$ is fixed and has finite third moments and is strictly positive almost surely. We then form a diagonal matrix

$$\Lambda_{\tau} = \operatorname{diag}(\lambda_{\tau, 1}, \dots, \lambda_{\tau, d}).$$

Thus the diagonal entries of Λ_{τ} are independent but could have different distributions, and Λ_{τ} is identically distributed for $\tau = 1, \dots, B$. Then, conditional on Λ_{τ} , we sample independent and identically distributed $x_{\tau, 1}, \dots, x_{\tau, N}, x_{\tau, \text{query}} \sim \mathbf{N}(0, \Lambda_{\tau})$. A training prompt is then given by $P_{\tau} = (x_{\tau, 1}, \langle w_{\tau}, x_{\tau, 1} \rangle, \dots, x_{\tau, N}, \langle w_{\tau}, x_{\tau, N} \rangle, x_{\tau, \text{query}})$. Notice that here, $x_{\tau, i}, x_{\tau, \text{query}}$ are conditionally independent given the covariance matrix Λ_{τ} , but not independent in general. We consider the same token embedding matrix as (4) and linear self-attention network, which forms the prediction $\hat{y}_{\text{query}, \tau}$ as in (6). The empirical risk is the same as before (see (7)), and as in (8), we then take $B \rightarrow \infty$ and consider the gradient flow on the population loss. The population loss now includes an expectation over the distribution of the covariance matrices in addition to the task weight w_{τ} and covariate distributions, and is given by

$$L(\theta) = \frac{1}{2} \mathbb{E}_{w_{\tau}, \Lambda_{\tau}, x_{\tau, 1}, \dots, x_{\tau, N}, x_{\tau, \text{query}}} [(\hat{y}_{\tau, \text{query}} - \langle w_{\tau}, x_{\tau, \text{query}} \rangle)^2]. \quad (20)$$

In the main result for this section, we show that gradient flow with a suitable initialization converges to a global minimum, and we characterize the limiting solution. The proof will be deferred to Appendix C.

Theorem 8 (Global convergence with random covariance) *Consider gradient flow of the linear self-attention network f_{LSA} defined in (3) over the population loss (20), where Λ_τ are diagonal with independent diagonal entries which are strictly positive a.s. and have finite third moments. Suppose the initialization satisfies Assumption 3, $\|\mathbb{E}\Lambda_\tau\Theta\|_F \neq 0$, with initialization scale $\sigma > 0$ satisfying*

$$\sigma^2 < \frac{2 \|\mathbb{E}\Lambda_\tau\Theta\|_F^2}{\sqrt{d} \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right]}. \quad (21)$$

Then gradient flow converges to a global minimum of the population loss (20). Moreover, W^{PV} and W^{KQ} converge to W_*^{PV} and W_*^{KQ} respectively, where

$$\begin{aligned} W_*^{KQ} &= \left\| \left[\mathbb{E}\Gamma_\tau \Lambda_\tau^2 \right]^{-1} \mathbb{E} \left[\Lambda_\tau^2 \right] \right\|_F^{-\frac{1}{2}} \cdot \begin{pmatrix} \left[\mathbb{E}\Gamma_\tau \Lambda_\tau^2 \right]^{-1} \left[\mathbb{E}\Lambda_\tau^2 \right] & 0_d \\ 0_d^\top & 0 \end{pmatrix}, \\ W_*^{PV} &= \left\| \left[\mathbb{E}\Gamma_\tau \Lambda_\tau^2 \right]^{-1} \mathbb{E} \left[\Lambda_\tau^2 \right] \right\|_F^{\frac{1}{2}} \cdot \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}, \end{aligned} \quad (22)$$

where $\Gamma_\tau = \frac{N+1}{N}\Lambda_\tau + \frac{1}{N}\text{tr}(\Lambda_\tau)I_d \in \mathbb{R}^{d \times d}$ and the expectations above are over the distribution of Λ_τ .

From this result, we can see why the trained transformer fails in the random covariance case. Suppose we have a new prompt corresponding to a weight matrix $w \in \mathbb{R}^d$ and covariance matrix Λ_{new} , sampled from the same distribution as the covariance matrices for training prompts, so that conditionally on Λ_{new} we have $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathbf{N}(0, \Lambda_{\text{new}})$. The ground-truth labels are given by $y_i = \langle w, x_i \rangle, i \in [M]$ and $y_{\text{query}} = \langle w, x_{\text{query}} \rangle$. At convergence, the prediction by the trained transformer on the new task will be

$$\begin{aligned} &\hat{y}_{\text{query}} \\ &= (0_d^\top \quad 1) \begin{pmatrix} \frac{1}{M} \sum_{i=1}^M x_i x_i^\top + \frac{1}{M} x_{\text{query}} x_{\text{query}}^\top & \frac{1}{M} \sum_{i=1}^M x_i y_i \\ \frac{1}{M} \sum_{i=1}^M x_i^\top y_i & \frac{1}{M} \sum_{i=1}^M y_i^2 \end{pmatrix} \begin{pmatrix} \left[\mathbb{E}\Gamma_\tau \Lambda_\tau^2 \right]^{-1} \mathbb{E}\Lambda_\tau^2 & 0_d \\ 0_d^\top & 0 \end{pmatrix} \begin{pmatrix} x_{\text{query}} \\ 0 \end{pmatrix} \\ &= x_{\text{query}}^\top \cdot \left[\mathbb{E}\Lambda_\tau^2 \right] \left[\mathbb{E}\Gamma_\tau \Lambda_\tau^2 \right]^{-1} \cdot \left[\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right] w \\ &\rightarrow x_{\text{query}}^\top \cdot \left[\mathbb{E}\Lambda_\tau^2 \right] \left[\mathbb{E}\Gamma_\tau \Lambda_\tau^2 \right]^{-1} \cdot \Lambda_{\text{new}} w \quad \text{almost surely when } M \rightarrow \infty. \end{aligned} \quad (23)$$

The last line comes from the strong law of large numbers. Thus, in order for the prediction on the query example to be close to the ground-truth $x_{\text{query}}^\top w$, we need $\left[\mathbb{E}\Lambda_\tau^2 \right] \left[\mathbb{E}\Gamma_\tau \Lambda_\tau^2 \right]^{-1} \cdot \Lambda_{\text{new}}$

to be close to the identity. When $\Lambda_\tau \equiv \Lambda_{\text{new}}$ is deterministic, this indeed is the case as we know from Theorem 5. However, this clearly does not hold in general when Λ_τ is random.

To make things concrete, let us assume for simplicity that $M, N \rightarrow \infty$ so that $\Gamma_\tau \rightarrow \Lambda_\tau$ and the identity (23) holds (conditionally on Λ_{new}). Then, taking expectation over Λ_{new} in (23), we obtain

$$\mathbb{E}[\widehat{y}_{\text{query}} | x_{\text{query}}, w] \rightarrow x_{\text{query}}^\top \cdot [\mathbb{E}\Lambda_\tau^2] [\mathbb{E}\Lambda_\tau^3]^{-1} \cdot [\mathbb{E}\Lambda_\tau] w.$$

If we consider the case $\lambda_{\tau,i} \stackrel{\text{i.i.d.}}{\sim} \text{Exponential}(1)$, so that $\mathbb{E}[\Lambda_\tau] = I_d$, $\mathbb{E}[\Lambda_\tau^2] = 2I_d$, and $\mathbb{E}[\Lambda_\tau^3] = 6I_d$, we get

$$\mathbb{E}\widehat{y}_{\text{query}} \rightarrow \frac{1}{3}\langle w, x_{\text{query}} \rangle.$$

This shows that for transformers with a single linear self-attention layer, training on in-context examples with random covariate distributions does not allow for in-context learning of a hypothesis class with varying covariate distributions.

Experiments with large, nonlinear transformers. We have shown that even when trained on prompts with random covariance matrices, transformers with a single linear self-attention layer fail to in-context learn linear models with random covariance matrices. We now investigate the behavior of more complex transformer architectures that are trained on in-context examples of linear models, both in the fixed-covariance case and in the random-covariance case.

We examine the performance of transformers with a GPT2 architecture (Radford et al., 2019) that are trained on linear regression tasks with mean-zero Gaussian features with either a fixed covariance matrix or random covariance matrices. For the fixed covariance case, the covariance matrix is fixed to the identity matrix across prompts. For the random covariance case, covariates are drawn from $x \sim \mathbf{N}(0, c\Lambda)$ where Λ is diagonal with $\lambda_i \stackrel{\text{i.i.d.}}{\sim} \text{Exponential}(1)$ and $c > 0$ is a scaling factor. We set $c = 1$ during training and vary this value at test time. The transformer is trained using the procedure of Garg et al. (2022) (see Appendix E for more details). We consider linear models in $d = 20$ dimensions and we train on prompt lengths of $N = 40, 70, 100$ with either fixed or random covariance matrices. The performance of these trained models, when tested on new data with fixed covariance or random covariance matrices ($c = 1, 4, 9$), is represented in six curves in Figure 1. Using the calculation (23), we can compare the prediction error for the linear self-attention networks in the $M \rightarrow \infty, N \rightarrow \infty$ limit (the black dash line) to those of GPT2 architectures. We additionally compare these models to the ordinary least-squares solution which is optimal for this task.

From the figure, we can see that the GPT2 model trained on fixed covariance succeeds in the random covariance setting if the variance is not too large, which shows that the larger nonlinear model is able to generalize better than the model with a single linear self-attention layer. However, when the variance is large ($c = 4, 9$ for the bottom two figures), the GPT2 model trained with fixed covariance is unsuccessful. When trained on random covariance, the model performs better for test prompts from higher-variance random covariance matrices, but still fails to match least squares when the scaling is largest ($c = 9$).

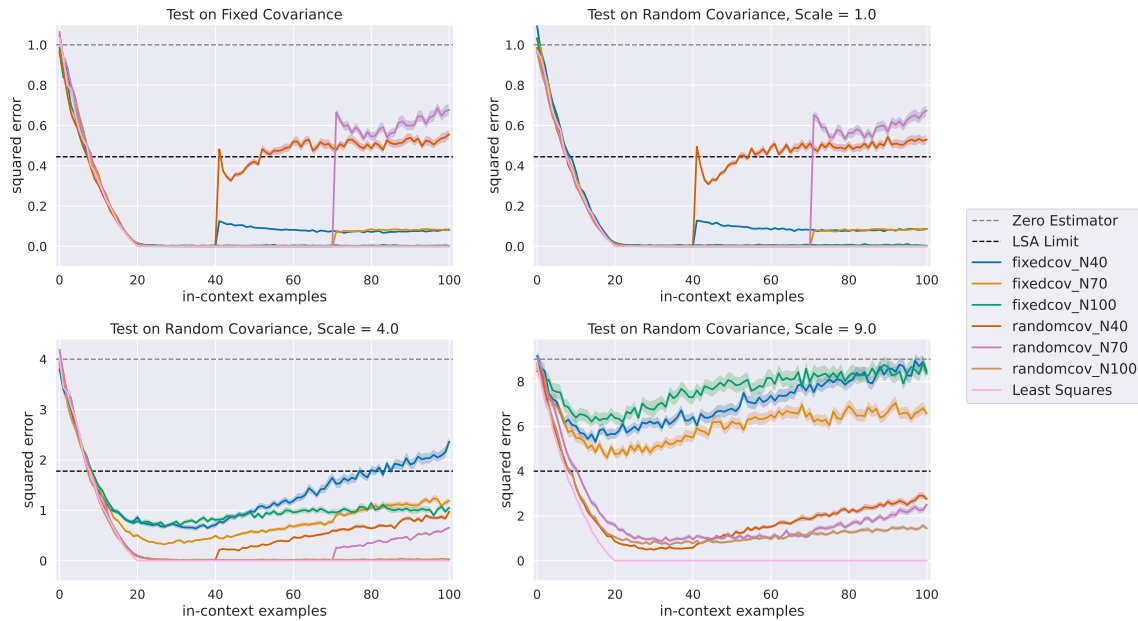


Figure 1: Normalized prediction error for transformers with GPT2 architectures as a function of the number of in-context test examples M when trained on in-context examples of linear models in $d = 20$ dimensions. Colored lines correspond to different training context lengths ($N \in \{40, 70, 100\}$) and different training procedures (either a fixed identity covariance matrix or random diagonal covariance matrices with each diagonal element sampled i.i.d. from the standard exponential distribution). The four figures correspond to evaluating on either fixed covariance or random covariance matrices of different scales. The gray dashed line shows the prediction error of zero estimator and the black dashed line the prediction error of LSA model when $M, N \rightarrow \infty$. The GPT2 models achieve smaller error when they are trained on random covariance matrices with larger contexts, but their prediction error spikes when evaluated on contexts larger than those they were trained on.

Furthermore, we notice some surprising behaviors when the test prompt length exceeds the training prompt length (i.e., $M > N$): there is an evident spike in prediction error, regardless of whether training and testing were performed on fixed or random covariance, and the spike appears to decrease when evaluated on prompts with higher variance. Although we are unsure of why the spike should decrease with higher-variance prompts, the failure of large language models to generalize to larger contexts than they were trained on is a well-known problem (Dai et al., 2019; Anil et al., 2022). In our setting, we conjecture that this spike in error comes from the absolute positional encodings in the GPT2 architecture. The

positional encodings are randomly-initialized and are learnable parameters but the encoding for position i is only updated if the transformer encounters a prompt which has a context of length i . Thus, when evaluating on prompts of length $M > N$, the model is relying upon random positional encodings for $M - N$ samples. We note that a concurrent work has explored the performance of transformers with GPT2 architectures for in-context learning of linear models and found that removing positional encoders improves performance when evaluating on larger contexts (Ahuja et al., 2023). We leave further investigation of this behavior for future work.

5. Proof ideas

In this section, we briefly outline the proof sketch of Theorem 4. The full proof of this theorem is left for Appendix A.

5.1 Equivalence to a quadratic optimization problem

We recall each task τ corresponds to a weight vector $w_\tau \sim \mathbf{N}(0, I_d)$. The prompt inputs for this task are $x_{\tau,j} \stackrel{\text{i.i.d.}}{\sim} \mathbf{N}(0, \Lambda)$, which are also independent of w_τ . The corresponding labels are $y_{\tau,j} = \langle w_\tau, x_{\tau,j} \rangle$. For each task τ , we can form the prompt into a token matrix $E_\tau \in \mathbb{R}^{(d+1) \times (N+1)}$ as in (4), with the right-bottom entry being zero.

The first key step in our proof is to recognize that the prediction $\hat{y}_{\text{query}}(E_\tau; \theta)$ in the linear self-attention model can be written as the output of a quadratic function $u^\top H_\tau u$ for some matrix H_τ depending on the token embedding matrix E_τ and for some vector u depending on $\theta = (W^{KQ}, W^{PV})$. This is shown in the following lemma, the proof of which is provided in Appendix A.1.

Lemma 9 *Let $E_\tau \in \mathbb{R}^{(d+1) \times (N+1)}$ be an embedding matrix corresponding to a prompt of length N and weight w_τ . Then the prediction $\hat{y}_{\text{query}}(E_\tau; \theta)$ for the query covariate can be written as the output of a quadratic function,*

$$\hat{y}_{\text{query}}(E_\tau; \theta) = u^\top H_\tau u,$$

where the matrix H_τ is defined as,

$$H_\tau = \frac{1}{2} X_\tau \otimes \left(\frac{E_\tau E_\tau^\top}{N} \right) \in \mathbb{R}^{(d+1)^2 \times (d+1)^2}, \quad X_\tau = \begin{pmatrix} 0_{d \times d} & x_{\tau, \text{query}} \\ (x_{\tau, \text{query}})^\top & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)} \quad (24)$$

and

$$u = \text{Vec}(U) \in \mathbb{R}^{(d+1)^2}, \quad U = \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)},$$

where $U_{11} = W_{11}^{KQ} \in \mathbb{R}^{d \times d}$, $u_{12} = w_{21}^{PV} \in \mathbb{R}^{d \times 1}$, $u_{21} = w_{21}^{KQ} \in \mathbb{R}^{d \times 1}$, $u_{-1} = w_{22}^{PV} \in \mathbb{R}$ correspond to particular components of W^{PV} and W^{KQ} , defined in (5).

This implies that we can write the original loss function (7) as

$$\widehat{L} = \frac{1}{2B} \sum_{\tau=1}^B \left(u^\top H_\tau u - w_\tau^\top x_{\tau, \text{query}} \right)^2. \quad (25)$$

Thus, our problem is reduced to *understanding the dynamics of an optimization algorithm defined in terms of a quadratic function*. We also note that this quadratic optimization problem is an instance of a rank-one matrix factorization problem, a problem well-studied in the deep learning theory literature (Gunasekar et al., 2017; Arora et al., 2019; Li et al., 2018; Chi et al., 2019; Belabbas, 2020; Li et al., 2020; Jin et al., 2023; Soltanolkotabi et al., 2023).

Note, however, this quadratic function is non-convex. To see this, we will show that H_τ has negative eigenvalues. By standard properties of the Kronecker product, the eigenvalues of $H_\tau = \frac{1}{2} X_\tau \otimes \left(\frac{E_\tau E_\tau^\top}{N} \right)$ are the products of the eigenvalues of $\frac{1}{2} X_\tau$ and the eigenvalues of $\frac{E_\tau E_\tau^\top}{N}$. Since $E_\tau E_\tau^\top$ is symmetric and positive semi-definite, all of its eigenvalues are nonnegative. Since $E_\tau E_\tau^\top$ is nonzero almost surely, it thus has at least one strictly positive eigenvalue. Thus, if X_τ has any negative eigenvalues, H_τ does as well. The characteristic polynomial of X_τ is given by,

$$\det(\mu I - X_\tau) = \det \begin{pmatrix} \mu I_d & -x_{\tau, \text{query}} \\ -x_{\tau, \text{query}}^\top & \mu \end{pmatrix} = \mu^{d-1} \left(\mu^2 - \|x_{\tau, \text{query}}\|_2^2 \right).$$

Therefore, we know almost surely, X_τ has one negative eigenvalue. Thus H_τ has at least $d + 1$ negative eigenvalues, and hence the quadratic form $u^\top H_\tau u$ is non-convex.

5.2 Dynamical system of gradient flow

We now describe the dynamical system for the coordinates of u above. We prove the following lemma in Appendix A.2.

Lemma 10 *Let $u = \text{Vec}(U) := \text{Vec} \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix}$ as in Lemma 9. Consider gradient flow over*

$$L := \frac{1}{2} \mathbb{E} \left(u^\top H_\tau u - w_\tau^\top x_{\tau, \text{query}} \right)^2 \quad (26)$$

with respect to u starting from an initial value satisfying Assumption 3. Then the dynamics of U follows

$$\begin{aligned} \frac{d}{dt} U_{11}(t) &= -u_{-1}^2 \Gamma \Lambda U_{11} \Lambda + u_{-1} \Lambda^2 \\ \frac{d}{dt} u_{-1}(t) &= -\text{tr} \left[u_{-1} \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - \Lambda^2 (U_{11})^\top \right], \end{aligned} \quad (27)$$

and $u_{12}(t) = 0_d, u_{21}(t) = 0_d$ for all $t \geq 0$, where $\Gamma = \left(1 + \frac{1}{N}\right) \Lambda + \frac{1}{N} \text{tr}(\Lambda) I_d \in \mathbb{R}^{d \times d}$.

We see that the dynamics are governed by a complex system of $d^2 + 1$ coupled differential equations. Moreover, basic calculus (for details, see Lemma 15) shows that these dynamics are the same as those of gradient flow on the following objective function:

$$\tilde{\ell} : \mathbb{R}^{d \times d} \times \mathbb{R} \rightarrow \mathbb{R}, \quad \tilde{\ell}(U_{11}, u_{-1}) = \text{tr} \left[\frac{1}{2} u_{-1}^2 \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - u_{-1} \Lambda^2 (U_{11})^\top \right]. \quad (28)$$

Actually, the loss function $\tilde{\ell}$ is simply the loss function L in (26) plus some constants that do not depend on the parameter u . Therefore our problem is reduced to studying the dynamics of gradient flow on the above objective function.

Our next key observation is that the set of global minima for $\tilde{\ell}$ satisfies the condition $u_{-1} U_{11} = \Gamma^{-1}$. Thus, if we can establish global convergence of gradient flow over the above objective function $\tilde{\ell}$, then we have that $u_{-1}(t) U_{11}(t) \rightarrow \Gamma^{-1} \approx_{N \rightarrow \infty} \Lambda^{-1}$.

Lemma 11 *For any global minimum of $\tilde{\ell}$, we have*

$$u_{-1} U_{11} = \Gamma^{-1}. \quad (29)$$

Putting this together with Lemma 10, we see that at those global minima of the population objective satisfying $U_{11} = (c\Gamma)^{-1}$, $u_{-1} = c$ and $u_{12} = u_{21} = 0_d$, the transformer's predictions for a new linear regression task prompt are given by

$$\hat{y}_{\text{query}}(E; \theta) = \frac{1}{M} \sum_{i=1}^M y_i x_i^\top \Gamma^{-1} x_{\text{query}} = w^\top \left(\frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right) \Gamma^{-1} x_{\text{query}} \approx w^\top x_{\text{query}}.$$

Thus, the only remaining task is to show global convergence when gradient flow has an initialization satisfying Assumption 3.

5.3 PL inequality and global convergence

We now show that although the optimization problem is non-convex, a Polyak-Łojasiewicz (PL) inequality holds, which implies that gradient flow converges to a global minimum. Moreover, we can exactly calculate the limiting value of U_{11} and u_{-1} .

Lemma 12 *Suppose the initialization of gradient flow satisfies Assumption 3 with initialization scale satisfying $\sigma^2 < \frac{2}{\sqrt{d} \|\Gamma\|_{\text{op}}}$ for $\Gamma = (1 + \frac{1}{N})\Lambda + \frac{\text{tr}(\Lambda)}{N} I_d$. If we define*

$$\mu := \frac{\sigma^2}{\sqrt{d} \|\Lambda\|_{\text{op}}^2 \text{tr}(\Gamma^{-1} \Lambda^{-1}) \text{tr}(\Lambda^{-1})} \|\Lambda \Theta\|_F^2 \left[2 - \sqrt{d} \sigma^2 \|\Gamma\|_{\text{op}} \right] > 0, \quad (30)$$

then gradient flow on $\tilde{\ell}$ with respect to U_{11} and u_{-1} satisfies, for any $t \geq 0$,

$$\begin{aligned} \left\| \nabla \tilde{\ell}(U_{11}(t), u_{-1}(t)) \right\|_2^2 &:= \left\| \frac{\partial \tilde{\ell}}{\partial U_{11}} \right\|_F^2 + \left| \frac{\partial \tilde{\ell}}{\partial u_{-1}} \right|^2 \\ &\geq \mu \left(\tilde{\ell}(U_{11}(t), u_{-1}(t)) - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) \right). \end{aligned}$$

Moreover, gradient flow converges to the global minimum of $\tilde{\ell}$, and U_{11} and u_{-1} satisfy

$$\lim_{t \rightarrow \infty} u_{-1}(t) = \|\Gamma^{-1}\|_F^{\frac{1}{2}} \quad \text{and} \quad \lim_{t \rightarrow \infty} U_{11}(t) = \|\Gamma^{-1}\|_F^{-\frac{1}{2}} \Gamma^{-1}.$$

With these observations, proving Theorem 4 becomes a direct application of Lemma 9, 10, 11, and Lemma 12. It then only requires translating U_{11} and u_{-1} back to the original parameterization using W^{PV} and W^{KQ} .

6. Conclusion and future work

In this work, we investigated the dynamics of in-context learning of transformers with a single linear self-attention layer under gradient flow on the population loss. In particular, we analyzed the dynamics of these transformers when trained on prompts consisting of random instances of noiseless linear models over anisotropic Gaussian marginals. We showed that despite non-convexity, gradient flow from a suitable random initialization converges to a global minimum of the population objective. We characterized the prediction error of the trained transformer when given a new prompt that consists of a training dataset where the responses are a nonlinear function of the inputs. We showed how the trained transformer is naturally robust to shifts in the task and query distributions but is brittle to distribution shifts between the covariates seen during training and the covariates seen at test time, matching the empirical observations on trained transformer models of Garg et al. (2022).

There are a number of natural directions for future research. First, our results hold for gradient flow on the population loss with a particular class of random initialization schemes. It is a natural question if similar results would hold for stochastic gradient descent with finite step sizes and for more general initializations. Further, we restricted our attention to transformers with a single linear self-attention layer. Although this model class is rich enough to allow for in-context learning of linear predictors, we are particularly interested in understanding the dynamics of in-context learning in nonlinear and deep transformers.

Finally, the framework of in-context learning introduced in prior work was restricted to the setting where the marginal distribution over the covariates (\mathcal{D}_x) was fixed across prompts. This allows for guarantees akin to distribution-specific PAC learning, where the trained transformer is able to achieve small prediction error when given a test prompt consisting of linear regression data when the marginals over the covariates are fixed. However, other learning algorithms (such as ordinary least squares) are able to achieve small prediction error for prompts corresponding to well-specified linear regression tasks for very general classes of distributions over the covariates. As we showed in Section 4.3, when transformers with a single linear self-attention layer are trained on prompts where the covariate distributions are themselves sampled from a distribution, they do not succeed on test prompts with covariate distributions sampled from the same distribution. By contrast, we demonstrated with experiments that larger, nonlinear transformer architectures appear to be more successful in this setting but are still sub-optimal. Developing a better understanding of the dynamics of in-context learning when the covariate distribution varies across prompts is an intriguing direction for future research.

Acknowledgments

We gratefully acknowledge the support of the NSF and the Simons Foundation for the Collaboration on the Theoretical Foundations of Deep Learning through awards DMS-2031883 and #814639, and of the NSF through grant DMS-2023505.

Contents

1	Introduction	2
2	Additional Related Work	3
3	Preliminaries	5
3.1	In-context learning	5
3.2	Linear self-attention networks	6
3.3	Training procedure	8
4	Main results	9
4.1	Convergence of gradient flow and prediction error for new tasks	10
4.2	Behavior of trained transformer under distribution shifts	13
4.3	Transformers trained on prompts with random covariate distributions	14
5	Proof ideas	19
5.1	Equivalence to a quadratic optimization problem	19
5.2	Dynamical system of gradient flow	20
5.3	PL inequality and global convergence	21
6	Conclusion and future work	22
A	Proof of Theorem 4	25
A.1	Proof of Lemma 9	25
A.2	Proof of Lemma 10	26
A.3	Proof of Lemma 11	32
A.4	Proof of Lemma 12	34
B	Proof of Theorem 5	39
C	Proof of Theorem 8	41
C.1	Dynamical system	42
C.2	Loss function and global minima	43
C.3	PL Inequality and global convergence	44
D	Technical lemmas	49
E	Experiment details	51

Appendix A. Proof of Theorem 4

In this section, we prove Lemma 9, Lemma 10, Lemma 11 and Lemma 12. Theorem 4 is a natural corollary of these four lemmas when we translate u_{-1} and U_{11} back to W^{PV} and W^{KQ} .

A.1 Proof of Lemma 9

For the reader's convenience, we restate the lemma below.

Lemma 13 *Let $E_\tau \in \mathbb{R}^{(d+1) \times (N+1)}$ be an embedding matrix corresponding to a prompt of length N and weight w_τ . Then the prediction $\hat{y}_{\text{query}}(E_\tau; \theta)$ for the query covariate can be written as the output of a quadratic function,*

$$\hat{y}_{\text{query}}(E_\tau; \theta) = u^\top H_\tau u,$$

where the matrix H_τ is defined as,

$$H_\tau = \frac{1}{2} X_\tau \otimes \left(\frac{E_\tau E_\tau^\top}{N} \right) \in \mathbb{R}^{(d+1)^2 \times (d+1)^2}, \quad X_\tau = \begin{pmatrix} 0_{d \times d} & x_{\tau, \text{query}} \\ (x_{\tau, \text{query}})^\top & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)} \quad (24)$$

and

$$u = \text{Vec}(U) \in \mathbb{R}^{(d+1)^2}, \quad U = \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)},$$

where $U_{11} = W_{11}^{KQ} \in \mathbb{R}^{d \times d}$, $u_{12} = w_{21}^{PV} \in \mathbb{R}^{d \times 1}$, $u_{21} = w_{21}^{KQ} \in \mathbb{R}^{d \times 1}$, $u_{-1} = w_{22}^{PV} \in \mathbb{R}$ correspond to particular components of W^{PV} and W^{KQ} , defined in (5).

Proof First, we decompose W_{PV} and W_{KQ} in the way above. From the definition, we know $\hat{y}_{\tau, \text{query}}$ is the right-bottom entry of $f_{\text{LSA}}(E_\tau)$, which is

$$\hat{y}_{\tau, \text{query}} = \left((u_{12})^\top \quad u_{-1} \right) \left(\frac{E_\tau E_\tau^\top}{N} \right) \begin{pmatrix} U_{11} \\ (u_{21})^\top \end{pmatrix} x_{\tau, \text{query}}.$$

We denote $u_i \in \mathbb{R}^{d+1}$ as the i -th column of $\begin{pmatrix} U_{11} \\ (u_{21})^\top \end{pmatrix}$ and $x_{\tau, \text{query}}^i$ as the i -th entry of $x_{\tau, \text{query}}$ for $i \in [d]$. Then, we have

$$\begin{aligned} & \hat{y}_{\tau, \text{query}} \\ &= \sum_{i=1}^d x_{\tau, \text{query}}^i \left((u_{12})^\top \quad u_{-1} \right) \left(\frac{E_\tau E_\tau^\top}{N} \right) u_i = \sum_{i=1}^d \text{tr} \left[u_i \left((u_{12})^\top \quad u_{-1} \right) \cdot x_{\tau, \text{query}}^i \left(\frac{E_\tau E_\tau^\top}{N} \right) \right] \\ &= \text{tr} \left[\text{Vec} \left[\begin{pmatrix} U_{11} \\ (u_{21})^\top \end{pmatrix} \right] \left((u_{12})^\top \quad u_{-1} \right) \cdot x_{\tau, \text{query}}^\top \otimes \left(\frac{E_\tau E_\tau^\top}{N} \right) \right] \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{2} \operatorname{tr} \left[\operatorname{Vec} \left[\begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix} \right] \operatorname{Vec}^\top \left[\begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix} \right] \right. \\
 &\quad \left. \times \begin{pmatrix} 0_{d(d+1) \times d(d+1)} & x_{\tau, \text{query}} \otimes \left(\frac{E_\tau E_\tau^\top}{N} \right) \\ x_{\tau, \text{query}}^\top \otimes \left(\frac{E_\tau E_\tau^\top}{N} \right) & 0_{(d+1) \times (d+1)} \end{pmatrix} \right] \\
 &= \frac{1}{2} \operatorname{tr} \left[uu^\top \cdot X_\tau \otimes \left(\frac{E_\tau E_\tau^\top}{N} \right) \right] \\
 &= \langle H_\tau, uu^\top \rangle.
 \end{aligned}$$

Here, we use some algebraic facts about matrix vectorization, Kronecker product and trace. For reference, we refer to (Petersen and Pedersen, 2008). \blacksquare

A.2 Proof of Lemma 10

For the reader's convenience, we restate the lemma below.

Lemma 14 *Let $u = \operatorname{Vec}(U) := \operatorname{Vec} \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix}$ as in Lemma 9. Consider gradient flow over*

$$L := \frac{1}{2} \mathbb{E} \left(u^\top H_\tau u - w_\tau^\top x_{\tau, \text{query}} \right)^2 \quad (26)$$

with respect to u starting from an initial value satisfying Assumption 3. Then the dynamics of U follows

$$\begin{aligned}
 \frac{d}{dt} U_{11}(t) &= -u_{-1}^2 \Gamma \Lambda U_{11} \Lambda + u_{-1} \Lambda^2 \\
 \frac{d}{dt} u_{-1}(t) &= -\operatorname{tr} \left[u_{-1} \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - \Lambda^2 (U_{11})^\top \right],
 \end{aligned} \quad (27)$$

and $u_{12}(t) = 0_d, u_{21}(t) = 0_d$ for all $t \geq 0$, where $\Gamma = (1 + \frac{1}{N}) \Lambda + \frac{1}{N} \operatorname{tr}(\Lambda) I_d \in \mathbb{R}^{d \times d}$.

Proof From the definition of L in (26) and the dynamics of gradient flow, we calculate the derivatives of u . Here, we use the chain rule and some facts about matrix derivatives. See Lemma 29 for reference.

$$\frac{du}{dt} = -\mathbb{E} \left(\langle H_\tau, uu^\top \rangle H_\tau \right) u + \mathbb{E} \left(w_\tau^\top x_{\tau, \text{query}} H_\tau \right) u. \quad (31)$$

Step One: Calculate the Second Term We first calculate the second term. From the definition of H_τ , we have

$$\mathbb{E} \left[w_\tau^\top x_{\tau, \text{query}} H_\tau \right] = \frac{1}{2} \sum_{i=1}^d \mathbb{E} \left[\left(x_{\tau, \text{query}}^i X_\tau \right) \otimes \left(w_\tau^i \frac{E_\tau E_\tau^\top}{N} \right) \right].$$

For ease of notation, we denote

$$\widehat{\Lambda}_\tau := \frac{1}{N} \sum_{i=1}^N x_{\tau,i} x_{\tau,i}^\top. \quad (32)$$

Then, from the definition of $\frac{E_\tau E_\tau^\top}{N}$, we know

$$\frac{E_\tau E_\tau^\top}{N} = \begin{pmatrix} \widehat{\Lambda}_\tau + \frac{1}{N} x_{\tau,\text{query}} \cdot x_{\tau,\text{query}}^\top & \widehat{\Lambda}_\tau w_\tau \\ w_\tau^\top \widehat{\Lambda}_\tau & w_\tau^\top \widehat{\Lambda}_\tau w_\tau \end{pmatrix}.$$

Since $w_\tau \sim \mathcal{N}(0, I_d)$ is independent of all prompt inputs and query input, we have

$$\begin{aligned} & \frac{1}{2} \sum_{i=1}^d \mathbb{E} \left[(x_{\tau,\text{query}}^i X_\tau) \otimes \left(\frac{w_\tau^i}{N} \begin{pmatrix} x_{\tau,\text{query}} \cdot x_{\tau,\text{query}}^\top & 0 \\ 0 & 0 \end{pmatrix} \right) \right] \\ &= \frac{1}{2} \sum_{i=1}^d \mathbb{E} \left[\mathbb{E} \left[(x_{\tau,\text{query}}^i X_\tau) \otimes \left(\frac{w_\tau^i}{N} \begin{pmatrix} x_{\tau,\text{query}} \cdot x_{\tau,\text{query}}^\top & 0 \\ 0 & 0 \end{pmatrix} \right) \middle| x_{\tau,\text{query}} \right] \right] \\ &= \frac{1}{2} \sum_{i=1}^d \mathbb{E} \left[(x_{\tau,\text{query}}^i X_\tau) \otimes \left(\frac{\mathbb{E}[w_\tau^i | x_{\tau,\text{query}}]}{N} \begin{pmatrix} x_{\tau,\text{query}} \cdot x_{\tau,\text{query}}^\top & 0 \\ 0 & 0 \end{pmatrix} \right) \right] = 0. \end{aligned}$$

Therefore, we have

$$\mathbb{E} \left[w_\tau^\top x_{\tau,\text{query}} H_\tau \right] = \frac{1}{2} \sum_{i=1}^d \mathbb{E} \left[(x_{\tau,\text{query}}^i X_\tau) \otimes \left(w_\tau^i \begin{pmatrix} \widehat{\Lambda}_\tau & \widehat{\Lambda}_\tau w_\tau \\ w_\tau^\top \widehat{\Lambda}_\tau & w_\tau^\top \widehat{\Lambda}_\tau w_\tau \end{pmatrix} \right) \right].$$

Since X_τ only depends on $x_{\tau,\text{query}}$ by definition, and $x_{\tau,\text{query}}$ is independent of w_τ and $x_{\tau,i}$, $i = 1, 2, \dots, N$, we have

$$\begin{aligned} \mathbb{E} \left[w_\tau^\top x_{\tau,\text{query}} H_\tau \right] &= \frac{1}{2} \sum_{i=1}^d \left[\mathbb{E} (x_{\tau,\text{query}}^i X_\tau) \otimes \mathbb{E} \left(w_\tau^i \begin{pmatrix} \widehat{\Lambda}_\tau & \widehat{\Lambda}_\tau w_\tau \\ w_\tau^\top \widehat{\Lambda}_\tau & w_\tau^\top \widehat{\Lambda}_\tau w_\tau \end{pmatrix} \right) \right] \\ &= \frac{1}{2} \sum_{i=1}^d \left[\begin{pmatrix} 0_{d \times d} & \Lambda_i \\ \Lambda_i^\top & 0 \end{pmatrix} \otimes \begin{pmatrix} \mathbb{E}(w_\tau^i) \Lambda & \Lambda \mathbb{E}(w_\tau^i w_\tau) \\ \mathbb{E}(w_\tau^i w_\tau^\top) \Lambda & \mathbb{E}(w_\tau^i w_\tau^\top \Lambda w_\tau) \end{pmatrix} \right] \\ &= \frac{1}{2} \sum_{i=1}^d \begin{pmatrix} 0_{d \times d} & \Lambda_i \\ \Lambda_i^\top & 0 \end{pmatrix} \otimes \begin{pmatrix} 0_{d \times d} & \Lambda_i \\ \Lambda_i^\top & 0 \end{pmatrix}, \end{aligned}$$

where Λ_i denotes $\Lambda_{\cdot i}$. Here, the second line comes from the fact that $\mathbb{E} \widehat{\Lambda}_\tau = \Lambda$, and that w_τ is independent of all prompt input and query input. The last line comes from the fact that $w_\tau \sim \mathcal{N}(0, I_d)$. Therefore, simple computation shows that

$$\mathbb{E} \left[w_\tau^\top x_{\tau,\text{query}} H_\tau \right] u = \frac{1}{2} \begin{pmatrix} \mathbf{0}_{d(d+1) \times d(d+1)} & A \\ A^\top & \mathbf{0}_{(d+1) \times (d+1)} \end{pmatrix} u, \quad (33)$$

where $A \in \mathbb{R}^{d(d+1) \times (d+1)}$ and $V_j \in \mathbb{R}^{(d+1) \times (d+1)}$ are defined by

$$A = \begin{pmatrix} V_1 + V_1^\top \\ V_2 + V_2^\top \\ \dots \\ V_d + V_d^\top \end{pmatrix}, \quad V_j = \begin{pmatrix} 0_{d \times d} & \sum_{i=1}^d \Lambda_{ij} \Lambda_i \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0_{d \times d} & \Lambda \Lambda_j \\ 0 & 0 \end{pmatrix}. \quad (34)$$

Step Two: Calculate the First Term Next, we compute the first term in (31), namely

$$D := 2\mathbb{E} \left(\langle H_\tau, uu^\top \rangle H_\tau u \right).$$

For simplicity, we denote $Z_\tau := \frac{1}{N} E_\tau E_\tau^\top$. Using the definition of H_τ in (24) and Lemma 29, we have

$$\begin{aligned} D &= 2\mathbb{E} \left(\langle H_\tau, uu^\top \rangle H_\tau u \right) && \text{(definition)} \\ &= \frac{1}{2} \mathbb{E} \left[\text{tr} \left(X_\tau \otimes Z_\tau \text{Vec}(U) \text{Vec}(U)^\top \right) (X_\tau \otimes Z_\tau) \text{Vec}(U) \right] && \text{(definition of } H_\tau \text{ in (24) and } u = \text{Vec}(U)) \\ &= \frac{1}{2} \mathbb{E} \left[\text{tr} \left(\text{Vec}(Z_\tau U X_\tau) \text{Vec}(U)^\top \right) \text{Vec}(Z_\tau U X_\tau) \right] && \text{(Vec}(AXB) = (B^\top \otimes A) \text{Vec}(X) \text{ in Lemma 29)} \\ &= \frac{1}{2} \mathbb{E} \left[\text{Vec}(U)^\top \text{Vec}(Z_\tau U X_\tau) \cdot \text{Vec}(Z_\tau U X_\tau) \right] && \text{(property of trace operator)} \\ &= \frac{1}{2} \mathbb{E} \left[\sum_{i,j=1}^{d+1} \left((Z_\tau U X_\tau)_{ij} U_{ij} \right) \text{Vec}(Z_\tau U X_\tau) \right]. \end{aligned}$$

Step Three: u_{12} and u_{21} Vanish We first prove that if $u_{12} = u_{21} = 0_d$, then $\frac{d}{dt} u_{12} = 0_d$ and $\frac{d}{dt} u_{21} = 0_d$. If this is true, then these two blocks will be zero all the time since we assume they are zero at initial time in Assumption 3. We denote $A_{k\cdot}$ and $A_{\cdot k}$ as the k -th row and k -th column of matrix A , respectively.

Under the assumption that $u_{12} = u_{21} = 0_d$, we first compute

$$(Z_\tau U X_\tau) = \begin{pmatrix} \widehat{\Lambda}_\tau w_\tau u_{-1} x_{\tau, \text{query}}^\top & \left(\widehat{\Lambda}_\tau + \frac{1}{N} x_{\tau, \text{query}} \cdot x_{\tau, \text{query}}^\top \right) U_{11} x_{\tau, \text{query}} \\ w_\tau^\top \left(\widehat{\Lambda}_\tau \right) w_\tau u_{-1} x_{\tau, \text{query}}^\top & w_\tau^\top \left(\widehat{\Lambda}_\tau \right) U_{11} x_{\tau, \text{query}} \end{pmatrix}.$$

Written in an entry-wise manner, it will be

$$(Z_\tau U X_\tau)_{kl} = \begin{cases} \left(\widehat{\Lambda}_\tau \right)_{k\cdot} w_\tau u_{-1} x_{\tau, \text{query}}^l & k, l \in [d] \\ \left(\widehat{\Lambda}_\tau + \frac{1}{N} x_{\tau, \text{query}} \cdot x_{\tau, \text{query}}^\top \right)_{k\cdot} U_{11} x_{\tau, \text{query}} & k \in [d], l = d+1 \\ w_\tau^\top \left(\widehat{\Lambda}_\tau \right) w_\tau u_{-1} x_{\tau, \text{query}}^l & l \in [d], k = d+1 \\ w_\tau^\top \left(\widehat{\Lambda}_\tau \right) U_{11} x_{\tau, \text{query}} & k = l = d+1 \end{cases}. \quad (35)$$

We use D_{ij} to denote the (i, j) -th entry of the $(d+1) \times (d+1)$ matrix \bar{D} such that $\text{Vec}(\bar{D}) = D$. Now we fix a $k \in [d]$, then

$$\begin{aligned} D_{k,d+1} &= \frac{1}{2} \mathbb{E} \left[\sum_{i,j=1}^{d+1} \left((Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{k,d+1} \right] \\ &= \frac{1}{2} \mathbb{E} \left[\sum_{i,j=1}^d \left((Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{k,d+1} \right] \\ &\quad + \frac{1}{2} \mathbb{E} \left[\left((Z_\tau U X_\tau)_{d+1,d+1} u_{-1} \right) (Z_\tau U X_\tau)_{k,d+1} \right], \end{aligned} \quad (36)$$

since $U_{i,d+1} = U_{d+1,i} = 0$ for any $i \in [d]$. For the first term in the right hand side of last equation, we fix $i, j \in [d]$ and have

$$\begin{aligned} &\mathbb{E} \left((Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{k,d+1} \\ &= \mathbb{E} \left(U_{ij} \left(\hat{\Lambda}_\tau \right)_i w_\tau u_{-1} x_{\tau,\text{query}}^j \cdot \left(\hat{\Lambda}_\tau + \frac{1}{N} x_{\tau,\text{query}} \cdot x_{\tau,\text{query}}^\top \right)_{k:} U_{11} x_{\tau,\text{query}} \right) = 0, \end{aligned}$$

since w_τ is independent with all prompt input and query input, namely all $x_{\tau,i}$ for $i \in [\text{query}]$, and w_τ is mean zero. Similarly, for the second term of (36), we have

$$\begin{aligned} &\mathbb{E} \left((Z_\tau U X_\tau)_{d+1,d+1} u_{-1} \right) (Z_\tau U X_\tau)_{k,d+1} \\ &= \mathbb{E} \left(u_{-1} w_\tau^\top \left(\hat{\Lambda}_\tau \right) U_{11} x_{\tau,\text{query}} \cdot \left(\hat{\Lambda}_\tau + \frac{1}{N} x_{\tau,\text{query}} \cdot x_{\tau,\text{query}}^\top \right)_{k:} U_{11} x_{\tau,\text{query}} \right) = 0 \end{aligned}$$

since $\mathbb{E}(w_\tau^\top) = 0$ and w_τ is independent of all $x_{\tau,i}$ for $i \in [\text{query}]$. Therefore, we have $D_{k,d+1} = 0$ for $k \in [d]$. Similar calculation shows that $D_{d+1,k} = 0$ for $k \in [d]$.

For $k \in [d]$, to calculate the derivative of $U_{k,d+1}$, it suffices to further calculate the inner product of the $d(d+1) + k$ th row of $\mathbb{E} [w_\tau^\top x_{\tau,\text{query}} H_\tau]$ and u . From (33), we know this is

$$\frac{1}{2} \sum_{j=1}^d \Lambda_k^\top \Lambda_j U_{d+1,j} = 0$$

given that $u_{12} = u_{21} = 0_d$. Therefore, we conclude that the derivative of $U_{k,d+1}$ will vanish given $u_{12} = u_{21} = 0_d$. Similarly, we conclude the same result for $U_{d+1,k}$ for $k \in [d]$. Therefore, we know $u_{12} = 0_d$ and $u_{21} = 0_d$ for all time $t \geq 0$.

Step Four: Dynamics of U_{11} Next, we calculate the derivatives of U_{11} given $u_{12} = u_{21} = 0_d$. For a fixed pair of $k, l \in [d]$, we have

$$D_{kl} = \frac{1}{2} \mathbb{E} \left[\sum_{i,j=1}^d \left((Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{kl} \right] + \frac{1}{2} \mathbb{E} \left[\left((Z_\tau U X_\tau)_{d+1,d+1} u_{-1} \right) (Z_\tau U X_\tau)_{kl} \right].$$

For fixed $i, j \in [d]$, we have

$$\begin{aligned}
 \mathbb{E} \left[\left((Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{kl} \right] &= U_{ij} u_{-1}^2 \mathbb{E} \left[\left(\widehat{\Lambda}_\tau \right)_{i:} w_\tau x_{\tau, \text{query}}^j x_{\tau, \text{query}}^l w_\tau^\top \left(\widehat{\Lambda}_\tau \right)_{:k} \right] \\
 &= U_{ij} u_{-1}^2 \mathbb{E} \left[x_{\tau, \text{query}}^j x_{\tau, \text{query}}^l \right] \cdot \mathbb{E} \left[\left(\widehat{\Lambda}_\tau \right)_{i:} \left(\widehat{\Lambda}_\tau \right)_{:k} \right] \\
 &= U_{ij} u_{-1}^2 \Lambda_{\tau, jl} \mathbb{E} \left[\left(\widehat{\Lambda}_\tau \right)_{i:} \left(\widehat{\Lambda}_\tau \right)_{:k} \right].
 \end{aligned}$$

Therefore, we sum over $i, j \in [d]$ to get

$$\frac{1}{2} \mathbb{E} \left[\sum_{i, j=1}^d \left((Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{kl} \right] = \frac{1}{2} u_{-1}^2 \mathbb{E} \left(\left(\widehat{\Lambda}_\tau \right)_{k:} \left(\widehat{\Lambda}_\tau \right) \right) U_{11} \Lambda_l$$

For the last term, we have

$$\frac{1}{2} \mathbb{E} \left[\left((Z_\tau U X_\tau)_{d+1, d+1} u_{-1} \right) (Z_\tau U X_\tau)_{kl} \right] = \frac{1}{2} u_{-1}^2 \mathbb{E} \left(\left(\widehat{\Lambda}_\tau \right)_{k:} \left(\widehat{\Lambda}_\tau \right) \right) U_{11} \Lambda_l.$$

So we have

$$D_{kl} = u_{-1}^2 \mathbb{E} \left(\left(\widehat{\Lambda}_\tau \right)_{k:} \left(\widehat{\Lambda}_\tau \right) \right) U_{11} \Lambda_l.$$

Additionally, we have

$$\begin{aligned}
 2 \left[\mathbb{E} \left(w_\tau^\top x_{\tau, \text{query}} H_\tau \right) u \right]_{(l-1)(d+1)+k} &= \left[\left(\begin{array}{cc} \mathbf{0}_{d(d+1) \times d(d+1)} & A \\ A^\top & \mathbf{0}_{(d+1) \times (d+1)} \end{array} \right) \cdot u \right]_{(l-1)(d+1)+k} \\
 &\hspace{15em} \text{(definition)} \\
 &= \left(\mathbf{0}_{(d+1) \times d(d+1)} \quad V_l + V_l^\top \right)_{k:} \cdot U \\
 &\hspace{15em} \text{(definition of } A \text{ in (34))} \\
 &= \Lambda_k^\top \Lambda_l u_{-1}. \hspace{15em} \text{(definition of } V_i \text{ in (34))}
 \end{aligned}$$

Therefore, we have that for $k, l \in [d]$, the dynamics of U_{kl} is

$$\frac{d}{dt} U_{kl} = -u_{-1}^2 \mathbb{E} \left(\left(\widehat{\Lambda}_\tau \right)_{k:} \left(\widehat{\Lambda}_\tau \right) \right) U_{11} \Lambda_l + u_{-1} \Lambda_k^\top \Lambda_l,$$

which implies

$$\frac{d}{dt} U_{11} = -u_{-1}^2 \mathbb{E} \left(\left(\widehat{\Lambda}_\tau \right)^2 \right) U_{11} \Lambda + u_{-1} \Lambda^2.$$

From the definition of $\widehat{\Lambda}_\tau$ (equation (32)), the independence and Gaussianity of $x_{\tau, i}$ and Lemma 30, we compute

$$\mathbb{E} \left(\left(\widehat{\Lambda}_\tau \right)^2 \right) = \mathbb{E} \left(\left(\frac{1}{N} \sum_{i=1}^N x_{\tau, i} x_{\tau, i}^\top \right)^2 \right) \hspace{10em} \text{(definition (32))}$$

$$\begin{aligned}
 &= \frac{N-1}{N} \left[\mathbb{E} \left(x_{\tau,1} x_{\tau,1}^\top \right) \right]^2 + \frac{1}{N} \mathbb{E} \left(x_{\tau,1} x_{\tau,1}^\top x_{\tau,1} x_{\tau,1}^\top \right) \\
 &\hspace{15em} \text{(independence between prompt input)} \\
 &= \frac{N+1}{N} \Lambda^2 + \frac{1}{N} \text{tr}(\Lambda) \Lambda. \hspace{10em} \text{(Lemma 30)}
 \end{aligned}$$

We define

$$\Gamma := \frac{N+1}{N} \Lambda + \frac{1}{N} \text{tr}(\Lambda) I_d. \tag{37}$$

Then, from (31), we know the dynamics of U_{11} is

$$\frac{d}{dt} U_{11} = -u_{-1}^2 \Gamma \Lambda U_{11} \Lambda + u_{-1} \Lambda^2. \tag{38}$$

Step Five: Dynamics of u_{-1} Finally, we compute the dynamics of u_{-1} . We have

$$\begin{aligned}
 D_{d+1,d+1} &= \frac{1}{2} \mathbb{E} \left[\sum_{i,j=1}^d \left((Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{d+1,d+1} \right] \\
 &\quad + \frac{1}{2} \mathbb{E} \left[\left((Z_\tau U X_\tau)_{d+1,d+1} u_{-1} \right) (Z_\tau U X_\tau)_{d+1,d+1} \right]. \tag{39}
 \end{aligned}$$

For the first term above, we have

$$\begin{aligned}
 &\mathbb{E} \left[\sum_{i,j=1}^d \left((Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{d+1,d+1} \right] \\
 &= u_{-1} \sum_{i,j=1}^d U_{ij} \mathbb{E} \left[\left(\widehat{\Lambda}_\tau \right)_i \cdot w_\tau w_\tau^\top \cdot \left(\widehat{\Lambda}_\tau \right) \cdot U_{11} x_{\tau,\text{query}} x_{\tau,\text{query}}^j \right] \hspace{5em} \text{(from (35))} \\
 &= u_{-1} \sum_{i,j=1}^d U_{ij} \mathbb{E} \left[\left(\widehat{\Lambda}_\tau \right)_i \cdot \left(\widehat{\Lambda}_\tau \right) \cdot U_{11} x_{\tau,\text{query}} x_{\tau,\text{query}}^j \right] \text{ (independence and distribution of } w_\tau) \\
 &= u_{-1} \sum_{i,j=1}^d U_{ij} \mathbb{E} \left[\left(\widehat{\Lambda}_\tau \right)_i \cdot \left(\widehat{\Lambda}_\tau \right) \cdot U_{11} \Lambda_j \right] \hspace{5em} \text{(independence between prompt covariates)} \\
 &= u_{-1} \mathbb{E} \text{tr} \left[\sum_{i,j=1}^d \Lambda_j U_{ij} \left(\widehat{\Lambda}_\tau \right)_i \cdot \left(\widehat{\Lambda}_\tau \right) U_{11} \right] = u_{-1} \mathbb{E} \text{tr} \left[\Lambda (U_{11})^\top \left(\widehat{\Lambda}_\tau \right)^2 U_{11} \right] \\
 &= u_{-1} \text{tr} \left[\mathbb{E} \left(\widehat{\Lambda}_\tau \right)^2 U_{11} \Lambda (U_{11})^\top \right].
 \end{aligned}$$

For the second term in (39), we have

$$\begin{aligned}
 &\mathbb{E} \left[\left((Z_\tau U X_\tau)_{d+1,d+1} u_{-1} \right) (Z_\tau U X_\tau)_{d+1,d+1} \right] \\
 &= u_{-1} \mathbb{E} \left[w_\tau^\top \left(\widehat{\Lambda}_\tau \right) U_{11} x_{\tau,\text{query}} x_{\tau,\text{query}}^\top (U_{11})^\top \left(\widehat{\Lambda}_\tau \right) w_\tau \right] \hspace{5em} \text{(from (35))} \\
 &= u_{-1} \mathbb{E} \text{tr} \left[w_\tau w_\tau^\top \left(\widehat{\Lambda}_\tau \right) U_{11} x_{\tau,\text{query}} x_{\tau,\text{query}}^\top (U_{11})^\top \left(\widehat{\Lambda}_\tau \right) \right]
 \end{aligned}$$

$$\begin{aligned}
 &= u_{-1} \mathbb{E} \operatorname{tr} \left[\left(\widehat{\Lambda}_\tau \right) U_{11} \Lambda(U_{11})^\top \left(\widehat{\Lambda}_\tau \right) \right] \\
 &= u_{-1} \operatorname{tr} \left[\mathbb{E} \left(\widehat{\Lambda}_\tau \right)^2 U_{11} \Lambda(U_{11})^\top \right].
 \end{aligned}$$

Therefore, we know

$$D_{d+1,d+1} = u_{-1} \operatorname{tr} \left[\mathbb{E} \left(\widehat{\Lambda}_\tau \right)^2 U_{11} \Lambda(U_{11})^\top \right].$$

Additionally, we have

$$\begin{aligned}
 2 \left[\mathbb{E} \left(w_\tau^\top x_{\tau, \text{query}} H_\tau \right) u \right]_{(d+1)^2} &= \left[\begin{pmatrix} \mathbf{0}_{d(d+1) \times d(d+1)} & A \\ A^\top & \mathbf{0}_{(d+1) \times (d+1)} \end{pmatrix} \cdot u \right]_{(d+1)^2} \quad (\text{from (33)}) \\
 &= \left(V_1 + V_1^\top \quad \dots \quad V_d + V_d^\top \quad \mathbf{0}_{(d+1) \times (d+1)} \right)_{d+1:} \cdot U \\
 &\quad (\text{definition of } A \text{ in (34)}) \\
 &= \sum_{i,j=1}^d \Lambda_i^\top \Lambda_j U_{ji} = \operatorname{tr} \left(\Lambda(U_{11})^\top \Lambda \right).
 \end{aligned}$$

Then, from (31), we have the dynamics of u_{-1} is

$$\frac{d}{dt} u_{-1} = - \operatorname{tr} \left[u_{-1} \Gamma \Lambda U_{11} \Lambda(U_{11})^\top - \Lambda^2(U_{11})^\top \right]. \quad (40)$$

■

A.3 Proof of Lemma 11

Lemma 11 gives the form of global minima of an equivalent loss function. First, we prove that gradient flow on L defined in (8) from the initial values satisfying Assumption 3 is equivalent to gradient flow on another loss function $\tilde{\ell}$ defined below. Then, we derive an expression for the global minima of this loss function.

First, from the dynamics of gradient flow, we can actually recover the loss function up to a constant. We have the following lemma.

Lemma 15 (Loss Function) *Consider gradient flow over L in (26) with respect to u starting from an initial value satisfying Assumption 3. This is equivalent to doing gradient flow with respect to U_{11} and u_{-1} on the loss function*

$$\tilde{\ell}(U_{11}, u_{-1}) = \operatorname{tr} \left[\frac{1}{2} u_{-1}^2 \Gamma \Lambda U_{11} \Lambda(U_{11})^\top - u_{-1} \Lambda^2(U_{11})^\top \right]. \quad (41)$$

Proof The proof is simply by taking gradient of the loss function in (41). For techniques in matrix derivatives, see Lemma 29. We take the gradient of $\tilde{\ell}$ on U_{11} to obtain

$$\frac{\partial \tilde{\ell}}{\partial U_{11}} = \frac{1}{2} u_{-1}^2 \Lambda^\top \Gamma^\top U_{11} \Lambda^\top + \frac{1}{2} u_{-1}^2 \Gamma \Lambda U_{11} \Lambda - u_{-1} \Lambda^2 = u_{-1}^2 \Gamma \Lambda U_{11} \Lambda - u_{-1} \Lambda^2,$$

since Γ and Λ are commutable. We take derivatives w.r.t. u_{-1} to get

$$\frac{\partial \tilde{\ell}}{\partial u_{-1}} = \text{tr} \left[u_{-1} \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - \Lambda^2 (U_{11})^\top \right].$$

Combining this with Lemma 10, we have

$$\frac{d}{dt} U_{11}(t) = -\frac{\partial \tilde{\ell}}{\partial U_{11}}, \quad \frac{d}{dt} u_{-1}(t) = -\frac{\partial \tilde{\ell}}{\partial u_{-1}}.$$

■

We remark that actually this is the loss function L up to some constant. This loss function $\tilde{\ell}$ can be negative. But we can still compute its global minima as follows.

Corollary 16 (Minimum of Loss Function) *The loss function $\tilde{\ell}$ in Lemma 15 satisfies*

$$\min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) = -\frac{1}{2} \text{tr} [\Lambda^2 \Gamma^{-1}]$$

and

$$\tilde{\ell}(U_{11}, u_{-1}) - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) = \frac{1}{2} \left\| \Gamma^{\frac{1}{2}} \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \right\|_F^2.$$

Proof First, we claim that

$$\tilde{\ell}(U_{11}, u_{-1}) = \frac{1}{2} \text{tr} \left[\Gamma \cdot \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right)^\top \right] - \frac{1}{2} \text{tr} [\Lambda^2 \Gamma^{-1}].$$

To calculate this, we just need to expand the terms in the brackets and notice that Γ and Λ commute:

$$\begin{aligned} & \text{tr} \left[\Gamma \cdot \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right)^\top \right] - \text{tr} [\Lambda^2 \Gamma^{-1}] \\ \stackrel{(i)}{=} & \text{tr} \left[\Gamma \cdot \left(u_{-1}^2 \Lambda^{\frac{1}{2}} U_{11} \Lambda (U_{11})^\top \Lambda^{1/2} - u_{-1} \Lambda \Gamma^{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{3}{2}} \Gamma^{-1} + \Gamma^{-2} \Lambda^2 \right) \right] \\ & - \text{tr} [\Lambda^2 \Gamma^{-1}] \\ = & \text{tr} \left[\Gamma \cdot \left(u_{-1}^2 \Lambda^{\frac{1}{2}} U_{11} \Lambda (U_{11})^\top \Lambda^{1/2} - u_{-1} \Lambda \Gamma^{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{3}{2}} \Gamma^{-1} \right) \right] \\ = & u_{-1}^2 \text{tr} \left[\Gamma \Lambda^{\frac{1}{2}} U_{11} \Lambda (U_{11})^\top \Lambda^{\frac{1}{2}} \right] - u_{-1} \text{tr} \left[\Gamma \Lambda \Gamma^{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Gamma \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{3}{2}} \Gamma^{-1} \right] \\ \stackrel{(ii)}{=} & u_{-1}^2 \text{tr} \left[\Gamma \Lambda U_{11} \Lambda (U_{11})^\top \right] - 2u_{-1} \text{tr} \left[\Lambda^2 U_{11} \Lambda^{\frac{1}{2}} \right] \\ = & 2\tilde{\ell}(U_{11}, u_{-1}). \end{aligned}$$

Equations (i) and (ii) use that Γ and Λ commute.

Since $\Gamma \succeq 0$ and $\left(u_{-1}\Lambda^{\frac{1}{2}}U_{11}\Lambda^{\frac{1}{2}} - \Lambda\Gamma^{-1}\right)\left(u_{-1}\Lambda^{\frac{1}{2}}U_{11}\Lambda^{\frac{1}{2}} - \Lambda\Gamma^{-1}\right)^\top \succeq 0$, we know from Lemma 32 that

$$\frac{1}{2} \operatorname{tr} \left[\Gamma \cdot \left(u_{-1}\Lambda^{\frac{1}{2}}U_{11}\Lambda^{\frac{1}{2}} - \Lambda\Gamma^{-1}\right) \left(u_{-1}\Lambda^{\frac{1}{2}}U_{11}\Lambda^{\frac{1}{2}} - \Lambda\Gamma^{-1}\right)^\top \right] \geq 0,$$

which implies

$$\tilde{\ell}(U_{11}, u_{-1}) \geq -\frac{1}{2} \operatorname{tr} [\Lambda^2 \Gamma^{-1}].$$

Equality holds when

$$U_{11} = \Gamma^{-1}, \quad u_{-1} = 1,$$

so the minimum of $\tilde{\ell}$ is $-\frac{1}{2} \operatorname{tr} [\Lambda^2 \Gamma^{-1}]$. The expression for $\tilde{\ell}(U_{11}, u_{-1}) - \min \tilde{\ell}(U_{11}, u_{-1})$ comes from the fact that $\operatorname{tr}(A^\top A) = \|A\|_F^2$ for any matrix A . \blacksquare

Lemma 11 is an immediate consequence of Corollary16, since the loss will keep the same when we replace (U_{11}, u_{-1}) by $(cU_{11}, c^{-1}u_{-1})$ for any non-zero constant c .

A.4 Proof of Lemma 12

In this section, we prove that the dynamical system in Lemma 10 satisfies a PL inequality. Then, the PL inequality naturally leads to the global convergence of this dynamical system. First, we prove a simple lemma, which says the parameters in the LSA model will keep 'balanced' in the whole trajectory. From the proof of this lemma, we can understand why we assume a balanced parameter at the initial time.

Lemma 17 (Balanced Parameters) *Consider gradient flow over L in (26) with respect to u starting from an initial value satisfying Assumption 3. For any $t \geq 0$, it holds that*

$$u_{-1}^2 = \operatorname{tr} [U_{11}(U_{11})^\top]. \quad (42)$$

Proof From Lemma 10, we multiply the first equation in (27) by $(U_{11})^\top$ from the right to get

$$\left(\frac{d}{dt}U_{11}(t)\right)(U_{11}(t))^\top = -u_{-1}^2\Gamma\Lambda U_{11}\Lambda(U_{11})^\top + u_{-1}\Lambda^2(U_{11})^\top.$$

Also we multiply the second equation in Lemma 10 by u_{-1} to obtain

$$\left(\frac{d}{dt}u_{-1}(t)\right)u_{-1}(t) = \operatorname{tr} \left[-u_{-1}^2\Gamma\Lambda U_{11}\Lambda(U_{11})^\top + u_{-1}\Lambda^2(U_{11})^\top\right].$$

Therefore, we have

$$\operatorname{tr} \left[\left(\frac{d}{dt}U_{11}(t)\right)(U_{11}(t))^\top\right] = \left(\frac{d}{dt}u_{-1}(t)\right)u_{-1}(t).$$

Taking the transpose of the equation above and adding to itself gives

$$\frac{d}{dt} \operatorname{tr} [U_{11}(t)(U_{11}(t))^\top] = \frac{d}{dt} (u_{-1}(t)^2).$$

Notice that from Assumption 3, we know that at $t = 0$,

$$u_{-1}(0)^2 = \sigma^2 = \sigma^2 \operatorname{tr} [\Theta \Theta^\top \Theta \Theta^\top] = \operatorname{tr} [U_{11}(0)(U_{11}(0))^\top].$$

So for any time $t \geq 0$, the equation holds. \blacksquare

In order to prove the PL inequality, we first prove an important property which says the trajectories of $u_{-1}(t)$ stay away from saddle point at origin. First, we prove that $u_{-1}(t)$ will stay positive along the whole trajectory.

Lemma 18 *Consider gradient flow over L in (26) with respect to u starting from an initial value satisfying Assumption 3. If the initial scale satisfies*

$$0 < \sigma < \sqrt{\frac{2}{\sqrt{d} \|\Gamma\|_{op}}}, \quad (43)$$

then, for any $t \geq 0$, it holds that

$$u_{-1} > 0.$$

Proof From Lemma 15, we are actually doing gradient flow on the loss $\tilde{\ell}$. The loss function is non-increasing, because

$$\frac{d\tilde{\ell}}{dt} = \left\langle \frac{dU_{11}}{dt}, \frac{\partial \tilde{\ell}}{\partial U_{11}} \right\rangle + \left\langle \frac{du_{-1}}{dt}, \frac{\partial \tilde{\ell}}{\partial u_{-1}} \right\rangle = - \left\| \frac{dU_{11}}{dt} \right\|_F^2 - \left\| \frac{du_{-1}}{dt} \right\|_F^2 \leq 0.$$

We notice that when $u_{-1} = 0$, $\tilde{\ell} = 0$. Therefore, as long as $\tilde{\ell}(U_{11}(0), u_{-1}(0)) < 0$, then for any time, u_{-1} will be non-zero. Further, since $u_{-1}(0) > 0$ and the trajectory of $u_{-1}(t)$ must be continuous, we know $u_{-1}(t) > 0$ for any $t \geq 0$.

Then, it suffices to prove when $0 < \sigma < \sqrt{\frac{2}{\sqrt{d} \|\Gamma\|_{op}}}$, it holds that $\tilde{\ell}(U_{11}(0), u_{-1}(0)) < 0$. From Assumption 3, we can calculate the loss function at the initial time:

$$\tilde{\ell}(U_{11}(0), u_{-1}(0)) = \frac{\sigma^4}{2} \operatorname{tr} [\Gamma \Lambda \Theta \Theta^\top \Lambda \Theta \Theta^\top] - \sigma^2 \operatorname{tr} [\Lambda^2 \Theta \Theta^\top].$$

From the property of trace, we know

$$\operatorname{tr} [\Lambda^2 \Theta \Theta^\top] = \operatorname{tr} [\Lambda \Theta \Theta^\top \Lambda^\top] = \|\Lambda \Theta\|_F^2.$$

From Von-Neumann's trace inequality (Lemma 31) and the fact that $\|\Theta \Theta^\top\|_F = 1$, we know

$$\begin{aligned} \operatorname{tr} [\Gamma \Lambda \Theta \Theta^\top \Lambda \Theta \Theta^\top] &\leq \sqrt{d} \left\| \Lambda \Theta \Theta^\top \Lambda \Theta \Theta^\top \right\|_F \|\Gamma\|_{op} \\ &\leq \sqrt{d} \|\Lambda \Theta\|_F^2 \left\| \Theta \Theta^\top \right\|_F \|\Gamma\|_{op} \\ &= \sqrt{d} \|\Lambda \Theta\|_F^2 \|\Gamma\|_{op}. \end{aligned}$$

Therefore, we have

$$\begin{aligned}\tilde{\ell}(U_{11}(0), u_{-1}(0)) &\leq \frac{\sqrt{d}\sigma^4}{2} \|\Lambda\Theta\|_F^2 \|\Gamma\|_{op} - \sigma^2 \|\Lambda\Theta\|_F^2 \\ &= \frac{\sigma^2}{2} \|\Lambda\Theta\|_F^2 \left[\sqrt{d}\sigma^2 \|\Gamma\|_{op} - 2 \right].\end{aligned}$$

From Assumption 3, we know $\|\Lambda\Theta\|_F \neq 0$. From (37), we know $\|\Gamma\|_{op} > 0$. Therefore, when

$$0 < \sigma < \sqrt{\frac{2}{\sqrt{d}\|\Gamma\|_{op}}},$$

we have

$$\tilde{\ell}(U_{11}(0), u_{-1}(0)) < 0. \quad \blacksquare$$

From the lemma above, we can actually further prove that the $u_{-1}(t)$ can be lower bounded by a positive constant for any $t \geq 0$. This will be a critical property to prove the PL inequality. We have the following lemma.

Lemma 19 *Consider gradient flow over L in (26) with respect to u starting from an initial value satisfying Assumption 3 with initial scale $0 < \sigma < \sqrt{\frac{2}{\sqrt{d}\|\Gamma\|_{op}}}$. For any $t \geq 0$, it holds that*

$$u_{-1} \geq \sqrt{\frac{\sigma^2}{2\sqrt{d}\|\Lambda\|_{op}^2} \|\Lambda\Theta\|_F^2 \left[2 - \sqrt{d}\sigma^2 \|\Gamma\|_{op} \right]} > 0. \quad (44)$$

Proof We prove by contradiction. Suppose the claim does not hold. From Lemma 17, we know $u_{-1}^2 = \text{tr} [U_{11}(U_{11})^\top] = \|U_{11}\|_F^2$. From Lemma 18, we know $u_{-1} = \|U_{11}\|_F$. Recall the definition of loss function:

$$\tilde{\ell}(U_{11}, u_{-1}) = \text{tr} \left[\frac{1}{2} u_{-1}^2 \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - u_{-1} \Lambda^2 (U_{11})^\top \right].$$

Since $\Gamma \succeq 0, \Lambda \succeq 0$, and they commute, we know from Lemma 32 that $\Gamma \Lambda \succeq 0$. Again, since $U_{11} \Lambda (U_{11})^\top = \left(U_{11} \Lambda^{\frac{1}{2}} \right) \left(U_{11} \Lambda^{\frac{1}{2}} \right)^\top \succeq 0$, from Lemma 32 we have $\text{tr} \left[\frac{1}{2} u_{-1}^2 \Gamma \Lambda U_{11} \Lambda (U_{11})^\top \right] \geq 0$. So

$$\tilde{\ell}(U_{11}, u_{-1}) \geq -\text{tr} \left[u_{-1} \Lambda^2 (U_{11})^\top \right].$$

From Von-Neumann's trace inequality, we know for any $t \geq 0$,

$$-\text{tr} \left[u_{-1} \Lambda^2 (U_{11})^\top \right] \geq -\sqrt{d} u_{-1} \|\Lambda^2\|_{op} \|U_{11}\|_F = -\sqrt{d} u_{-1}^2 \|\Lambda\|_{op}^2.$$

Therefore, under our assumption that the claim does not hold, we have

$$\tilde{\ell}(U_{11}, u_{-1}) \geq -\sqrt{d} u_{-1}^2 \|\Lambda\|_{op}^2 > -\frac{\sigma^2}{2} \|\Lambda\Theta\|_F^2 \left[2 - \sqrt{d}\sigma^2 \|\Gamma\|_{op} \right] \geq \tilde{\ell}(U_{11}(0), u_{-1}(0)).$$

Here, the last inequality comes from the proof of Lemma 18. This contradicts the non-increasing property of the loss function in gradient flow. \blacksquare

Finally, let's prove the PL inequality and further, the global convergence of gradient flow on the loss function $\tilde{\ell}$. We recall the stated lemma from the main text.

Lemma 20 *Suppose the initialization of gradient flow satisfies Assumption 3 with initialization scale satisfying $\sigma^2 < \frac{2}{\sqrt{d}\|\Gamma\|_{op}}$ for $\Gamma = (1 + \frac{1}{N})\Lambda + \frac{\text{tr}(\Lambda)}{N}I_d$. If we define*

$$\mu := \frac{\sigma^2}{\sqrt{d}\|\Lambda\|_{op}^2 \text{tr}(\Gamma^{-1}\Lambda^{-1}) \text{tr}(\Lambda^{-1})} \|\Lambda\Theta\|_F^2 \left[2 - \sqrt{d}\sigma^2 \|\Gamma\|_{op}\right] > 0, \quad (30)$$

then gradient flow on $\tilde{\ell}$ with respect to U_{11} and u_{-1} satisfies, for any $t \geq 0$,

$$\begin{aligned} \left\| \nabla \tilde{\ell}(U_{11}(t), u_{-1}(t)) \right\|_2^2 &:= \left\| \frac{\partial \tilde{\ell}}{\partial U_{11}} \right\|_F^2 + \left| \frac{\partial \tilde{\ell}}{\partial u_{-1}} \right|^2 \\ &\geq \mu \left(\tilde{\ell}(U_{11}(t), u_{-1}(t)) - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) \right). \end{aligned}$$

Moreover, gradient flow converges to the global minimum of $\tilde{\ell}$, and U_{11} and u_{-1} satisfy

$$\lim_{t \rightarrow \infty} u_{-1}(t) = \|\Gamma^{-1}\|_F^{\frac{1}{2}} \quad \text{and} \quad \lim_{t \rightarrow \infty} U_{11}(t) = \|\Gamma^{-1}\|_F^{-\frac{1}{2}} \Gamma^{-1}.$$

Proof From the definition and Lemma 19, we have

$$\begin{aligned} &\|\nabla \ell(U_{11}, u_{-1})\|_2^2 \\ &\geq \left\| \frac{\partial \ell}{\partial U_{11}} \right\|_F^2 = \|u_{-1}^2 \Gamma \Lambda U_{11} \Lambda - u_{-1} \Lambda^2\|_F^2 \\ &= u_{-1}^2 \left\| \Gamma \Lambda^{\frac{1}{2}} \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \Lambda^{\frac{1}{2}} \right\|_F^2 \\ &\geq \frac{\sigma^2}{2\sqrt{d}\|\Lambda\|_{op}^2} \|\Lambda\Theta\|_F^2 \left[2 - \sqrt{d}\sigma^2 \|\Gamma\|_{op}\right] \left\| \Gamma \Lambda^{\frac{1}{2}} \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \Lambda^{\frac{1}{2}} \right\|_F^2. \quad (45) \end{aligned}$$

To see why the second line is true, recall that $u_{-1} \in \mathbb{R}$ and Γ and Λ commute. The last line comes from the lower bound of u_{-1} in Lemma 19. From Corollary 16, we know

$$\begin{aligned} \ell - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \ell(U_{11}, u_{-1}) &= \frac{1}{2} \text{tr} \left[\Gamma \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right)^\top \right] \\ &= \frac{1}{2} \left\| \Gamma^{\frac{1}{2}} \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \right\|_F^2. \end{aligned}$$

Therefore, we know that

$$\begin{aligned} &\ell - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \ell(U_{11}, u_{-1}) \\ &\leq \frac{1}{2} \left\| \Gamma \Lambda^{\frac{1}{2}} \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \Lambda^{\frac{1}{2}} \right\|_F^2 \cdot \left\| \Gamma^{-\frac{1}{2}} \Lambda^{-\frac{1}{2}} \right\|_F^2 \left\| \Lambda^{-\frac{1}{2}} \right\|_F^2 \end{aligned}$$

$$= \frac{1}{2} \left\| \Gamma \Lambda^{\frac{1}{2}} \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \Lambda^{\frac{1}{2}} \right\|_F^2 \cdot \text{tr}(\Gamma^{-1} \Lambda^{-1}) \text{tr}(\Lambda^{-1}) \quad (46)$$

We compare (45) and (46) to obtain that in order to make the PL condition hold, one needs to let

$$\mu := \frac{\sigma^2}{\sqrt{d} \|\Lambda\|_{op}^2 \text{tr}(\Gamma^{-1} \Lambda^{-1}) \text{tr}(\Lambda^{-1})} \|\Lambda \Theta\|_F^2 \left[2 - \sqrt{d} \sigma^2 \|\Gamma\|_{op} \right] > 0.$$

Once we set this μ , we get the PL inequality. The μ is positive due to the assumption for σ in the lemma.

From the dynamics of gradient flow and the PL condition, we know

$$\begin{aligned} \frac{d}{dt} \left(\tilde{\ell} - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) \right) &= \left\langle \frac{dU_{11}}{dt}, \frac{\partial \tilde{\ell}}{\partial U_{11}} \right\rangle + \left\langle \frac{du_{-1}}{dt}, \frac{\partial \tilde{\ell}}{\partial u_{-1}} \right\rangle \\ &= - \left\| \frac{dU_{11}}{dt} \right\|_F^2 - \left| \frac{du_{-1}}{dt} \right|^2 \\ &\leq -\mu \left(\tilde{\ell} - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) \right). \end{aligned}$$

Therefore, we have when $t \rightarrow \infty$,

$$\begin{aligned} 0 &\leq \tilde{\ell} - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) \\ &\leq \exp(-\mu t) \left[\tilde{\ell}(U_{11}(0), u_{-1}(0)) - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) \right] \rightarrow 0, \end{aligned}$$

which implies

$$\lim_{t \rightarrow \infty} \left[\tilde{\ell} - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) \right] = 0.$$

From Corollary 16, we know this is

$$\left\| \Gamma^{\frac{1}{2}} \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \right\|_F^2 \rightarrow 0.$$

Since Γ and Λ are non-singular and positive definite, and they commute, we know

$$\|u_{-1} U_{11} - \Gamma^{-1}\|_F^2 \leq \left\| \Gamma^{-\frac{1}{2}} \Lambda^{-\frac{1}{2}} \right\|_F^2 \left\| \Gamma^{\frac{1}{2}} \left(u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \right\|_F^2 \left\| \Lambda^{-\frac{1}{2}} \right\|_F^2 \rightarrow 0.$$

This implies $u_{-1} U_{11} - \Gamma^{-1} \rightarrow 0_{d \times d}$ entry-wise. Since $u_{-1} = \|U_{11}\|_F$, we know

$$u_{-1}^2 = \|u_{-1} U_{11}\|_F \rightarrow \|\Gamma^{-1}\|_F.$$

Therefore, we know

$$\lim_{t \rightarrow \infty} u_{-1}(t) = \|\Gamma^{-1}\|_F^{\frac{1}{2}} \text{ and } \lim_{t \rightarrow \infty} U_{11}(t) = \|\Gamma^{-1}\|_F^{-\frac{1}{2}} \Gamma^{-1}.$$

■

Appendix B. Proof of Theorem 5

In this section, we prove Theorem 5, which characterizes the excess risk of the prediction of a trained LSA layer with respect to the risk of best linear predictor, on a new task which is possibly non-linear. First, we restate the theorem.

Theorem 5 *Let \mathcal{D} be a distribution over $(x, y) \in \mathbb{R}^d \times \mathbb{R}$, whose marginal distribution on x is $\mathcal{D}_x = \mathcal{N}(0, \Lambda)$. Assume $\mathbb{E}_{\mathcal{D}}[y], \mathbb{E}_{\mathcal{D}}[xy], \mathbb{E}_{\mathcal{D}}[y^2xx^\top]$ exist and are finite. Assume the test prompt is of the form $P = (x_1, y_1, \dots, x_M, y_M, x_{\text{query}}, y_{\text{query}})$, where $(x_i, y_i), (x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$. Let f_{LSA}^* be the LSA model with parameters W_*^{PV} and W_*^{KQ} in (11), and \hat{y}_{query} is the prediction for x_{query} given the prompt. If we define*

$$a := \Lambda^{-1} \mathbb{E}_{(x,y) \sim \mathcal{D}} [xy], \quad \Sigma := \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[(xy - \mathbb{E}(xy))(xy - \mathbb{E}(xy))^\top \right], \quad (15)$$

then, for $\Gamma = \Lambda + \frac{1}{N} \Lambda + \frac{1}{N} \text{tr}(\Lambda) I_d$. we have,

$$\begin{aligned} \mathbb{E} (\hat{y}_{\text{query}} - y_{\text{query}})^2 &= \underbrace{\min_{w \in \mathbb{R}^d} \mathbb{E} (\langle w, x_{\text{query}} \rangle - y_{\text{query}})^2}_{\text{Error of best linear predictor}} \\ &+ \frac{1}{M} \text{tr} [\Sigma \Gamma^{-2} \Lambda] + \frac{1}{N^2} \left[\|a\|_{\Gamma^{-2} \Lambda^3}^2 + 2 \text{tr}(\Lambda) \|a\|_{\Gamma^{-2} \Lambda^2}^2 + \text{tr}(\Lambda)^2 \|a\|_{\Gamma^{-2} \Lambda}^2 \right], \end{aligned} \quad (16)$$

where the expectation is over $(x_i, y_i), (x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$.

Proof Unless otherwise specified, we use \mathbb{E} to denote the expectation over (x_i, y_i) and $(x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$. Since when $(x, y) \sim \mathcal{D}$, we assume $\mathbb{E}[x], \mathbb{E}[y], \mathbb{E}[xy], \mathbb{E}[xx^\top], \mathbb{E}[y^2xx^\top]$ exist, we know that $\mathbb{E} (\langle w, x_{\text{query}} \rangle - y_{\text{query}})^2$ exists for each $w \in \mathbb{R}^d$. We denote

$$a := \arg \min_{w \in \mathbb{R}^d} \mathbb{E} (\langle w, x_{\text{query}} \rangle - y_{\text{query}})^2$$

as the weight of the best linear approximator. Actually, if we denote the function inside the minimum above as $R(w)$, we can write it as

$$R(w) = w^\top \Lambda w - 2 \mathbb{E} \left(y_{\text{query}} \cdot x_{\text{query}}^\top \right) w + \mathbb{E} y_{\text{query}}^2.$$

Since the Hessian matrix $\frac{\partial^2}{\partial w \partial w^\top} R(w)$ is Λ , which is positive definite, we know that this function is strictly convex and hence, the global minimum can be achieved at the unique first-order stationary point. This is

$$a = \Lambda^{-1} \mathbb{E} (y_{\text{query}} \cdot x_{\text{query}}). \quad (47)$$

We also define a similar vector for ease of computation:

$$b = \Gamma^{-1} \mathbb{E} (y_{\text{query}} \cdot x_{\text{query}}). \quad (48)$$

Therefore, we can decompose the risk as

$$\begin{aligned}
 \mathbb{E}(\widehat{y}_{\text{query}} - y_{\text{query}})^2 &= \underbrace{\mathbb{E}(\langle a, x_{\text{query}} \rangle - y_{\text{query}})^2}_{\text{I}} + \underbrace{\mathbb{E}(\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle)^2}_{\text{II}} \\
 &+ \underbrace{\mathbb{E}(\langle b, x_{\text{query}} \rangle - \langle a, x_{\text{query}} \rangle)^2}_{\text{III}} + \underbrace{2\mathbb{E}(\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle)(\langle a, x_{\text{query}} \rangle - y_{\text{query}})}_{\text{IV}} \\
 &+ \underbrace{2\mathbb{E}(\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle)(\langle b, x_{\text{query}} \rangle - \langle a, x_{\text{query}} \rangle)}_{\text{V}} \\
 &+ \underbrace{2\mathbb{E}(\langle b, x_{\text{query}} \rangle - \langle a, x_{\text{query}} \rangle)(\langle a, x_{\text{query}} \rangle - y_{\text{query}})}_{\text{VI}}
 \end{aligned}$$

The term I is the first term on the right hand side of (16). So it suffices to calculate II to VI.

First, from the tower property of conditional expectation, we have

$$\begin{aligned}
 \text{V} &= 2\mathbb{E} \left[\mathbb{E} \left((\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle)(\langle b, x_{\text{query}} \rangle - \langle a, x_{\text{query}} \rangle) \middle| x_{\text{query}} \right) \right] \\
 &= 2\mathbb{E} \left[\mathbb{E} \left(\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle \middle| x_{\text{query}} \right) (\langle b, x_{\text{query}} \rangle - \langle a, x_{\text{query}} \rangle) \right] = 0,
 \end{aligned}$$

since

$$\mathbb{E} \left(\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle \middle| x_{\text{query}} \right) = \left(\mathbb{E} \frac{1}{M} \sum_{i=1}^M y_i \Gamma^{-1} x_i - b \right)^\top x_{\text{query}} = 0.$$

Similarly, for IV, we have

$$\begin{aligned}
 \text{IV} &= 2\mathbb{E}(\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle)(\langle a, x_{\text{query}} \rangle - y_{\text{query}}) \\
 &= 2\mathbb{E} \left[\mathbb{E} \left((\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle)(\langle a, x_{\text{query}} \rangle - y_{\text{query}}) \middle| x_{\text{query}}, y_{\text{query}} \right) \right] \\
 &= 2\mathbb{E} \left[\mathbb{E} \left(\widehat{y}_{\text{query}} - \langle b, x_{\text{query}} \rangle \middle| x_{\text{query}}, y_{\text{query}} \right) (\langle a, x_{\text{query}} \rangle - y_{\text{query}}) \right] \\
 &= 0.
 \end{aligned}$$

For VI, we have

$$\begin{aligned}
 \text{VI} &= 2\mathbb{E} \text{tr} \left[(b - a)(\langle a, x_{\text{query}} \rangle - y_{\text{query}}) x_{\text{query}}^\top \right] \\
 &= 2 \text{tr} \left[(b - a) a^\top \Lambda \right] - 2 \text{tr} \left[(b - a) \mathbb{E} \left(y_{\text{query}} x_{\text{query}}^\top \right) \right] = 0,
 \end{aligned}$$

where the last line comes from the definition of a . Therefore, all cross terms vanish and it suffices to consider II and III.

From the definition II is equal to

$$\begin{aligned}
 & \mathbb{E} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i - \mathbb{E}(y_{\text{query}} \cdot x_{\text{query}}) \right)^\top \Gamma^{-1} x_{\text{query}} x_{\text{query}}^\top \Gamma^{-1} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i - \mathbb{E}(y_{\text{query}} \cdot x_{\text{query}}) \right) \\
 &= \mathbb{E} \operatorname{tr} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i - \mathbb{E}(y_{\text{query}} \cdot x_{\text{query}}) \right) \left(\frac{1}{M} \sum_{i=1}^M y_i x_i - \mathbb{E}(y_{\text{query}} \cdot x_{\text{query}}) \right)^\top \Gamma^{-2} \Lambda \\
 & \quad \text{(property of trace and the fact that } \Gamma \text{ and } \Lambda \text{ commute)} \\
 &= \frac{1}{M^2} \sum_{i,j=1}^M \mathbb{E} \operatorname{tr} \left\{ (y_i x_i - \mathbb{E}(y_{\text{query}} \cdot x_{\text{query}})) (y_j x_j - \mathbb{E}(y_{\text{query}} \cdot x_{\text{query}}))^\top \Gamma^{-2} \Lambda \right\} \\
 &= \frac{1}{M} \mathbb{E} \operatorname{tr} \left\{ (y_1 x_1 - \mathbb{E}(y_{\text{query}} \cdot x_{\text{query}})) (y_1 x_1 - \mathbb{E}(y_{\text{query}} \cdot x_{\text{query}}))^\top \Gamma^{-2} \Lambda \right\} \\
 & \quad \text{(all cross terms vanish due to the independence of } x_i \text{)} \\
 &= \frac{1}{M} \operatorname{tr} [\Sigma \Gamma^{-2} \Lambda].
 \end{aligned}$$

The last line comes from the definition of Σ .

For III, we have

$$\begin{aligned}
 \text{III} &= \mathbb{E}(b - a)^\top x_{\text{query}} x_{\text{query}}^\top (b - a) = a^\top \Lambda (\Gamma^{-1} - \Lambda^{-1}) \Lambda (\Gamma^{-1} - \Lambda^{-1}) \Lambda a \\
 &= \operatorname{tr} \left[(I - \Gamma \Lambda^{-1})^2 \Gamma^{-2} \Lambda^3 a a^\top \right] \\
 & \quad \text{(property of trace and the fact that } \Gamma \text{ and } \Lambda \text{ commute)} \\
 &= \frac{1}{N^2} \operatorname{tr} \left[(I_d + \operatorname{tr}(\Lambda) \Lambda^{-1})^2 \Gamma^{-2} \Lambda^3 a a^\top \right] \\
 &= \frac{1}{N^2} \left[\operatorname{tr}(\Gamma^{-2} \Lambda^3 a a^\top) + 2 \operatorname{tr}(\Lambda) \operatorname{tr}(\Gamma^{-2} \Lambda^2 a a^\top) + \operatorname{tr}(\Lambda)^2 \operatorname{tr}(\Gamma^{-2} \Lambda a a^\top) \right].
 \end{aligned}$$

Combining all terms above, we conclude. \blacksquare

Appendix C. Proof of Theorem 8

The proof of Theorem 8 is very similar to that of Theorem 4. The first step is to explicitly write out the dynamical system. In order to do so, we notice that the Lemma 9 does not depend on the training data and data-generating distribution and hence, it still holds in the case of a random covariance matrix. Therefore, we know when we input the embedding matrix E_τ to the linear self-attention layer with parameter $\theta = (W^{KQ}, W^{PV})$, the prediction will be

$$\widehat{y}_{\text{query}}(E_\tau; \theta) = u^\top H_\tau u,$$

where the matrix H_τ is defined as,

$$H_\tau = \frac{1}{2} X_\tau \otimes \left(\frac{E_\tau E_\tau^\top}{N} \right) \in \mathbb{R}^{(d+1)^2 \times (d+1)^2}, \quad X_\tau = \begin{pmatrix} 0_{d \times d} & x_{\tau, \text{query}} \\ (x_{\tau, \text{query}})^\top & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}$$

and

$$u = \text{Vec}(U) \in \mathbb{R}^{(d+1)^2}, \quad U = \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)},$$

where $U_{11} = W_{11}^{KQ} \in \mathbb{R}^{d \times d}$, $u_{12} = w_{21}^{PV} \in \mathbb{R}^{d \times 1}$, $u_{21} = w_{21}^{KQ} \in \mathbb{R}^{d \times 1}$, $u_{-1} = w_{22}^{PV} \in \mathbb{R}$ correspond to particular components of W^{PV} and W^{KQ} , defined in (5).

C.1 Dynamical system

The next lemma gives the dynamical system when the covariance matrices in the prompts are i.i.d. sampled from some distribution. Notice that in the lemma below, we do not assume Λ_τ are almost surely diagonal. The case when the covariance matrices are diagonal can be viewed as a special case of the following lemma.

Lemma 21 *Consider gradient flow on (20) with respect to u starting from an initial value that satisfies Assumption 3. We assume the covariance matrices Λ_τ are sampled from some distribution with finite third moment and Λ_τ are positive definite almost surely. We denote*

$$u = \text{Vec}(U) := \text{Vec} \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix} \text{ and define}$$

$$\Gamma_\tau = \left(1 + \frac{1}{N}\right) \Lambda_\tau + \frac{1}{N} \text{tr}(\Lambda_\tau) I_d \in \mathbb{R}^{d \times d}.$$

Then the dynamics of U follows

$$\begin{aligned} \frac{d}{dt} U_{11}(t) &= -u_{-1}^2 \mathbb{E} [\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau] + u_{-1} \mathbb{E} [\Lambda_\tau^2] \\ \frac{d}{dt} u_{-1}(t) &= -u_{-1} \text{tr} \mathbb{E} [\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top] + \text{tr} \left(\mathbb{E} [\Lambda_\tau^2] (U_{11})^\top \right), \end{aligned} \tag{49}$$

and $u_{12}(t) = 0_d, u_{21}(t) = 0_d$ for all $t \geq 0$.

Proof This lemma is a natural corollary of Lemma 10. Notice that Lemma 10 holds for any fixed positive definite Λ_τ . So when Λ_τ is random, if we condition on Λ_τ , the dynamical system will be

$$\begin{aligned} \frac{d}{dt} U_{11}(t) &= -u_{-1}^2 [\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau] + u_{-1} [\Lambda_\tau^2] \\ \frac{d}{dt} u_{-1}(t) &= -u_{-1} \text{tr} [\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top] + \text{tr} \left([\Lambda_\tau^2] (U_{11})^\top \right), \end{aligned} \tag{50}$$

and $u_{12}(t) = 0_d, u_{21}(t) = 0_d$ for all $t \geq 0$. Then, we conclude by simply taking expectation over Λ_τ . \blacksquare

The lemma above gives the dynamical system with general random covariance matrix. When Λ_τ are diagonal almost surely, we can actually simplify the dynamical system above. In this case, we have the following corollary.

Corollary 22 *Under the assumptions of Lemma 21, we further assume the covariance matrix Λ_τ to be diagonal almost surely. We denote $u_{ij}(t) \in \mathbb{R}$ as the (i, j) -th entry of $U_{11}(t)$, and further denote*

$$\begin{aligned}\gamma_i &= \mathbb{E} \left[\frac{N+1}{N} \lambda_{\tau,i}^3 + \frac{1}{N} \lambda_{\tau,i}^2 \cdot \sum_{j=1}^d \lambda_{\tau,j} \right], \\ \xi_i &= \mathbb{E} [\lambda_{\tau,i}^2], \\ \zeta_{ij} &= \mathbb{E} \left[\frac{N+1}{N} \lambda_{\tau,i}^2 \lambda_{\tau,j} + \frac{1}{N} \lambda_{\tau,i} \lambda_{\tau,j} \cdot \sum_{k=1}^d \lambda_{\tau,k} \right]\end{aligned}\tag{51}$$

for $i, j \in [d]$, where the expectation is over the distribution of Λ_τ . Then, the dynamical system (49) is equivalent to

$$\begin{aligned}\frac{d}{dt} u_{ii}(t) &= -\gamma_i u_{-1}^2 u_{ii} + \xi_i u_{-1} \quad \forall i \in [d], \\ \frac{d}{dt} u_{ij}(t) &= -\zeta_{ij} u_{-1}^2 u_{ij} \quad \forall i \neq j \in [d], \\ \frac{d}{dt} u_{-1}(t) &= -\sum_{i=1}^d [\gamma_i u_{-1} u_{ii}^2] - \sum_{i \neq j} \zeta_{ij} u_{-1} u_{ij}^2 + \sum_{i=1}^d [\xi_i u_{ii}].\end{aligned}\tag{52}$$

Proof This is directly obtained by rewriting the equation for each entry of U_{11} and recalling the assumption that Λ_τ (and hence Γ_τ) is diagonal almost surely. \blacksquare

C.2 Loss function and global minima

As in the proof of Theorem 4, we can actually recover the loss function in the random covariance case, up to a constant.

Lemma 23 *The differential equations in (52) are equivalent to gradient flow on the loss function*

$$\begin{aligned}\ell_{\text{rdm}}(U_{11}, u_{-1}) &= \mathbb{E} \operatorname{tr} \left[\frac{1}{2} u_{-1}^2 \Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top - u_{-1} \Lambda_\tau^2 (U_{11})^\top \right] \\ &= \frac{1}{2} \sum_{i=1}^d [\gamma_i u_{-1}^2 u_{ii}^2] + \frac{1}{2} \sum_{i \neq j} \zeta_{ij} u_{-1}^2 u_{ij}^2 - \sum_{i=1}^d [\xi_i u_{ii} u_{-1}]\end{aligned}\tag{53}$$

with respect to $u_{ij} \forall i, j \in [d]$ and u_{-1} , from an initial value that satisfies Assumption 3.

Proof This can be verified by simply taking gradient of ℓ_{rdm} to show that

$$\frac{d}{dt} u_{ii} = -\frac{\partial \ell_{\text{rdm}}}{\partial u_{ii}} \quad \forall i \in [d], \quad \frac{d}{dt} u_{ij} = -\frac{\partial \ell_{\text{rdm}}}{\partial u_{ij}} \quad \forall i \neq j \in [d], \quad \frac{d}{dt} u_{-1} = -\frac{\partial \ell_{\text{rdm}}}{\partial u_{-1}}.$$

\blacksquare

Next, we solve for the minimum of ℓ_{rdm} and give the expression for all global minima.

Lemma 24 *Let ℓ_{rdm} be the loss function in (53). We denote*

$$\min \ell_{\text{rdm}} := \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \ell_{\text{rdm}}(U_{11}, u_{-1}).$$

Then, we have

$$\min \ell_{\text{rdm}} = -\frac{1}{2} \sum_{i=1}^d \frac{\xi_i^2}{\gamma_i} \quad (54)$$

and

$$\ell_{\text{rdm}}(U_{11}, u_{-1}) - \min \ell_{\text{rdm}} = \frac{1}{2} \sum_{i=1}^d \gamma_i \left(u_{ii} u_{-1} - \frac{\xi_i}{\gamma_i} \right)^2 + \frac{1}{2} \sum_{i \neq j} \zeta_{ij} u_{-1}^2 u_{ij}^2. \quad (55)$$

Moreover, denoting u_{ij} as the (i, j) -entry of U_{11} , all global minima of ℓ_{rdm} satisfy

$$u_{-1} \cdot u_{ij} = \mathbb{I}(i = j) \cdot \frac{\xi_i}{\gamma_i}. \quad (56)$$

Proof From the definition of ℓ_{rdm} , we have

$$\ell_{\text{rdm}} = \frac{1}{2} \sum_{i=1}^d \gamma_i \left(u_{ii} u_{-1} - \frac{\xi_i}{\gamma_i} \right)^2 + \frac{1}{2} \sum_{i \neq j} \zeta_{ij} u_{-1}^2 u_{ij}^2 - \frac{1}{2} \sum_{i=1}^d \frac{\xi_i^2}{\gamma_i} \geq -\frac{1}{2} \sum_{i=1}^d \frac{\xi_i^2}{\gamma_i}.$$

The equation holds when $u_{ij} = 0$ for $i \neq j \in [d]$ and $u_{-1} u_{ii} = \frac{\xi_i}{\gamma_i}$ for each $i \in [d]$. This can be achieved by simply letting $u_{-1} = 1$ and $u_{ii} = \frac{\xi_i}{\gamma_i}$ for $i \in [d]$. Of course, when we replace (u_{-1}, u_{ii}) with $(c u_{-1}, c^{-1} u_{ii})$ for any constant $c \neq 0$, we can also achieve this global minimum. \blacksquare

C.3 PL Inequality and global convergence

Finally, to end the proof, we prove a Polyak-Lojasiewicz Inequality on the loss function ℓ_{rdm} , and then prove global convergence. Before that, let's first prove the balanced condition of parameters will hold during the whole trajectory.

Lemma 25 (Balanced condition) *Under the assumptions of Lemma 21, for any $t \geq 0$, it holds that*

$$u_{-1}^2 = \text{tr} \left[U_{11} (U_{11})^\top \right]. \quad (57)$$

Proof The proof is similar to the proof of Lemma 17. From Lemma 10, we multiply the first equation in (49) by $(U_{11})^\top$ from the right to get

$$\left[\frac{d}{dt} U_{11}(t) \right] (U_{11})^\top = -u_{-1}^2 \mathbb{E} \left[\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top \right] + u_{-1} \mathbb{E} \left[\Lambda_\tau^2 (U_{11})^\top \right].$$

Also we multiply the second equation in Lemma 49 by u_{-1} to obtain

$$\left(\frac{d}{dt}u_{-1}(t)\right)u_{-1}(t) = -u_{-1}^2 \operatorname{tr} \mathbb{E} \left[\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top \right] + u_{-1} \operatorname{tr} \left(\mathbb{E} [\Lambda_\tau^2] (U_{11})^\top \right),$$

Therefore, we have

$$\operatorname{tr} \left[\left(\frac{d}{dt} U_{11}(t) \right) (U_{11}(t))^\top \right] = \left(\frac{d}{dt} u_{-1}(t) \right) u_{-1}(t).$$

Taking the transpose of the equation above and adding to itself gives

$$\frac{d}{dt} \operatorname{tr} \left[U_{11}(t) (U_{11}(t))^\top \right] = \frac{d}{dt} (u_{-1}(t)^2).$$

Notice that from Assumption 3, we know that

$$u_{-1}(0)^2 = \sigma^2 = \sigma^2 \operatorname{tr} \left[\Theta \Theta^\top \Theta \Theta^\top \right] = \operatorname{tr} \left[U_{11}(0) (U_{11}(0))^\top \right].$$

So for any time $t \geq 0$, the equation holds. ■

Next, similar to the proof of Theorem 4, we prove that, as long as the initial scale is small enough, u_{-1} will be positive along the whole trajectory and can be lower bounded by a positive constant, which implies that the trajectories will be away from the saddle point at the origin.

Lemma 26 *We do gradient flow on ℓ_{rdm} with respect to $u_{i,j}$ ($\forall i, j \in [d]$) and u_{-1} . Suppose the initialization satisfies Assumption 3 with initial scale*

$$0 < \sigma < \sqrt{\frac{2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2}{\sqrt{d} \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right]}}, \quad (58)$$

then for any $t \geq 0$, it holds that

$$u_{-1}(t) > 0. \quad (59)$$

Proof From the dynamics of gradient flow, we know the loss function ℓ_{rdm} is non-increasing:

$$\frac{d\ell_{\text{rdm}}}{dt} = \sum_{i,j=1}^d \frac{\partial \ell_{\text{rdm}}}{\partial u_{ij}} \cdot \frac{du_{ij}}{dt} + \frac{\partial \ell_{\text{rdm}}}{\partial u_{-1}} \cdot \frac{du_{-1}}{dt} = - \sum_{i,j=1}^d \left[\frac{\partial \ell_{\text{rdm}}}{\partial u_{ij}} \right]^2 - \left[\frac{\partial \ell_{\text{rdm}}}{\partial u_{-1}} \right]^2 \leq 0.$$

Since we assume $U_{11}(0) = \Theta \Theta^\top$, we know the loss function at $t = 0$ is

$$\ell_{\text{rdm}}(U_{11}(0), u_{-1}(0)) = \mathbb{E} \operatorname{tr} \left[\frac{\sigma^4}{2} \Gamma_\tau \Lambda_\tau \Theta \Theta^\top \Lambda_\tau \Theta \Theta^\top - \sigma^2 \Lambda_\tau^2 \Theta \Theta^\top \right].$$

From the property of trace, we know

$$\mathbb{E} \operatorname{tr} \left[\sigma^2 \Lambda_\tau^2 \Theta \Theta^\top \right] = \sigma^2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2.$$

From Von-Neumann's trace inequality and the assumption that $\|\Theta\Theta^\top\|_F = 1$, we know

$$\begin{aligned} & \mathbb{E} \operatorname{tr} \left[\frac{\sigma^4}{2} \Gamma_\tau \Lambda_\tau \Theta \Theta^\top \Lambda_\tau \Theta \Theta^\top \right] \\ & \leq \frac{\sigma^4 \sqrt{d}}{2} \mathbb{E} \|\Gamma_\tau\|_{op} \left\| \Lambda_\tau \Theta \Theta^\top \Lambda_\tau \Theta \Theta^\top \right\|_F \\ & \leq \frac{\sigma^4 \sqrt{d}}{2} \frac{\|\Theta\Theta^\top\|_F^2}{\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2} \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right] = \frac{\sigma^4 \sqrt{d}}{2} \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right]. \end{aligned}$$

From the assumptions on Θ and Λ_τ we know $\mathbb{E} \Lambda_\tau \Theta \neq 0_{d \times d}$ and $\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 > 0$. Therefore, comparing the two displays above, we know when (58) holds, we must have $\ell_{\text{rdm}}(0) < 0$. So from the non-increasing property of the loss function, we know $\ell_{\text{rdm}}(t) < 0$ for any time $t \geq 0$. Notice that when $u_{-1} = 0$, the loss function is also zero, which suggests that $u_{-1}(t) \neq 0$ for any time $t \geq 0$. Since $u_{-1}(0) > 0$ and the trajectory of u_{-1} must be continuous, we know that it stays positive at all times. \blacksquare

Lemma 27 *We do gradient flow on ℓ_{rdm} with respect to $u_{i,j}$ ($\forall i, j \in [d]$) and u_{-1} . Suppose the initialization satisfies Assumption 3 and the initial scale satisfies (58). Then, for any $t \geq 0$, it holds that*

$$u_{-1}(t) \geq \sqrt{\frac{\sigma^2}{2\sqrt{d} \|\mathbb{E} \Lambda_\tau^2\|_{op}} \left[2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2 - \sqrt{d} \sigma^2 \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right] \right]} > 0. \quad (60)$$

Proof From the dynamics of gradient flow, we know ℓ_{rdm} is non-increasing (see the proof of Lemma 26). Recall the definition of the loss function:

$$\ell_{\text{rdm}}(U_{11}, u_{-1}) = \mathbb{E} \operatorname{tr} \left[\frac{1}{2} u_{-1}^2 \Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top - u_{-1} \Lambda_\tau^2 (U_{11})^\top \right].$$

Since Λ_τ commutes with Γ_τ and they are both positive definite almost surely, we know that $\Gamma_\tau \Lambda_\tau \succeq 0_{d \times d}$ almost surely from Lemma 29. Again, since $U_{11} \Lambda_\tau (U_{11})^\top \succeq 0_{d \times d}$ almost surely, from Lemma 29 we have $\operatorname{tr} \left[\frac{1}{2} u_{-1}^2 \Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top \right] \geq 0$ almost surely. Therefore, we have

$$\ell_{\text{rdm}}(U_{11}, u_{-1}) \geq -\mathbb{E} \operatorname{tr} \left[u_{-1} \Lambda_\tau^2 (U_{11})^\top \right] = -\operatorname{tr} \left[u_{-1} (\mathbb{E} \Lambda_\tau^2) (U_{11})^\top \right].$$

From Von Neumann's trace inequality (Lemma 31) and the fact that $u_{-1}(t) > 0$ for any $t \geq 0$ (Lemma 26), we know $\ell_{\text{rdm}}(U_{11}(t), u_{-1}(t)) \geq -\sqrt{d} u_{-1} \|\mathbb{E} \Lambda_\tau^2\|_{op} \|U_{11}\|_F$. From Lemma 25, we know $u_{-1}^2 = \operatorname{tr}(U_{11}(U_{11})^\top) = \|U_{11}\|_F^2$. Since $u_{-1}(t) > 0$ for any time, we know actually $u_{-1}(t) = \|U_{11}(t)\|_F$. So we have

$$\ell_{\text{rdm}}(U_{11}(t), u_{-1}(t)) \geq -\sqrt{d} u_{-1}(t)^2 \|\mathbb{E} \Lambda_\tau^2\|_{op}.$$

From the proof of Lemma 26, we know

$$\ell_{\text{rdm}}(U_{11}(t), u_{-1}(t)) \leq \ell_{\text{rdm}}(U_{11}(0), u_{-1}(0)) \leq \frac{\sigma^4 \sqrt{d}}{2} \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right] - \sigma^2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2.$$

Combine the two preceding displays above, we have

$$u_{-1}(t) \geq \sqrt{\frac{\sigma^2}{2\sqrt{d} \|\mathbb{E} \Lambda_\tau^2\|_{op}} \left[2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2 - \sqrt{d} \sigma^2 \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right] \right]} > 0.$$

The last inequality comes from Lemma 26. ■

Finally, we prove the PL Inequality, which naturally leads to the global convergence.

Lemma 28 *We do gradient flow on ℓ_{rdm} with respect to $u_{i,j}$ ($\forall i, j \in [d]$) and u_{-1} . Suppose the initialization satisfies Assumption 3 and the initial scale satisfies (58). If we denote*

$$\eta = \min \{ \gamma_i, i \in [d]; \zeta_{ij}, i \neq j \in [d] \}$$

and

$$\nu := \frac{\eta \cdot \sigma^2}{2\sqrt{d} \|\mathbb{E} \Lambda_\tau^2\|_{op}} \left[2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2 - \sqrt{d} \sigma^2 \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right] \right] > 0, \quad (61)$$

then for any $t \geq 0$, it holds that

$$\|\nabla \ell_{\text{rdm}}(U_{11}, u_{-1})\|_2^2 := \sum_{i,j=1}^d \left| \frac{\partial \ell_{\text{rdm}}}{\partial u_{ij}} \right|^2 + \left| \frac{\partial \ell_{\text{rdm}}}{\partial u_{-1}} \right|^2 \geq \nu (\ell_{\text{rdm}} - \min \ell_{\text{rdm}}). \quad (62)$$

Additionally, ℓ_{rdm} converges to the global minimal value, u_{ij} and u_{-1} converge to the following limits,

$$\lim_{t \rightarrow \infty} u_{ij}(t) = \mathbb{I}(i = j) \cdot \left[\sum_{i=1}^d \frac{\xi_i^2}{\gamma_i^2} \right]^{-\frac{1}{4}} \cdot \frac{\xi_i}{\gamma_i} \quad \forall i \in [d], \quad \lim_{t \rightarrow \infty} u_{-1}(t) = \left[\sum_{i=1}^d \frac{\xi_i}{\gamma_i} \right]^{\frac{1}{4}}. \quad (63)$$

Translating back to the original parameterization, we have this is equivalent to

$$\begin{aligned} \lim_{t \rightarrow \infty} W^{KQ}(t) &= \begin{pmatrix} \left\| [\mathbb{E} \Gamma_\tau \Lambda_\tau^2]^{-1} \mathbb{E} [\Lambda_\tau^2] \right\|_F^{-\frac{1}{2}} \cdot [\mathbb{E} \Gamma_\tau \Lambda_\tau^2]^{-1} \mathbb{E} [\Lambda_\tau^2] & 0_d \\ 0_d^\top & 0 \end{pmatrix}, \\ \lim_{t \rightarrow \infty} W^{PV}(t) &= \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & \left\| [\mathbb{E} \Gamma_\tau \Lambda_\tau^2]^{-1} \mathbb{E} [\Lambda_\tau^2] \right\|_F^{\frac{1}{2}} \end{pmatrix}, \end{aligned}$$

where $\Gamma_\tau = \frac{N+1}{N} \Lambda_\tau + \frac{1}{N} \text{tr}(\Lambda_\tau) I_d \in \mathbb{R}^{d \times d}$ and \mathbb{E} is over Λ_τ .

Proof First, we prove the PL Inequality. From Lemma 24, we know

$$\ell_{\text{rdm}}(U_{11}, u_{-1}) - \min \ell_{\text{rdm}} = \frac{1}{2} \sum_{i=1}^d \gamma_i \left(u_{ii} u_{-1} - \frac{\xi_i}{\gamma_i} \right)^2 + \frac{1}{2} \sum_{i \neq j} \zeta_{ij} u_{-1}^2 u_{ij}^2,$$

where $\xi_i, \zeta_{ij}, \gamma_i$ are defined in (51). Meanwhile, we calculate the square norm of the gradient of ℓ_{rdm} :

$$\begin{aligned} \|\nabla \ell_{\text{rdm}}(U_{11}, u_{-1})\|_2^2 &:= \sum_{i,j=1}^d \left| \frac{\partial \ell_{\text{rdm}}}{\partial u_{ij}} \right|^2 + \left| \frac{\partial \ell_{\text{rdm}}}{\partial u_{-1}} \right|^2 \geq \sum_{i,j=1}^d \left| \frac{\partial \ell_{\text{rdm}}}{\partial u_{ij}} \right|^2 \\ &= \sum_{i=1}^d \gamma_i^2 u_{-1}^2 \left(u_{ii} u_{-1} - \frac{\xi_i}{\gamma_i} \right)^2 + \sum_{i \neq j} \zeta_{ij}^2 u_{-1}^4 u_{ij}^2. \end{aligned}$$

Comparing the two displays above, we know that in order to ensure that $\|\nabla \ell_{\text{rdm}}\|_2^2 \geq \nu (\ell_{\text{rdm}} - \min \ell_{\text{rdm}})$, it suffices to make

$$\begin{aligned} \gamma_i u_{-1}(t)^2 &\geq \frac{\nu}{2} \quad \forall i \in [d], \\ \zeta_{ij} u_{-1}(t)^2 &\geq \frac{\nu}{2} \quad \forall i \neq j \in [d]. \end{aligned}$$

We define $\eta := \min \{\gamma_i, \zeta_{ij}, i \neq j \in [d]\}$, then it is sufficient to make

$$\eta u_{-1}(t)^2 \geq \frac{\nu}{2}.$$

From Lemma 27, we know that we can actually lower bound u_{-1} from below by a positive constant. Then, the inequality holds if we take

$$\nu := \frac{\eta \cdot \sigma^2}{2\sqrt{d} \|\mathbb{E} \Lambda_\tau^2\|_{op}} \left[2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2 - \sqrt{d} \sigma^2 \left[\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right] \right] > 0.$$

Therefore, as long as we take ν as above, a PL inequality holds for ℓ_{rdm} .

With an abuse of notation, let us write $\ell_{\text{rdm}}(t) = \ell_{\text{rdm}}(U_{11}(t), u_{-1}(t))$. Then, from the dynamics of gradient flow and the PL Inequality ((62)), we know

$$\frac{d}{dt} [\ell_{\text{rdm}}(t) - \min \ell_{\text{rdm}}] = -\|\nabla \ell_{\text{rdm}}(t)\|_2^2 \leq -\nu (\ell_{\text{rdm}}(t) - \min \ell_{\text{rdm}}),$$

which by Grönwall's inequality implies

$$0 \leq \ell_{\text{rdm}}(t) - \min \ell_{\text{rdm}} \leq \exp(-\nu t) [\ell_{\text{rdm}}(0) - \min \ell_{\text{rdm}}] \rightarrow 0$$

when $t \rightarrow \infty$. From Lemma 24, we know

$$\sum_{i=1}^d \gamma_i \left(u_{ii} u_{-1} - \frac{\xi_i}{\gamma_i} \right)^2 + \sum_{i \neq j} \zeta_{ij} u_{-1}^2 u_{ij}^2 \rightarrow 0 \text{ when } t \rightarrow \infty.$$

This implies

$$\begin{aligned} u_{ii}u_{-1} &\rightarrow \frac{\xi_i}{\gamma_i} \quad \forall i \in [d], \\ u_{ij}u_{-1} &\rightarrow 0 \quad \forall i \neq j \in [d]. \end{aligned} \tag{64}$$

We take square of $u_{ii}(t)u_{-1}(t)$ and $u_{ij}(t)u_{-1}(t)$, then sum over all $i, j \in [d]$. Then, we get $u_{-1}^2 \sum_{i,j=1}^d u_{ij}^2 \rightarrow \sum_{i=1}^d \frac{\xi_i^2}{\gamma_i^2}$. From Lemma 25, we know for any $t \geq 0$, $u_{-1}(t)^2 = \text{tr}(U_{11}(U_{11})^\top) = \sum_{i,j=1}^d u_{ij}^2$. So we have

$$u_{-1}(t)^4 = u_{-1}^2 \sum_{i,j=1}^d u_{ij}^2 \rightarrow \sum_{i=1}^d \frac{\xi_i^2}{\gamma_i^2},$$

which implies

$$u_{-1}(t) \rightarrow \left[\sum_{i=1}^d \frac{\xi_i^2}{\gamma_i^2} \right]^{\frac{1}{4}} \tag{65}$$

when $t \rightarrow \infty$. Combining (64) and (65), we conclude

$$u_{ij}(t) \rightarrow 0 \quad \forall i \neq j \in [d], \quad u_{ii}(t) \rightarrow \left[\sum_{i=1}^d \frac{\xi_i^2}{\gamma_i^2} \right]^{-\frac{1}{4}} \cdot \frac{\xi_i}{\gamma_i} \quad \forall i \in [d].$$

■

Appendix D. Technical lemmas

Lemma 29 (Petersen and Pedersen, 2008) *We denote $\mathbf{A}, \mathbf{B}, \mathbf{X}$ as matrices and \mathbf{x} as vectors. Then, we have*

- $\frac{\partial \mathbf{x}^\top \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^\top) \mathbf{x}$.
- $\text{Vec}(\mathbf{A} \mathbf{X} \mathbf{B}) = (\mathbf{B}^\top \otimes \mathbf{A}) \text{Vec}(\mathbf{X})$.
- $\text{tr}(\mathbf{A}^\top \mathbf{B}) = \text{Vec}(\mathbf{A})^\top \text{Vec}(\mathbf{B})$.
- $\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{X} \mathbf{B} \mathbf{X}^\top) = \mathbf{X} \mathbf{B}^\top + \mathbf{X} \mathbf{B}$.
- $\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{A} \mathbf{X}^\top) = \mathbf{A}$.
- $\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^\top \mathbf{C}) = \mathbf{A}^\top \mathbf{C}^\top \mathbf{X} \mathbf{B}^\top + \mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B}$.

Lemma 30 *If X is Gaussian random vector of d dimension, mean zero and covariance matrix Λ , and $A \in \mathbb{R}^{d \times d}$ is a fixed matrix. Then*

$$\mathbb{E} \left[X X^\top A X X^\top \right] = \Lambda \left(A + A^\top \right) \Lambda + \text{tr}(A \Lambda) \Lambda.$$

Proof We denote $X = (X_1, \dots, X_d)^\top$. Then,

$$XX^\top AXX^\top = X(X^\top AX)X^\top = \left(\sum_{i,j=1}^d A_{ij} X_i X_j \right) XX^\top.$$

So we know $(XX^\top AXX^\top)_{k,l} = \left(\sum_{i,j=1}^d A_{ij} X_i X_j \right) X_k X_l$. From Isserlis' Theorem in probability theory (Theorem 1.1 in Michalowicz et al. (2009), originally proposed in Wick (1950)), we know for any $i, j, k, l \in [d]$, it holds that

$$\mathbb{E}[X_i X_j X_k X_l] = \Lambda_{ij} \Lambda_{kl} + \Lambda_{ik} \Lambda_{jl} + \Lambda_{il} \Lambda_{jk}.$$

Then, we have for any fixed $k, l \in [d]$,

$$\begin{aligned} \mathbb{E}(XX^\top AXX^\top)_{k,l} &= \sum_{i,j=1}^d A_{ij} \Lambda_{ij} \Lambda_{kl} + A_{ij} \Lambda_{ik} \Lambda_{jl} + A_{ij} \Lambda_{il} \Lambda_{jk} \\ &= \text{tr}(A\Lambda) \Lambda_{kl} + \Lambda_k^\top (A + A^\top) \Lambda_l. \end{aligned}$$

Therefore, we know

$$\mathbb{E}(XX^\top AXX^\top) = \Lambda (A + A^\top) \Lambda + \text{tr}(A\Lambda) \Lambda. \quad \blacksquare$$

Lemma 31 (Von-Neumann's Trace Inequality) *Let $U, V \in \mathbb{R}^{d \times n}$ with $d \leq n$. We have*

$$\text{tr}(U^\top V) \leq \sum_{i=1}^d \sigma_i(U) \sigma_i(V) \leq \|U\|_{\text{op}} \times \sum_{i=1}^d \sigma_i(V) \leq \sqrt{d} \cdot \|U\|_{\text{op}} \|V\|_F$$

where $\sigma_1(X) \geq \sigma_2(X) \geq \dots \geq \sigma_d(X)$ are the ordered singular values of $X \in \mathbb{R}^{d \times n}$.

Lemma 32 ((Meenakshi and Rajian, 1999)) *For any two positive semi-definitive matrices $A, B \in \mathbb{R}^{d \times d}$, we have*

- $\text{tr}[AB] \geq 0$.
- $AB \succeq 0$ if and only if A and B commute.

Appendix E. Experiment details

In this section, we provide more details for the experiment in Figure 1. Our experimental setup is based on the codebase provided by Garg et al. (2022), with a modification that allows for the possibility that the covariate distribution changes across prompts. We use the standard GPT2 architecture with embedding size 256, 12 layers and 8 heads (Radford et al., 2018) as implemented by HuggingFace (Wolf et al., 2020). For the GPT2 models, we use the embedding method proposed by Garg et al. (2022), where instead of concatenating x and y into a single token, they are treated as separate tokens. It is also worth noting that the training objective function for the GPT2 model is different than those we consider for the linear self-attention network: for the GPT2 model, the objective function is the average over the full length of the context sequence (predictions for each x_i using $(x_k, y_k)_{k < i}$), while in our setting the objective function is only for the final query point. However, in the figure, for both GPT2 and the linear self-attention model the error plotted corresponds to the error for predicting the final query point.

In all experiments, covariates are sampled from a mean-zero Gaussian in $d = 20$ dimensions with either fixed or random covariance matrix. For the fixed covariance case, we fix the covariance matrix to be identity; for the random case, the covariance matrices are restricted to be diagonal and all diagonal entries are i.i.d. sampled from the standard exponential distribution. The linear weights in all tasks are i.i.d. sampled from standard Gaussian distribution and also independently from all covariates. We trained the model for 500000 steps using Adam (Kingma and Ba, 2014) with a batch size of 64 and learning rate of 0.0001. We use the same curriculum strategy of Garg et al. (2022) for acceleration.

For testing the trained model, we used ordinary least squares as a baseline which is optimal for noiseless linear regression tasks. For prompts at test time, covariates are sampled i.i.d. from a mean-zero Gaussian distribution. For the fixed-covariance evaluation, the covariance is the identity matrix. In the random-covariance evaluation, the covariance is a random diagonal matrix with diagonal entries sampled from the standard exponential distribution, multiplied by a scaling coefficient $c \in \{1, 4, 9\}$, i.e. for each task τ , the covariance matrix in the random case is

$$\Lambda_\tau = c \cdot \text{diag}(\lambda_{\tau,1}, \dots, \lambda_{\tau,d})$$

where $\lambda_{\tau,i} \stackrel{\text{i.i.d.}}{\sim} \text{Exponential}(1)$ for any τ and $i \in [d]$. The plots in Figure 1 show the error averaged over 64^2 prompts, where we sample 64 covariance matrices for each curve and 64 prompts for each covariance matrix. We compute 90% confidence intervals over 1000 bootstrap trials for each test.

References

- Jacob Abernethy, Alekh Agarwal, Teodor V. Marinov, and Manfred K. Warmuth. A mechanism for sample-efficient in-context learning for sparse retrieval tasks. *Preprint, arXiv:2305.17040*, 2023.
- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Preprint, arXiv:2306.00297*, 2023.

- Kabir Ahuja, Madhur Panwar, and Navin Goyal. In-context learning through the Bayesian prism. *arXiv preprint arXiv:2306.04891*, 2023.
- Kartik Ahuja and David Lopez-Paz. A closer look at in-context learning under distribution shifts. *Preprint, arXiv:2305.16704*, 2023.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? Investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Cem Anil, Yuhuai Wu, Anders Johan Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Venkatesh Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *International Conference on Machine Learning*, pages 244–253, 2018.
- Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Shahar Azulay, Edward Moroshko, Mor Shpigel Nacson, Blake E Woodworth, Nathan Srebro, Amir Globerson, and Daniel Soudry. On the implicit bias of initialization shape: Beyond infinitesimal mirror descent. In *International Conference on Machine Learning*, pages 468–477, 2021.
- Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *Preprint, arXiv:2306.04637*, 2023.
- Mohamed Ali Belabbas. On implicit regularization: Morse functions and applications to matrix factorization. *arXiv preprint arXiv:2001.04264*, 2020.
- Satwik Bhattamishra, Arkil Patel, and Navin Goyal. On the computational power of transformers and its implications in sequence modeling. *arXiv preprint arXiv:2006.09286*, 2020.
- Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Zhifang Sui, and Furu Wei. Why can GPT learn in-context? Language models secretly perform gradient descent as meta optimizers. *arXiv preprint arXiv:2212.10559*, 2022.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Association for Computational Linguistics (ACL)*, 2019.

- Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal transformers. *arXiv preprint arXiv:1807.03819*, 2019.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. *Advances in Neural Information Processing Systems*, 31, 2018.
- Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, 2022.
- Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn in-context? A case study of simple function classes. *arXiv preprint arXiv:2208.01066*, 2022.
- Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. *Advances in Neural Information Processing Systems*, 30, 2017.
- Chi Han, Ziqi Wang, Han Zhao, and Heng Ji. In-context learning of large language models explained as kernel regression, 2023.
- Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. *Advances in Neural Information Processing Systems*, 35:37822–37836, 2022.
- Jikai Jin, Zhiyuan Li, Kaifeng Lyu, Simon S Du, and Jason D Lee. Understanding incremental learning of gradient descent: A fine-grained analysis of matrix sensing. *arXiv preprint arXiv:2301.11500*, 2023.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Shuai Li, Zhao Song, Yu Xia, Tong Yu, and Tianyi Zhou. The closeness of in-context learning and weight shifting for softmax regression. *arXiv preprint arXiv:2304.13276*, 2023a.
- Yingcong Li, M Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. *arXiv preprint arXiv:2301.07067*, 2023b.
- Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47, 2018.

- Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do transformers learn topic structure: Towards a mechanistic understanding. *arXiv preprint arXiv:2303.04245*, 2023c.
- Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. *arXiv preprint arXiv:2012.09839*, 2020.
- Valerii Likhoshesterov, Krzysztof Choromanski, and Adrian Weller. On the expressive power of self-attention matrices. *arXiv preprint arXiv:2106.03764*, 2021.
- Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers learn shortcuts to automata. In *International Conference on Learning Representations (ICLR)*, 2023.
- AR Meenakshi and C Rajian. On a product of positive semidefinite matrices. *Linear algebra and its applications*, 295(1-3):3–6, 1999.
- JV Michalowicz, JM Nichols, F Bucholtz, and CC Olson. An Isserlis’ theorem for mixed Gaussian variables: Application to the auto-bispectral density. *Journal of Statistical Physics*, 136:89–102, 2009.
- Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Jorge Pérez, Javier Marinković, and Pablo Barceló. On the Turing completeness of modern neural network architectures. *arXiv preprint arXiv:1901.03429*, 2019.
- Kaare Brandt Petersen and Michael Syskind Pedersen. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. Technical report, OpenAI, 2018. URL https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. Technical report, OpenAI, 2019. URL https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf.
- Mahdi Soltanolkotabi, Dominik Stöger, and Changzhi Xie. Implicit balancing and regularization: Generalization and convergence guarantees for overparameterized asymmetric matrix sensing. *arXiv preprint arXiv:2303.14244*, 2023.
- Asher Trockman and J Zico Kolter. Mimetic initialization of self-attention layers. *arXiv preprint arXiv:2305.09828*, 2023.

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 2017.
- Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. *arXiv preprint arXiv:2212.07677*, 2022.
- Xinyi Wang, Wanrong Zhu, and William Yang Wang. Large language models are implicitly topic models: Explaining and finding good demonstrations for in-context learning. *arXiv preprint arXiv:2301.11916*, 2023.
- Gian-Carlo Wick. The evaluation of the collision matrix. *Physical review*, 80(2):268, 1950.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, and Morgan Funtowicz. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45, 2020.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit Bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions? *arXiv preprint arXiv:1912.10077*, 2019.
- Chulhee Yun, Yin-Wen Chang, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. $O(n)$ connections are expressive enough: Universal approximability of sparse transformers. *Advances in Neural Information Processing Systems*, 33:13783–13794, 2020.
- Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does in-context learning learn? Bayesian model averaging, parameterization, and generalization. *Preprint, arXiv:2305.19420*, 2023.