

---

# Trained Transformers Learn Linear Models In-Context

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Attention-based neural network sequence models such as transformers have the  
2 capacity to act as supervised learning algorithms: They can take as input a sequence  
3 of labeled examples and output predictions for unlabeled test examples. Indeed,  
4 recent work by Garg et al. has shown that when training GPT2 architectures  
5 over random instances of linear regression problems, these models' predictions  
6 mimic those of ordinary least squares. Towards understanding the mechanisms  
7 underlying this phenomenon, we investigate the dynamics of in-context learning of  
8 linear predictors for a transformer with a single linear self-attention layer trained  
9 by gradient flow. We show that despite the non-convexity of the underlying  
10 optimization problem, gradient flow with a random initialization finds a global  
11 minimum of the objective function. Moreover, when given a prompt of labeled  
12 examples from a new linear prediction task, the trained transformer achieves small  
13 prediction error on unlabeled test examples. We further characterize the behavior  
14 of the trained transformer under distribution shifts.

## 15 1 Introduction

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

22 Recently, Garg et al. [2022], von Oswald et al. [2022], Akyürek et al. [2022] initiated the investigation  
23 of ICL from the perspective of learning particular function classes. At a high-level, this refers to when  
24 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  
25  $x_i, x_{\text{query}}$  are sampled i.i.d. from a distribution  $\mathcal{D}_x$  and  $h$  is sampled independently from a distribution  
26 over functions in a function class  $\mathcal{H}$ . The transformer succeeds at in-context learning if when given a  
27 new prompt  $(x'_1, h'(x'_1), \dots, x'_N, h'(x'_N), x'_{\text{query}})$  corresponding to an independently sampled  $h'$  it is  
28 able to formulate a prediction for  $x'_{\text{query}}$  that is close to  $h'(x'_{\text{query}})$  given a sufficiently large number of  
29 examples  $N$ . However, this leaves open the question of how it is that *gradient-based optimization*  
30 *algorithms* over transformer architectures produce models which are capable of in-context learning.

31 In this work, we investigate the learning dynamics of gradient flow in a simplified transformer  
32 architecture when the training prompts consists of random instances of linear regression datasets. We  
33 establish that for a class of transformers with a single layer and with a linear self-attention module  
34 (LSAs), gradient flow on the population loss with a suitable random initialization converges to a global  
35 minimum of the population objective, despite the non-convexity of the underlying objective function.

36 Next, we characterize the learning algorithm that is encoded by the transformer at convergence,  
 37 as well as the prediction error achieved when the model is given a test prompt corresponding to  
 38 a new (and possibly nonlinear) prediction task. Then, we use this to conclude that transformers  
 39 trained by gradient flow indeed in-context learn the class of linear models. Moreover, we characterize  
 40 the robustness of the trained transformer to a variety of distribution shifts. We show that although  
 41 a number of shifts can be tolerated, shifts in the covariate distribution of the features  $x_i$  can not.  
 42 Motivated by this failure under covariate shift, we consider a generalized setting of in-context learning  
 43 where the covariate distribution can vary across prompts. We provide global convergence guarantees  
 44 for LSAs trained by gradient flow in this setting and show that even when trained on a variety of  
 45 covariate distributions, LSAs still fail under covariate shift. We then empirically investigate the  
 46 behavior of large, nonlinear transformers when trained on linear regression prompts. We find that  
 47 these more complex models are able to generalize better under covariate shift, especially when trained  
 48 on prompts with varying covariate distributions.

## 49 2 Preliminaries

50 **In-context learning** We begin by describing a framework for in-context learning of function classes,  
 51 as initiated by Garg et al. [2022]. In-context learning refers to the behavior of models that operate on  
 52 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  
 53 some (unknown) function  $h$  and examples  $x_i$  and query  $x_{\text{query}}$ . The goal for an in-context learner is  
 54 to use the prompt to form a prediction  $\hat{y}(x_{\text{query}})$  for the query such that  $\hat{y}(x_{\text{query}}) \approx h(x_{\text{query}})$ .

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

61 **Definition 2.1** (Trained on in-context examples). *Let  $\mathcal{D}_x$  be a distribution over an input space  $\mathcal{X}$ ,  
 62  $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$  a set of functions  $\mathcal{X} \rightarrow \mathcal{Y}$ , and  $\mathcal{D}_{\mathcal{H}}$  a distribution over functions in  $\mathcal{H}$ . Let  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$   
 63 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-  
 64 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  
 65 parameterized by  $\theta$  in some set  $\Theta$ . For  $N > 0$ , we say that a model  $f : \mathcal{S} \times \mathcal{X} \rightarrow \mathcal{Y}$  is trained on  
 66 in-context examples of functions in  $\mathcal{H}$  under loss  $\ell$  w.r.t.  $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$  if  $f = f_{\theta^*}$  where  $\theta^* \in \Theta$  satisfies*

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

67 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  
 68 during training.

69 As mentioned above, this definition naturally leads to a method for *learning a learning algorithm*  
 70 *from data*: Sample independent prompts by sampling a random function  $h \sim \mathcal{D}_{\mathcal{H}}$  and feature vectors  
 71  $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$ , and then minimize the objective function appearing in (1) using stochastic gradient  
 72 descent or other stochastic optimization algorithms. This procedure returns a model that is learned  
 73 from in-context examples and can form predictions for test (query) examples given a sequence of  
 74 training data. This leads to the following natural definition that quantifies how well such a model  
 75 performs on in-context examples corresponding to a particular hypothesis class.

76 **Definition 2.2** (In-context learning of a hypothesis class). *Let  $\mathcal{D}_x$  be a distribution over an input*  
 77 *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 over functions in  $\mathcal{H}$ . Let*  
 78  *$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  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*  
 79 *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*  
 80 *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*  
 81  *$\ell$  with respect to  $(\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x)$  if there exists a function  $M_{\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x}(\varepsilon) : (0, 1) \rightarrow \mathbb{N}$  such that for every*  
 82  *$\varepsilon \in (0, 1)$ , and for every prompt  $P$  of length  $M \geq M_{\mathcal{D}_{\mathcal{H}}, \mathcal{D}_x}(\varepsilon)$ ,*

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

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

84 Note that in order for a model to in-context learn a hypothesis class, it must be expressive enough  
 85 to achieve arbitrarily small error when sampling a random prompt whose labels are governed by  
 86 some hypothesis  $h$ . With these two definitions in hand, we can formulate the following questions:  
 87 suppose a function class  $\mathcal{F}_\Theta$  is given and  $\mathcal{D}_\mathcal{H}$  corresponds to random instances of hypotheses in a  
 88 hypothesis class  $\mathcal{H}$ . Can a model from  $\mathcal{F}_\Theta$  that is trained on in-context examples of functions in  
 89  $\mathcal{H}$  w.r.t.  $(\mathcal{D}_\mathcal{H}, \mathcal{D}_x)$  in-context learn the hypothesis class  $\mathcal{H}$  w.r.t.  $(\mathcal{D}_\mathcal{H}, \mathcal{D}_x)$ ? How large must the  
 90 training prompts be in order for this to occur? Do standard gradient-based optimization algorithms  
 91 suffice for training the model from in-context examples? How many in-context examples  $M_{\mathcal{D}_\mathcal{H}, \mathcal{D}_x}(\varepsilon)$   
 92 are needed to achieve error  $\varepsilon$ ? In the remaining sections, we shall answer these questions for the  
 93 case of one-layer transformers with linear self-attention modules when the hypothesis class is linear  
 94 models, the loss of interest is the squared loss, and the marginals are (possibly anisotropic) Gaussian  
 95 marginals.

96 **Linear self-attention networks** In this work, we consider a simplified version of the single-layer  
 97 self-attention module [Vaswani et al., 2017]. Let  $x_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$  denote the feature vector  
 98 and its label, and  $E \in \mathbb{R}^{(d+1) \times (N+1)}$  be an embedding matrix that is formed using a prompt  
 99  $(x_1, y_1, \dots, x_N, y_N, x_{\text{query}})$  of length  $N$ . The specific expression of token matrix and the linear  
 100 self-attention(LSA) layer are defined as

$$E = \begin{pmatrix} x_1 & x_2 & \cdots & x_N & x_{\text{query}} \\ y_1 & y_2 & \cdots & y_N & 0 \end{pmatrix}, \quad f_{\text{LSA}}(E; \theta) = E + W^{PV} E \cdot \frac{E^\top W^{KQ} E}{\rho}; \quad (3)$$

101 Here, we have  $\theta = (W^{KQ}, W^{PV})$ , where  $W^{KQ}$  is the merged key-query matrix and  $W^{PV}$  the  
 102 merged projection-value matrix.  $\rho$  is the normalizer which is the width of token matrix  $E$  minus  
 103 one. Under the above token embedding, we take  $\rho = N$ . The prediction for the token  $x_{\text{query}}$  is the  
 104 bottom-right entry of the output matrix, namely,  $\hat{y}_{\text{query}} = \hat{y}_{\text{query}}(E; \theta) = [f_{\text{LSA}}(E; \theta)]_{(d+1), (N+1)}$ .

105 **Training procedure** We assume training prompts are sampled as follows. Let  $\Lambda$  be a posi-  
 106 tive definite covariance matrix. Each training prompt, indexed by  $\tau \in \mathbb{N}$ , takes the form of  
 107  $P_\tau = (x_{\tau,1}, h_\tau(x_{\tau,1}), \dots, x_{\tau,N}, h_\tau(x_{\tau,N}), x_{\tau,\text{query}})$ , where task weights  $w_\tau \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$ , inputs  
 108  $x_{\tau,i}, x_{\tau,\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \Lambda)$ , and labels  $h_\tau(x) = \langle w_\tau, x \rangle$ . Each prompt corresponds to an embedding  
 109 matrix  $E_\tau$ , formed using the transformation (3). We denote the prediction of the LSA model on the  
 110 query label in the task  $\tau$  as  $\hat{y}_{\tau,\text{query}}$ . In this paper, we consider the gradient flow over the population  
 111 loss, which captures the behavior of gradient descent with infinitesimal step size and has dynamics  
 112 given by the following differential equation:

$$\frac{d}{dt} \theta = -\nabla L(\theta), \quad L(\theta) = \frac{1}{2} \mathbb{E}_{w_\tau, x_{\tau,1}, \dots, x_{\tau,N}, x_{\tau,\text{query}}} [(\hat{y}_{\tau,\text{query}}(E; \theta) - \langle w_\tau, x_{\tau,\text{query}} \rangle)^2]. \quad (4)$$

113 For the initialization, we assume

$$W^{PV}(0) = \sigma \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}, \quad W^{KQ}(0) = \sigma \begin{pmatrix} \Theta \Theta^\top & 0_d \\ 0_d^\top & 0 \end{pmatrix}, \quad (5)$$

114 where  $\sigma > 0$  is a parameter, and let  $\Theta \in \mathbb{R}^{d \times d}$  be any matrix satisfying  $\|\Theta \Theta^\top\|_F = 1$  and  
 115  $\Theta \Lambda \neq 0_{d \times d}$ . This initialization is satisfied for a particular class of random initialization schemes: if  
 116  $M$  has i.i.d. entries from a continuous distribution, then by setting  $\Theta \Theta^\top = MM^\top / \|MM^\top\|_F$ , the  
 117 assumption is satisfied almost surely. At a high-level, this initialization allows for the layers to be  
 118 ‘balanced’ throughout the gradient flow trajectory. Random initializations that induce this balanced-  
 119 ness condition have been utilized in a number of theoretical works on deep linear networks [Du et al.,  
 120 2018, Arora et al., 2018, 2019, Azulay et al., 2021]. We leave the question of convergence under  
 121 alternative random initialization schemes for future work.

## 122 3 Main results

### 123 3.1 Global convergence and prediction error for new tasks

124 In this section, we prove that under suitable initialization, gradient flow will converge to a global  
 125 optimum. Due to the space limit, we leave the rigorous proof in the appendix.

126 **Theorem 3.1** (Convergence and limits). *Consider gradient flow of the linear self-attention network*  
 127  *$f_{\text{LSA}}$  over the population loss (4). Suppose in (5) the initialization scale  $\sigma > 0$  satisfies  $\sigma^2 \|\Gamma\|_{\text{op}} \sqrt{d} <$*   
 128 *2. Then, the gradient flow converges to a global minimum of the population loss in (4). Moreover,*  
 129  *$W^{PV}$  and  $W^{KQ}$  converge to  $W_*^{PV}$  and  $W_*^{KQ}$  respectively, where*

$$W_*^{KQ} = c^{-1} \begin{pmatrix} \Gamma^{-1} & 0_d \\ 0_d^\top & 0 \end{pmatrix}, \quad W_*^{PV} = c \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}, \quad \Gamma := \left(1 + \frac{1}{N}\right) \Lambda + \frac{1}{N} \text{tr}(\Lambda) I_d, \quad (6)$$

130 where  $c = [\text{tr}(\Gamma^{-2})]^{1/4}$  is a constant.

131 We note that if we restrict our setting to  $\Lambda = I_d$ , then the limiting solution described found by  
 132 gradient flow is quite similar to the construction of von Oswald et al. [2022].

133 Next, we would like to characterize the prediction error of the trained network described above  
 134 when the network is given a new prompt. In fact, we can generalize to test prompts which could  
 135 take a significantly different form than the training prompts. Consider prompts that are of the form  
 136  $(x_1, y_1, \dots, x_M, y_M, x_{\text{query}})$  where, for some joint distribution  $\mathcal{D}$  over  $(x, y)$  pairs with marginal  
 137 distribution  $\mathcal{D}_x \sim \mathcal{N}(0, \Lambda)$ , we have  $(x_i, y_i) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$  and  $x_{\text{query}} \sim \mathcal{N}(0, \Lambda)$  independently. Note that  
 138 this allows for a label  $y_i$  to be a nonlinear function of the input  $x_i$ . The prediction of the trained  
 139 transformer for this prompt is then

$$\hat{y}_{\text{query}} = x_{\text{query}}^\top \Gamma^{-1} \left( \frac{1}{M} \sum_{i=1}^M y_i x_i \right) \approx x_{\text{query}}^\top \Lambda^{-1} \mathbb{E}[yx] = x_{\text{query}}^\top \left( \underset{w \in \mathbb{R}^d}{\text{argmin}} \mathbb{E}[(y - \langle w, x \rangle)^2] \right). \quad (7)$$

140 Here, when  $N$  and  $M$  are large, the approximation comes from  $\Gamma^{-1} \approx \Lambda^{-1}$  and strong law of large  
 141 numbers. The expectation above is over  $(x, y) \sim \mathcal{D}$ . This result suggests that trained transformers  
 142 in-context learn the *best linear predictor* over a distribution when the test prompt consists of i.i.d.  
 143 samples from a joint distribution over feature-response pairs. In the following theorem, we formalize  
 144 the above and characterize the prediction error when prompts take this form.

145 **Theorem 3.2** (Generalization error). *Let  $\mathcal{D}$  be a distribution over  $(x, y) \in \mathbb{R}^d \times \mathbb{R}$ , whose*  
 146 *marginal distribution on  $x$  is  $\mathcal{D}_x = \mathcal{N}(0, \Lambda)$ . Assume the test prompt is of the form  $P =$*   
 147  *$(x_1, y_1, \dots, x_M, y_M, x_{\text{query}})$ , where  $(x_i, y_i), (x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$ . Let  $f_{\text{LSA}}^*$  be the LSA model*  
 148 *with parameters  $W_*^{PV}$  and  $W_*^{KQ}$  in (6), and  $\hat{y}_{\text{query}}$  is the prediction for  $x_{\text{query}}$  given the prompt.*  
 149 *Assume  $\mathbb{E}_{\mathcal{D}}[y], \mathbb{E}_{\mathcal{D}}[xy], \mathbb{E}_{\mathcal{D}}[y^2 x x^\top]$  exist and are finite. Then, we have*

$$\mathbb{E}(\hat{y}_{\text{query}} - y_{\text{query}})^2 = \min_{w \in \mathbb{R}^d} \mathbb{E}(\langle w, x_{\text{query}} \rangle - y_{\text{query}})^2 + O\left(\frac{1}{M} + \frac{1}{N^2}\right), \quad (8)$$

150 where the expectation is over  $(x_i, y_i), (x_{\text{query}}, y_{\text{query}}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$  and  $O(\cdot)$  hides problem-dependent  
 151 quantities such as  $d$  and  $\Lambda$ .

152 This theorem shows that, provided the length of prompts seen during training ( $N$ ) and the length of  
 153 the test prompt ( $M$ ) is large enough, *a transformer trained by gradient flow from in-context examples*  
 154 *achieves prediction error competitive with the best linear model.* Moreover, our bound shows that  
 155 the length of prompts seen during training and the length of prompts seen at test-time have different  
 156 effects on the expected prediction error: ignoring dimension and covariance-dependent factors, the  
 157 prediction error is at most  $O(1/M + 1/N^2)$ , decreasing more rapidly as a function of the training  
 158 prompt length  $N$  compared to the test prompt length  $M$ . When  $\mathcal{D}$  corresponds to noiseless linear  
 159 models, the error for the best linear predictor vanishes, and a simpler expression for the generalization  
 160 risk is given in Appendix E.

### 161 3.2 Behavior of trained transformer under distribution shifts

162 Using the identity (7), it is straightforward to characterize the behavior of the trained transformer  
 163 under a variety of distribution shifts. In this section, we shall examine a number of shifts that were first  
 164 explored empirically for transformer architectures by Garg et al. [2022]. Although their experiments  
 165 were for transformers trained by gradient descent, we find that (in the case of linear models) many of  
 166 the behaviors of the trained transformers under distribution shift are identical to those predicted by

167 our theoretical characterizations of the performance of transformers with a single linear self-attention  
 168 layer trained by gradient flow on the population.

169 Following Garg et al. [2022], for training prompts of the form  $(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})$ ,  
 170 let us assume  $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x^{\text{train}}$  and  $h \sim \mathcal{D}_h^{\text{train}}$ , while for test prompts let us assume  $x_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x^{\text{test}}$ ,  
 171  $x_{\text{query}} \sim \mathcal{D}_{\text{query}}^{\text{test}}$ , 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:  $\mathcal{D}_x^{\text{train}} \neq \mathcal{D}_x^{\text{test}}$ .

172 In the following, we shall fix  $\mathcal{D}_x^{\text{train}} = \mathcal{N}(0, \Lambda)$  and vary the other distributions. Recall from (7) that  
 173 the prediction for a test prompt  $(x_1, y_1, \dots, x_N, y_N, x_{\text{query}})$  is given by (for  $N$  large), it holds that

$$\hat{y}_{\text{query}} = x_{\text{query}}^\top \Gamma^{-1} \left( \frac{1}{M} \sum_{i=1}^M y_i x_i \right) \approx x_{\text{query}}^\top \Lambda^{-1} \left( \frac{1}{M} \sum_{i=1}^M y_i x_i \right). \quad (9)$$

174 **Task shifts.** These shifts are tolerated easily by the trained transformer. As Theorem E.1 shows,  
 175 the trained transformer is competitive with the best linear model provided the prompt length during  
 176 training and at test time is large enough. In particular, even if the prompt is such that the labels  $y_i$  are  
 177 not given by  $\langle w, x_i \rangle$  for some  $w \sim \mathcal{N}(0, I_d)$ , the trained transformer will compute a prediction which  
 178 has error competitive with the best linear model that fits the test prompt.

179 For example, consider a prompt corresponding to a noisy linear model, so that the prompt consists of a  
 180 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  
 181 sub-Gaussian noise  $\varepsilon_i$ . Then from (7), the prediction of the transformer on query examples is

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

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

186 **Query shifts.** Continuing from (9), it holds that  $\hat{y}_{\text{query}} \approx x_{\text{query}}^\top \Lambda^{-1} \left( \frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right) w$  since  
 187  $y_i = \langle w, x_i \rangle$ . From this we see that whether query shifts can be tolerated hinges upon the distribution  
 188 of the  $x_i$ 's. Since  $\mathcal{D}_x^{\text{train}} = \mathcal{D}_x^{\text{test}}$ , if  $M$  is large then

$$\hat{y}_{\text{query}} \approx x_{\text{query}}^\top \Lambda^{-1} \Lambda w = x_{\text{query}}^\top w. \quad (10)$$

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

193 **Covariate shifts.** In contrast to task and query shifts, covariate shifts cannot be fully tolerated  
 194 in the transformer. This can be easily seen due to the identity (9): when  $\mathcal{D}_x^{\text{train}} \neq \mathcal{D}_x^{\text{test}}$ , then the  
 195 approximation in (10) does not hold as  $\frac{1}{M} \sum_{i=1}^M x_i x_i^\top$  will not cancel  $\Gamma^{-1}$  when  $M$  and  $N$  are large.  
 196 For instance, if we consider test prompts where the covariates are scaled by a constant  $c \neq 1$ , then

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

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

202 It may seem surprising that a transformer trained on linear regression tasks fails in settings where  
 203 ordinary least squares performs well. However, both the linear self-attention transformer we consider

204 and the transformers considered by Garg et al. [2022] were trained on instances of linear regression  
 205 when the covariate distribution  $\mathcal{D}_x$  over the features was fixed across instances. This leads to the  
 206 natural question of what happens if the transformers instead are trained on prompts where the  
 207 covariate distribution varies across instances, which we explore in the following section.

### 208 3.3 Transformers trained on prompts with random covariate distributions

209 The linear self-attention transformer we considered was trained on instances of linear regression  
 210 when the covariate distribution  $\mathcal{D}_x$  over the features was fixed across instances. This leads to  
 211 the natural question of what happens if the transformers instead are trained on prompts where  
 212 the covariate distribution varies across instances. Let us assume that the covariate distribution  
 213  $\mathcal{D}_x$  for each task is sampled from a distribution  $\Delta$ , and training prompts for each task are  
 214  $(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})$  where  $x_i, x_{\text{query}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_x$  and  $h \sim \mathcal{D}_H$ . In this paper, the covariate  
 215 distributions are sampled by first sampling a diagonal matrix  $\Lambda_\tau = \text{diag}(\lambda_{\tau,i} : i \in [d])$  where  $\lambda_{\tau,i}$   
 216 are independent, strictly positive a.s. and have finite third moments. We then sample  $x_i, x_{\text{query}} \sim$   
 217  $N(0, \Lambda_\tau)$  and  $w \sim N(0, I_d)$  with  $y_{\tau,i} = \langle w, x_{\tau,i} \rangle$  and form the token embedding matrix and  
 218 linear self-attention network (3) as before, and again consider gradient flow on the population loss.  
 219

220 We show that in this setting, gradient flow with a suitable random  
 221 initialization converges to a global minimum of the population loss. However, at this global  
 222 minimum, the transformer does not in-context learn the hypothesis class with varying covariate  
 223 distributions, even when the prompt length in the training and test time go to infinity (See Theorem F.2  
 224 and the following discussion). We further examined this random covariance case empirically on standard GPT2  
 225 architecture. We found that when trained on fixed covariance data, the GPT2 model will struggle with  
 226 the random covariance prompt at test time if the variance is large.  
 227  
 228  
 229  
 230  
 231  
 232  
 233  
 234  
 235  
 236  
 237  
 238  
 239

240 When trained on random covariance data however, the model performs better for test prompts from  
 241 higher-variance random covariance matrices, but still fails to match the performance of least squares.  
 242 More details about random covariance case and experiments on GPT2 are in Appendix F and G.

## 243 4 Conclusion and future work

244 In this work, we investigated the dynamics of in-context learning of transformers with a single linear  
 245 self-attention layer under gradient flow on the population loss, when trained on prompts consisting  
 246 of random instances of noiseless linear models over anisotropic Gaussian marginals. Despite non-  
 247 convexity, suitable random initialization leads to convergence to a specific global minimum. We  
 248 found that the trained transformer is robust to task and some query distribution shifts but brittle to  
 249 distribution shifts between training and test covariates, aligning with empirical observations from Garg  
 250 et al. [2022]. Future directions include exploring whether similar results apply to stochastic gradient  
 251 descent with more general initializations and finite step sizes. There’s also interest in understanding  
 252 in-context learning dynamics in deep, nonlinear transformers beyond the single linear self-attention  
 253 layer studied. Another intriguing direction is to determine how those more complex models like GPT2  
 254 provably show robustness against certain types of distribution shifts, especially over linguistic data.  
 255 Additionally, while current in-context learning focuses on fixed covariate distributions, understanding  
 256 its dynamics when these distributions vary across prompts, especially as larger transformers show  
 257 promise but remain sub-optimal, is a compelling research avenue.



Figure 1: Normalized prediction error for GPT2 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 gray dashed line shows the prediction error of zero estimator and the black dashed line that 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.

## 258 References

- 259 Jacob Abernethy, Alekh Agarwal, Teodor V. Marinov, and Manfred K. Warmuth. A mechanism for  
260 sample-efficient in-context learning for sparse retrieval tasks. *Preprint, arXiv:2305.17040*, 2023.
- 261 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement  
262 preconditioned gradient descent for in-context learning. *Preprint, arXiv:2306.00297*, 2023.
- 263 Kabir Ahuja, Madhur Panwar, and Navin Goyal. In-context learning through the bayesian prism.  
264 *arXiv preprint arXiv:2306.04891*, 2023.
- 265 Kartik Ahuja and David Lopez-Paz. A closer look at in-context learning under distribution shifts.  
266 *Preprint, arXiv:2305.16704*, 2023.
- 267 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algo-  
268 rithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*,  
269 2022.
- 270 Cem Anil, Yuhuai Wu, Anders Johan Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Venkatesh  
271 Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length  
272 generalization in large language models. In *Advances in Neural Information Processing Systems*  
273 (*NeurIPS*), 2022.
- 274 Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit  
275 acceleration by overparameterization. In *International Conference on Machine Learning*, pages  
276 244–253, 2018.
- 277 Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix  
278 factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- 279 Shahar Azulay, Edward Moroshko, Mor Shpigel Nacson, Blake E Woodworth, Nathan Srebro, Amir  
280 Globerson, and Daniel Soudry. On the implicit bias of initialization shape: Beyond infinitesimal  
281 mirror descent. In *International Conference on Machine Learning*, pages 468–477, 2021.
- 282 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:  
283 Provable in-context learning with in-context algorithm selection. *Preprint, arXiv:2306.04637*,  
284 2023.
- 285 Mohamed Ali Belabbas. On implicit regularization: Morse functions and applications to matrix  
286 factorization. *arXiv preprint arXiv:2001.04264*, 2020.
- 287 Satwik Bhattamishra, Arkil Patel, and Navin Goyal. On the computational power of transformers and  
288 its implications in sequence modeling. *arXiv preprint arXiv:2006.09286*, 2020.
- 289 Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization:  
290 An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- 291 Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Zhifang Sui, and Furu Wei. Why can gpt learn in-  
292 context? language models secretly perform gradient descent as meta optimizers. *arXiv preprint*  
293 *arXiv:2212.10559*, 2022.
- 294 Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov.  
295 Transformer-xl: Attentive language models beyond a fixed-length context. In *Association for*  
296 *Computational Linguistics (ACL)*, 2019.
- 297 Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal  
298 transformers, 2019.
- 299 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
300 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,  
301 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.  
302 In *International Conference on Learning Representations (ICLR)*, 2021.

- 303 Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous  
304 models: Layers are automatically balanced. *Advances in neural information processing systems*,  
305 31, 2018.
- 306 Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable  
307 creation in self-attention mechanisms. In *International Conference on Machine Learning*, 2022.
- 308 Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn  
309 in-context? a case study of simple function classes. *arXiv preprint arXiv:2208.01066*, 2022.
- 310 Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro.  
311 Implicit regularization in matrix factorization. *Advances in Neural Information Processing Systems*,  
312 30, 2017.
- 313 Chi Han, Ziqi Wang, Han Zhao, and Heng Ji. In-context learning of large language models explained  
314 as kernel regression, 2023.
- 315 Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure.  
316 *Advances in Neural Information Processing Systems*, 35:37822–37836, 2022.
- 317 Jikai Jin, Zhiyuan Li, Kaifeng Lyu, Simon S Du, and Jason D Lee. Understanding incremental learning  
318 of gradient descent: A fine-grained analysis of matrix sensing. *arXiv preprint arXiv:2301.11500*,  
319 2023.
- 320 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
321 *arXiv:1412.6980*, 2014.
- 322 Shuai Li, Zhao Song, Yu Xia, Tong Yu, and Tianyi Zhou. The closeness of in-context learning and  
323 weight shifting for softmax regression. *arXiv preprint arXiv:2304.13276*, 2023a.
- 324 Yingcong Li, M Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as  
325 algorithms: Generalization and stability in in-context learning. *arXiv preprint arXiv:2301.07067*,  
326 2023b.
- 327 Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized  
328 matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*,  
329 pages 2–47, 2018.
- 330 Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do transformers learn topic structure: Towards a  
331 mechanistic understanding. *arXiv preprint arXiv:2303.04245*, 2023c.
- 332 Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent  
333 for matrix factorization: Greedy low-rank learning. *arXiv preprint arXiv:2012.09839*, 2020.
- 334 Valerii Likhoshesterov, Krzysztof Choromanski, and Adrian Weller. On the expressive power of  
335 self-attention matrices. *arXiv preprint arXiv:2106.03764*, 2021.
- 336 Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers  
337 learn shortcuts to automata. In *International Conference on Learning Representations (ICLR)*,  
338 2023.
- 339 AR Meenakshi and C Rajian. On a product of positive semidefinite matrices. *Linear algebra and its*  
340 *applications*, 295(1-3):3–6, 1999.
- 341 JV Michalowicz, JM Nichols, F Bucholtz, and CC Olson. An isserlis’ theorem for mixed gaussian  
342 variables: Application to the auto-bispectral density. *Journal of Statistical Physics*, 136:89–102,  
343 2009.
- 344 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke  
345 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv*  
346 *preprint arXiv:2202.12837*, 2022.
- 347 OpenAI. Gpt-4 technical report, 2023.



- 348 Jorge Pérez, Javier Marinković, and Pablo Barceló. On the turing completeness of modern neural  
349 network architectures. *arXiv preprint arXiv:1901.03429*, 2019.
- 350 Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University  
351 of Denmark*, 7(15):510, 2008.
- 352 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language  
353 understanding by generative pre-training. 2018.
- 354 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
355 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 356 Mahdi Soltanolkotabi, Dominik Stöger, and Changzhi Xie. Implicit balancing and regularization:  
357 Generalization and convergence guarantees for overparameterized asymmetric matrix sensing.  
358 *arXiv preprint arXiv:2303.14244*, 2023.
- 359 Asher Trockman and J Zico Kolter. Mimetic initialization of self-attention layers. *arXiv preprint  
360 arXiv:2305.09828*, 2023.
- 361 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
362 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing  
363 Systems*, 30, 2017.
- 364 Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,  
365 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent.  
366 *arXiv preprint arXiv:2212.07677*, 2022.
- 367 Xinyi Wang, Wanrong Zhu, and William Yang Wang. Large language models are implicitly topic  
368 models: Explaining and finding good demonstrations for in-context learning. *arXiv preprint  
369 arXiv:2301.11916*, 2023.
- 370 Gian-Carlo Wick. The evaluation of the collision matrix. *Physical review*, 80(2):268, 1950.
- 371 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
372 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art  
373 natural language processing. In *Proceedings of the 2020 conference on empirical methods in  
374 natural language processing: system demonstrations*, pages 38–45, 2020.
- 375 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context  
376 learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- 377 Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J Reddi, and Sanjiv Kumar.  
378 Are transformers universal approximators of sequence-to-sequence functions? *arXiv preprint  
379 arXiv:1912.10077*, 2019.
- 380 Chulhee Yun, Yin-Wen Chang, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv  
381 Kumar. O (n) connections are expressive enough: Universal approximability of sparse transformers.  
382 *Advances in Neural Information Processing Systems*, 33:13783–13794, 2020.
- 383 Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does in-  
384 context learning learn? bayesian model averaging, parameterization, and generalization. *Preprint,  
385 arXiv:2305.19420*, 2023.

## 386 A Notations

387 In this section, we briefly describe the notation we use in the paper. We write  $[n] = \{1, 2, \dots, n\}$ . We  
388 use  $\otimes$  to denote the Kronecker product, and  $\text{Vec}$  the vectorization operator in column-wise order.  
389 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  $A, B \in \mathbb{R}^{m \times n}$   
390 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  $n$   
391 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,  
392 respectively. We denote the matrix operator norm and Frobenius norm as  $\|\cdot\|_{op}$  and  $\|\cdot\|_F$ . We use  $I_d$   
393 to denote the  $d$ -dimensional identity matrix and sometimes we also use  $I$  when the dimension is clear  
394 from the context. For a positive semi-definite matrix  $A$ , we write  $\|x\|_A^2 := x^\top Ax$ . Unless otherwise  
395 defined, we use lower case letters for scalars and vectors, and use upper case letters for matrices.

## 396 B Additional related works

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

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

408 Akyürek et al. [2022] and von Oswald et al. [2022] examined the behavior of transformers when  
409 trained on random instances of linear regression, as we do in this work. They considered the setting  
410 of isotropic Gaussian data with isotropic Gaussian weight vectors, and showed that the trained  
411 transformer's predictions mimic those of a single step of gradient descent. They also provided a  
412 construction of transformers which implement this single step of gradient descent. By contrast, we  
413 explicitly show that gradient flow provably converges to transformers which learn linear models  
414 in-context. Moreover, our analysis holds when the covariates are anisotropic Gaussians, for which a  
415 single step of vanilla gradient descent is unable to achieve small prediction error.<sup>1</sup>

416 Let us briefly mention a number of other works on understanding in-context learning in transformers  
417 and other sequence-based models. Han et al. [2023] suggests that Bayesian inference on prompts can  
418 be asymptotically interpreted as kernel regression. Dai et al. [2022] interprets ICL as implicit fine-  
419 tuning, viewing large language models as meta-optimizers performing gradient-based optimization.  
420 Xie et al. [2021] regards ICL as implicit Bayesian inference, with transformers learning a shared  
421 latent concept between prompts and test data, and they prove the ICL property when the training  
422 distribution is a mixture of HMMs. Similarly, Wang et al. [2023] perceives ICL as a Bayesian  
423 selection process, implicitly inferring information pertinent to the designated tasks. Li et al. [2023a]  
424 explores the functional resemblance between a single layer of self-attention and gradient descent on  
425 a softmax regression problem, offering upper bounds on their difference. Min et al. [2022] notes  
426 that the alteration of label parts in prompts does not drastically impair the ICL ability. They contend  
427 that ICL is invoked when prompts reveal information about the label space, input distribution, and  
428 sequence structure.

429 Another collection of works have sought to understand transformers from an approximation theoretic  
430 perspective. Yun et al. [2019, 2020] established that transformers can universally approximate any  
431 sequence-to-sequence function under some assumptions. Investigations by Edelman et al. [2022],  
432 Likhoshervstov et al. [2021] indicate that a single-layer self-attention can learn sparse functions of  
433 the input sequence, where sample complexity and hidden size are only logarithmic relative to the

---

<sup>1</sup>To see this, suppose  $(x_i, y_i)$  are i.i.d. with  $x \sim \mathcal{N}(0, \Lambda)$  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 is not close to  $x^\top w$  when  $\Lambda \neq I_d$ .

434 sequence length. Further studies by Pérez et al. [2019], Dehghani et al. [2019], Bhattamishra et al.  
435 [2020] indicate that the vanilla transformer and its variants exhibit Turing completeness. Liu et al.  
436 [2023] showed that transformers can approximate finite-state automata with few layers. Bai et al.  
437 [2023] showed that transformers can implement a variety of statistical machine learning algorithms  
438 as well as model selection procedures. Abernethy et al. [2023] showed that a pretrained transformer  
439 can be used to define a transformer that segments a prompt into examples and labels and learns to  
440 solve a sparse retrieval task. Zhang et al. [2023] interpreted in-context learning via a Bayesian model  
441 averaging process.

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

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

## 457 C Linear self-attention and training procedure

### 458 C.1 Linear self-attention and the prediction

459 Before describing the particular transformer models we analyze in this work, we first recall the  
 460 definition of the softmax-based single-head self-attention module [Vaswani et al., 2017]. Let  $E \in$   
 461  $\mathbb{R}^{d_e \times d_N}$  be an embedding matrix that is formed using a prompt  $(x_1, y_1, \dots, x_N, y_N, x_{\text{query}})$  of  
 462 length  $N$ . The user has the freedom to determine how this embedding matrix is formed from the  
 463 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  
 464 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  
 465  $d_N = N + 1$ . 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  
 466 matrices,  $W^P \in \mathbb{R}^{d_e \times d_v}$  the projection matrix, and  $\rho > 0$  a normalization factor. The softmax  
 467 self-attention module takes as input an embedding matrix  $E$  of width  $d_N$  and outputs a matrix of the  
 468 same size,

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

469 where softmax is applied column-wise and, given a vector input of  $v$ , the  $i$ -th entry of  $\text{softmax}(v)$  is  
 470 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  
 471 the *self-attention matrix*. Note that  $f_{\text{Attn}}$  can take as its input a sequence of arbitrary length.

472 In this work, we consider a simplified version of the single-layer self-attention module, which is  
 473 more amenable to theoretical analysis and yet is still capable of in-context learning linear models.  
 474 In particular, we consider a single-layer linear self-attention (LSA) model, which is a modified  
 475 version of  $f_{\text{Attn}}$  where we remove the softmax nonlinearity, merge the projection and value matrices  
 476 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  
 477  $W^{KQ} \in \mathbb{R}^{d_e \times d_e}$ . We concatenate these matrices into  $\theta = (W^{KQ}, W^{PV})$  and denote

$$f_{\text{LSA}}(E; \theta) = E + W^{PV} E \cdot \frac{E^\top W^{KQ} E}{\rho}. \quad (11)$$

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

483 The user has the flexibility to determine the method for constructing the embedding matrix from a  
 484 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  
 485 following embedding, which stacks  $(x_i, y_i)^\top \in \mathbb{R}^{d+1}$  into the first  $N$  columns with  $(x_{\text{query}}, 0)^\top \in$   
 486  $\mathbb{R}^{d+1}$  as the last column:

$$E = E(P) = \begin{pmatrix} x_1 & x_2 & \cdots & x_N & x_{\text{query}} \\ y_1 & y_2 & \cdots & y_N & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (N+1)}. \quad (12)$$

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

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

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

496 Here and after, we may occasionally suppress dependence on  $\theta$  and write  $\hat{y}_{\text{query}}(E; \theta)$  as  $\hat{y}_{\text{query}}$ . Since  
 497 the prediction takes only the right-bottom entry of the token matrix output by the LSA layer, actually  
 498 only part of  $W^{PV}$  and  $W^{KQ}$  affect the prediction. To see how, let us denote

$$W^{PV} = \begin{pmatrix} W_{11}^{PV} & w_{12}^{PV} \\ (w_{21}^{PV})^\top & w_{22}^{PV} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}, \quad W^{KQ} = \begin{pmatrix} W_{11}^{KQ} & w_{12}^{KQ} \\ (w_{21}^{KQ})^\top & w_{22}^{KQ} \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}, \quad (13)$$

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

$$\hat{y}_{\text{query}} = \left( (w_{21}^{PV})^\top \quad w_{22}^{PV} \right) \cdot \left( \frac{EE^\top}{N} \right) \begin{pmatrix} W_{11}^{KQ} \\ (w_{21}^{KQ})^\top \end{pmatrix} x_{\text{query}}, \quad (14)$$

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

## 503 C.2 Training procedure and the initialization

504 In this work, we will consider the task of in-context learning linear predictors. We will assume  
505 training prompts are sampled as follows. Let  $\Lambda$  be a positive definite covariance matrix. Each training  
506 prompt, indexed by  $\tau \in \mathbb{N}$ , takes the form of  $P_\tau = (x_{\tau,1}, h_\tau(x_{\tau,1}), \dots, x_{\tau,N}, h_\tau(x_{\tau,N}), x_{\tau,\text{query}})$ ,  
507 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  $h_\tau(x) = \langle w_\tau, x \rangle$ .

508 Each prompt corresponds to an embedding matrix  $E_\tau$ , formed using the transformation (3):

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

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

$$\hat{L}(\theta) = \frac{1}{2B} \sum_{\tau=1}^B \left( \hat{y}_{\tau,\text{query}} - \langle w_\tau, x_{\tau,\text{query}} \rangle \right)^2. \quad (15)$$

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

$$L(\theta) = \lim_{B \rightarrow \infty} \hat{L}(\theta) = \frac{1}{2} \mathbb{E}_{w_\tau, x_{\tau,1}, \dots, x_{\tau,N}, x_{\tau,\text{query}}} \left[ (\hat{y}_{\tau,\text{query}} - \langle w_\tau, x_{\tau,\text{query}} \rangle)^2 \right] \quad (16)$$

514 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  
515 vector  $w_\tau$ , 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  
516 of gradient descent with infinitesimal step size and has dynamics given by the following differential  
517 equation:

$$\frac{d}{dt} \theta = -\nabla L(\theta). \quad (17)$$

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

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

$$W^{PV}(0) = \sigma \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}, \quad W^{KQ}(0) = \sigma \begin{pmatrix} \Theta\Theta^\top & 0_d \\ 0_d^\top & 0 \end{pmatrix}. \quad (18)$$

521 This initialization is satisfied for a particular class of random initialization schemes: if  $M$  has i.i.d.  
522 entries from a continuous distribution, then by setting  $\Theta\Theta^\top = MM^\top / \|MM^\top\|_F$ , the assumption  
523 is satisfied almost surely. The reason we use this particular initialization scheme will be made more  
524 clear when we describe the proof, but at a high-level this is due to the fact that the predictions (14) can  
525 be viewed as the output of a two-layer linear network, and initializations satisfying Assumption C.1  
526 allow for the layers to be ‘balanced’ throughout the gradient flow trajectory. Random initializations  
527 that induce this balancedness condition have been utilized in a number of theoretical works on deep  
528 linear networks [Du et al., 2018, Arora et al., 2018, 2019, Azulay et al., 2021]. We leave the question  
529 of convergence under alternative random initialization schemes for future work.

530 **D Theorem 3.1 and the proof**

531 We first formally describe the theorem on global convergence and the expression for the limits:

532 **Theorem D.1** (Convergence and limits). *Consider gradient flow of the linear self-attention network*  
 533  *$f_{\text{LSA}}$  defined in (3) over the population loss (16). Suppose the initialization satisfies Assumption C.1*  
 534 *with initialization scale  $\sigma > 0$  satisfying  $\sigma^2 \|\Gamma\|_{\text{op}} \sqrt{d} < 2$  where we have defined*

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

535 *Then gradient flow converges to a global minimum of the population loss (16). Moreover,  $W^{PV}$  and*  
 536  *$W^{KQ}$  converge to  $W_*^{PV}$  and  $W_*^{KQ}$  respectively, where*

$$W_*^{KQ} = [\text{tr}(\Gamma^{-2})]^{-\frac{1}{4}} \begin{pmatrix} \Gamma^{-1} & 0_d \\ 0_d^\top & 0 \end{pmatrix}, \quad W_*^{PV} = [\text{tr}(\Gamma^{-2})]^{\frac{1}{4}} \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}. \quad (19)$$

537 **D.1 Proof of Theorem D.1**

538 In this section, we briefly outline the proof sketch of Theorem D.1.

539 **D.1.1 Equivalence to a quadratic optimization problem**

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

544 The first key step in our proof is to recognize that the prediction  $\hat{y}_{\text{query}}(E_\tau; \theta)$  in the linear self-  
 545 attention model can be written as the output of a quadratic function  $u^\top H_\tau u$  for some matrix  $H_\tau$  de-  
 546 pending on the token embedding matrix  $E_\tau$  and for some vector  $u$  depending on  $\theta = (W^{KQ}, W^{PV})$ .  
 547 This is shown in the following lemma, the proof of which is provided in Appendix D.2.1.

548 **Lemma D.2.** *Let  $E_\tau \in \mathbb{R}^{(d+1) \times (N+1)}$  be an embedding matrix corresponding to a prompt of length*  
 549  *$N$  and weight  $w_\tau$ . Then the prediction  $\hat{y}_{\text{query}}(E_\tau; \theta)$  for the query covariate can be written as the*  
 550 *output of a quadratic function,*

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

551 *where the matrix  $H_\tau$  is defined as,*

$$H_\tau = \frac{1}{2} X_\tau \otimes \left( \frac{E_\tau E_\tau^\top}{N} \right) \in \mathbb{R}^{(d+1)^2 \times (d+1)^2}, \quad X_\tau = \begin{pmatrix} 0_{d \times d} & x_{\tau, \text{query}} \\ (x_{\tau, \text{query}})^\top & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)} \quad (20)$$

552 *and*

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

553 *where  $U_{11} = W_{11}^{KQ} \in \mathbb{R}^{d \times d}$ ,  $u_{12} = w_{21}^{PV} \in \mathbb{R}^{d \times 1}$ ,  $u_{21} = w_{21}^{KQ} \in \mathbb{R}^{d \times 1}$ ,  $u_{-1} = w_{22}^{PV} \in \mathbb{R}$*   
 554 *correspond to particular components of  $W^{PV}$  and  $W^{KQ}$ , defined in (13).*

555

556 This implies that we can write the original loss function (15) as

$$\hat{L} = \frac{1}{2B} \sum_{\tau=1}^B (u^\top H_\tau u - w_\tau^\top x_{\tau, \text{query}})^2. \quad (21)$$

557 Thus, our problem is reduced to *understanding the dynamics of an optimization algorithm defined in*  
 558 *terms of a quadratic function.* We also note that this quadratic optimization problem is an instance of

559 a rank-one matrix factorization problem, a problem well-studied in the deep learning theory literature  
 560 [Gunasekar et al., 2017, Arora et al., 2019, Li et al., 2018, Chi et al., 2019, Belabbas, 2020, Li et al.,  
 561 2020, Jin et al., 2023, Soltanolkotabi et al., 2023].

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

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

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

### 570 D.1.2 Dynamical system of gradient flow

571 We now describe the dynamical system for the coordinates of  $u$  above. We prove the following lemma  
 572 in Appendix D.2.2.

573 **Lemma D.3.** *Let  $u = \text{Vec}(U) := \text{Vec} \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix}$  as in Lemma D.2. Consider gradient flow*

574 *over*

$$L := \frac{1}{2} \mathbb{E} (u^\top H_\tau u - w_\tau^\top x_{\tau, \text{query}})^2 \quad (22)$$

575 *with respect to  $u$  starting from an initial value satisfying Assumption C.1. Then the dynamics of  $U$*   
 576 *follows*

$$\begin{aligned} \frac{d}{dt} U_{11}(t) &= -u_{-1}^2 \Gamma \Lambda U_{11} \Lambda + u_{-1} \Lambda^2 \\ \frac{d}{dt} u_{-1}(t) &= -\text{tr} [u_{-1} \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - \Lambda^2 (U_{11})^\top], \end{aligned} \quad (23)$$

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

578

579 We see that the dynamics are governed by a complex system of  $d^2 + 1$  coupled differential equations.  
 580 Moreover, basic calculus (for details, see Lemma D.6) shows that these dynamics are the same as  
 581 those of gradient flow on the following objective function:

$$\tilde{\ell} : \mathbb{R}^{d \times d} \times \mathbb{R} \rightarrow \mathbb{R}, \quad \tilde{\ell}(U_{11}, u_{-1}) = \text{tr} \left[ \frac{1}{2} u_{-1}^2 \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - u_{-1} \Lambda^2 (U_{11})^\top \right]. \quad (24)$$

582 Actually, the loss function  $\tilde{\ell}$  is simply the loss function  $L$  in (22) plus some constants that do not  
 583 depend on the parameter  $u$ . Therefore our problem is reduced to studying the dynamics of gradient  
 584 flow on the above objective function.

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

588 **Lemma D.4.** *For any global minimum of  $\tilde{\ell}$ , we have*

$$u_{-1} U_{11} = \Gamma^{-1}. \quad (25)$$

589 Putting this together with Lemma D.3, we see that at those global minima of the population objective  
 590 satisfying  $U_{11} = (c\Gamma)^{-1}$ ,  $u_{-1} = c$  and  $u_{12} = u_{21} = 0_d$ , the transformer's predictions for a new  
 591 linear regression task prompt are given by

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

592 Thus, the only remaining task is to show global convergence when gradient flow has an initialization  
 593 satisfying Assumption C.1.

### 594 **D.1.3 PL inequality and global convergence**

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

598 **Lemma D.5.** *Suppose the initialization of gradient flow satisfies Assumption C.1 with initialization*  
 599 *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*

$$\mu := \frac{\sigma^2}{\sqrt{d}\|\Lambda\|_{op}^2 \text{tr}(\Gamma^{-1}\Lambda^{-1}) \text{tr}(\Lambda^{-1})} \|\Lambda\Theta\|_F^2 \left[2 - \sqrt{d}\sigma^2 \|\Gamma\|_{op}\right] > 0, \quad (26)$$

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

$$\left\|\nabla\tilde{\ell}(U_{11}(t), u_{-1}(t))\right\|_2^2 := \left\|\frac{\partial\tilde{\ell}}{\partial U_{11}}\right\|_F^2 + \left|\frac{\partial\tilde{\ell}}{\partial u_{-1}}\right|^2 \geq \mu \left(\tilde{\ell}(U_{11}(t), u_{-1}(t)) - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1})\right). \quad (27)$$

601 *Moreover, gradient flow converges to the global minimum of  $\tilde{\ell}$ , and  $U_{11}$  and  $u_{-1}$  converge to the*  
 602 *following,*

$$\lim_{t \rightarrow \infty} u_{-1}(t) = \|\Gamma^{-1}\|_F^{\frac{1}{2}} \text{ and } \lim_{t \rightarrow \infty} U_{11}(t) = \|\Gamma^{-1}\|_F^{-\frac{1}{2}} \Gamma^{-1}. \quad (28)$$

603

604 With these observations, proving Theorem D.1 becomes a direct application of Lemma D.2, D.3, D.4,  
 605 and Lemma D.5. It then only requires translating  $U_{11}$  and  $u_{-1}$  back to the original parameterization  
 606 using  $W^{PV}$  and  $W^{KQ}$ .

## 607 **D.2 Proof for supporting lemmas**

608 In this section, we prove Lemma D.2, Lemma D.3, Lemma D.4 and Lemma D.5. Theorem D.1 is a  
 609 natural corollary of these four lemmas when we translate  $u_{-1}$  and  $U_{11}$  back to  $W^{PV}$  and  $W^{KQ}$ .

### 610 **D.2.1 Proof of Lemma D.2**

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

612 **Lemma D.2.** *Let  $E_\tau \in \mathbb{R}^{(d+1) \times (N+1)}$  be an embedding matrix corresponding to a prompt of length*  
 613  *$N$  and weight  $w_\tau$ . Then the prediction  $\hat{y}_{\text{query}}(E_\tau; \theta)$  for the query covariate can be written as the*  
 614 *output of a quadratic function,*

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

615 *where the matrix  $H_\tau$  is defined as,*

$$H_\tau = \frac{1}{2} X_\tau \otimes \left(\frac{E_\tau E_\tau^\top}{N}\right) \in \mathbb{R}^{(d+1)^2 \times (d+1)^2}, \quad X_\tau = \begin{pmatrix} 0_{d \times d} & x_{\tau, \text{query}} \\ (x_{\tau, \text{query}})^\top & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)} \quad (20)$$

616 *and*

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

617 *where  $U_{11} = W_{11}^{KQ} \in \mathbb{R}^{d \times d}$ ,  $u_{12} = w_{21}^{PV} \in \mathbb{R}^{d \times 1}$ ,  $u_{21} = w_{21}^{KQ} \in \mathbb{R}^{d \times 1}$ ,  $u_{-1} = w_{22}^{PV} \in \mathbb{R}$*   
 618 *correspond to particular components of  $W^{PV}$  and  $W^{KQ}$ , defined in (13).*

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

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



621 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  
622  $i \in [d]$ . Then, we have

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

623 Here, we use some algebraic facts about matrix vectorization, Kronecker product and trace. For  
624 reference, we refer to [Petersen et al., 2008].  $\square$

### 625 D.2.2 Proof of Lemma D.3

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

627 **Lemma D.3.** *Let  $u = \text{Vec}(U) := \text{Vec} \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix}$  as in Lemma D.2. Consider gradient flow*

628 *over*

$$628 \quad L := \frac{1}{2} \mathbb{E} (u^\top H_\tau u - w_\tau^\top x_{\tau, \text{query}})^2 \quad (22)$$

629 *with respect to  $u$  starting from an initial value satisfying Assumption C.1. Then the dynamics of  $U$*   
630 *follows*

$$\begin{aligned}
\frac{d}{dt} U_{11}(t) &= -u_{-1}^2 \Gamma \Lambda U_{11} \Lambda + u_{-1} \Lambda^2 \\
\frac{d}{dt} u_{-1}(t) &= -\text{tr} [u_{-1} \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - \Lambda^2 (U_{11})^\top],
\end{aligned} \quad (23)$$

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

632 *Proof.* From the definition of  $L$  in (22) and the dynamics of gradient flow, we calculate the derivatives  
633 of  $u$ . Here, we use the chain rule and some facts about matrix derivatives. See Lemma H.1 for  
634 reference.

$$634 \quad \frac{du}{dt} = -2\mathbb{E} (\langle H_\tau, uu^\top \rangle H_\tau) u + 2\mathbb{E} (w_\tau^\top x_{\tau, \text{query}} H_\tau) u. \quad (29)$$

635 **Step One: Calculate the Second Term** We first calculate the second term. From the definition of  
636  $H_\tau$ , we have

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

637 For ease of notation, we denote

$$637 \quad \widehat{\Lambda}_\tau := \frac{1}{N} \sum_{i=1}^N x_{\tau, i} x_{\tau, i}^\top. \quad (30)$$

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

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

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

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

640 Therefore, we have

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

641 Since  $X_\tau$  only depends on  $x_{\tau, \text{query}}$  by definition, and  $x_{\tau, \text{query}}$  is independent of  $w_\tau$  and  $x_{\tau, i}$ ,  $i =$   
642  $1, 2, \dots, N$ , we have

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

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

$$\mathbb{E} [w_\tau^\top x_{\tau, \text{query}} H_\tau] u = \frac{1}{2} \begin{pmatrix} \mathbf{0}_{d(d+1) \times d(d+1)} & A \\ A^\top & \mathbf{0}_{(d+1) \times (d+1)} \end{pmatrix} \cdot u, \quad (31)$$

646 where

$$A = \begin{pmatrix} V_1 + V_1^\top \\ V_2 + V_2^\top \\ \dots \\ V_d + V_d^\top \end{pmatrix} \in \mathbb{R}^{d(d+1) \times (d+1)}, \quad V_j = \begin{pmatrix} 0_{d \times d} & \sum_{i=1}^d \Lambda_{ij} \Lambda_i \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0_{d \times d} & \Lambda \Lambda_j \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}. \quad (32)$$

**Step Two: Calculate the First Term** Next, we compute the first term in (29), namely

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

647 For simplicity, we denote  $Z_\tau := \frac{1}{N} E_\tau E_\tau^\top$ . Using the definition of  $H_\tau$  in (20) and Lemma H.1, we  
648 have

$$\begin{aligned} D &= 2\mathbb{E} (\langle H_\tau, uu^\top \rangle H_\tau u) && \text{(definition)} \\ &= \frac{1}{2} \mathbb{E} \left[ \text{tr} \left( X_\tau \otimes Z_\tau \text{Vec}(U) \text{Vec}(U)^\top \right) (X_\tau \otimes Z_\tau) \text{Vec}(U) \right] \\ & && \text{(definition of } H_\tau \text{ in (20) and } u = \text{Vec}(U)) \\ &= \frac{1}{2} \mathbb{E} \left[ \text{tr} \left( \text{Vec}(Z_\tau U X_\tau) \text{Vec}(U)^\top \right) \text{Vec}(Z_\tau U X_\tau) \right] \\ & && \text{(Vec}(AXB) = (B^\top \otimes A) \text{Vec}(X) \text{ in Lemma H.1)} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \mathbb{E} \left[ \text{Vec}(U)^\top \cdot \text{Vec}(Z_\tau U X_\tau) \cdot \text{Vec}(Z_\tau U X_\tau) \right] \\
&\hspace{20em} \text{(property of trace operator)} \\
&= \frac{1}{2} \mathbb{E} \left[ \sum_{i,j=1}^{d+1} \left( (Z_\tau U X_\tau)_{ij} U_{ij} \right) \text{Vec}(Z_\tau U X_\tau) \right].
\end{aligned}$$

649 **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  
650  $\frac{d}{dt} u_{21} = 0_d$ . If this is true, then these two blocks will be zero all the time since we assume they are  
651 zero at initial time in Assumption C.1. We denote  $A_{k\cdot}$  and  $A_{\cdot k}$  as the  $k$ -th row and  $k$ -th column of  
652 matrix  $A$ , respectively.

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

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

654 Written in an entry-wise manner, it will be

$$(Z_\tau U X_\tau)_{kl} = \begin{cases} \left( \widehat{\Lambda}_\tau \right)_{k:} w_\tau u_{-1} x_{\tau, \text{query}}^l & k, l \in [d] \\ \left( \widehat{\Lambda}_\tau + \frac{1}{N} x_{\tau, \text{query}} \cdot x_{\tau, \text{query}}^\top \right)_{k:} U_{11} x_{\tau, \text{query}} & k \in [d], l = d+1 \\ w_\tau^\top \left( \widehat{\Lambda}_\tau \right) w_\tau u_{-1} x_{\tau, \text{query}}^l & l \in [d], k = d+1 \\ w_\tau^\top \left( \widehat{\Lambda}_\tau \right) U_{11} x_{\tau, \text{query}} & k = l = d+1 \end{cases}. \quad (33)$$

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

$$\begin{aligned}
D_{k,d+1} &= \frac{1}{2} \mathbb{E} \left[ \sum_{i,j=1}^{d+1} \left( (Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{k,d+1} \right] \\
&= \frac{1}{2} \mathbb{E} \left[ \sum_{i,j=1}^d \left( (Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{k,d+1} \right] + \frac{1}{2} \mathbb{E} \left[ \left( (Z_\tau U X_\tau)_{d+1,d+1} u_{-1} \right) (Z_\tau U X_\tau)_{k,d+1} \right],
\end{aligned} \quad (34)$$

657 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,  
658 we fix  $i, j \in [d]$  and have

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

659 since  $w_\tau$  is independent with all prompt input and query input, namely all  $x_{\tau,i}$  for  $i \in [\text{query}]$ , and  
660  $w_\tau$  is mean zero. Similarly, for the second term of (34), we have

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

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

663

664 For  $k \in [d]$ , to calculate the derivative of  $U_{k,d+1}$ , it suffices to further calculate the inner product of  
 665 the  $d(d+1) + k$  th row of  $\mathbb{E} [w_\tau^\top x_{\tau,\text{query}} H_\tau]$  and  $u$ . From (31), we know this

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

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

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

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

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

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

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

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

673 For the last term, we have

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

674 So we have

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

675 Additionally, we have

$$\begin{aligned} 2 \left[ \mathbb{E} \left( w_\tau^\top x_{\tau,\text{query}} H_\tau \right) u \right]_{(l-1)(d+1)+k} &= \left[ \begin{pmatrix} \mathbf{0}_{d(d+1) \times d(d+1)} & A \\ A^\top & \mathbf{0}_{(d+1) \times (d+1)} \end{pmatrix} \cdot u \right]_{(l-1)(d+1)+k} \\ &= \left( \mathbf{0}_{(d+1) \times d(d+1)} \quad V_l + V_l^\top \right)_{k:} \cdot U \\ &= \Lambda_k^\top \Lambda_l u_{-1}. \end{aligned} \quad \begin{array}{l} \text{(definition)} \\ \text{(definition of } A \text{ in (32))} \\ \text{(definition of } V_i \text{ in (32))} \end{array}$$

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

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

677 which implies

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

678

679 From the definition of  $\widehat{\Lambda}_\tau$  (equation (30)), the independence and Gaussianity of  $x_{\tau,i}$  and Lemma H.2,  
 680 we compute

$$\begin{aligned} \mathbb{E} \left( \left( \widehat{\Lambda}_\tau \right)^2 \right) &= \mathbb{E} \left( \left( \frac{1}{N} \sum_{i=1}^N x_{\tau,i} x_{\tau,i}^\top \right)^2 \right) && \text{(definition (30))} \\ &= \frac{N-1}{N} \left[ \mathbb{E} \left( x_{\tau,1} x_{\tau,1}^\top \right) \right]^2 + \frac{1}{N} \mathbb{E} \left( x_{\tau,1} x_{\tau,1}^\top x_{\tau,1} x_{\tau,1}^\top \right) \\ &&& \text{(independence between prompt input)} \\ &= \frac{N+1}{N} \Lambda^2 + \frac{1}{N} \text{tr}(\Lambda) \Lambda. && \text{(Lemma H.2)} \end{aligned}$$

681 We define

$$\Gamma := \frac{N+1}{N} \Lambda + \frac{1}{N} \text{tr}(\Lambda) I_d. \quad (35)$$

682 Then, from (29), we know the dynamics of  $U_{11}$  is

$$\frac{d}{dt} U_{11} = -u_{-1}^2 \Gamma \Lambda U_{11} \Lambda + u_{-1} \Lambda^2. \quad (36)$$

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

$$\begin{aligned} D_{d+1,d+1} &= \frac{1}{2} \mathbb{E} \left[ \sum_{i,j=1}^d \left( (Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{d+1,d+1} \right] \\ &\quad + \frac{1}{2} \mathbb{E} \left[ \left( (Z_\tau U X_\tau)_{d+1,d+1} u_{-1} \right) (Z_\tau U X_\tau)_{d+1,d+1} \right]. \end{aligned} \quad (37)$$

684 For the first term above, we have

$$\begin{aligned} &\mathbb{E} \left[ \sum_{i,j=1}^d \left( (Z_\tau U X_\tau)_{ij} U_{ij} \right) (Z_\tau U X_\tau)_{d+1,d+1} \right] \\ &= u_{-1} \sum_{i,j=1}^d U_{ij} \mathbb{E} \left[ \left( \widehat{\Lambda}_\tau \right)_{i:} \cdot w_\tau w_\tau^\top \cdot \left( \widehat{\Lambda}_\tau \right) \cdot U_{11} x_{\tau,\text{query}} x_{\tau,\text{query}}^j \right] && \text{(from (33))} \\ &= u_{-1} \sum_{i,j=1}^d U_{ij} \mathbb{E} \left[ \left( \widehat{\Lambda}_\tau \right)_{i:} \cdot \left( \widehat{\Lambda}_\tau \right) \cdot U_{11} x_{\tau,\text{query}} x_{\tau,\text{query}}^j \right] && \text{(independence and distribution of } w_\tau) \\ &= u_{-1} \sum_{i,j=1}^d U_{ij} \mathbb{E} \left[ \left( \widehat{\Lambda}_\tau \right)_{i:} \cdot \left( \widehat{\Lambda}_\tau \right) \cdot U_{11} \Lambda_j \right] && \text{(independence between prompt covariates)} \\ &= u_{-1} \mathbb{E} \text{tr} \left[ \sum_{i,j=1}^d \Lambda_j U_{ij} \left( \widehat{\Lambda}_\tau \right)_{i:} \cdot \left( \widehat{\Lambda}_\tau \right) U_{11} \right] = u_{-1} \mathbb{E} \text{tr} \left[ \Lambda (U_{11})^\top \left( \widehat{\Lambda}_\tau \right)^2 U_{11} \right] \\ &= u_{-1} \text{tr} \left[ \mathbb{E} \left( \widehat{\Lambda}_\tau \right)^2 U_{11} \Lambda (U_{11})^\top \right]. \end{aligned}$$

685 For the second term in (37), we have

$$\begin{aligned} \mathbb{E} \left[ \left( (Z_\tau U X_\tau)_{d+1,d+1} u_{-1} \right) (Z_\tau U X_\tau)_{d+1,d+1} \right] &= u_{-1} \mathbb{E} \left[ w_\tau^\top \left( \widehat{\Lambda}_\tau \right) U_{11} x_{\tau,\text{query}} x_{\tau,\text{query}}^\top (U_{11})^\top \left( \widehat{\Lambda}_\tau \right) w_\tau \right] \\ &&& \text{(from (33))} \\ &= u_{-1} \mathbb{E} \text{tr} \left[ w_\tau w_\tau^\top \left( \widehat{\Lambda}_\tau \right) U_{11} x_{\tau,\text{query}} x_{\tau,\text{query}}^\top (U_{11})^\top \left( \widehat{\Lambda}_\tau \right) \right] \\ &= u_{-1} \mathbb{E} \text{tr} \left[ \left( \widehat{\Lambda}_\tau \right) U_{11} \Lambda (U_{11})^\top \left( \widehat{\Lambda}_\tau \right) \right] \\ &= u_{-1} \text{tr} \left[ \mathbb{E} \left( \widehat{\Lambda}_\tau \right)^2 U_{11} \Lambda (U_{11})^\top \right]. \end{aligned}$$

686 Therefore, we know

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

687 Additionally, we have

$$\begin{aligned} 2 \left[ \mathbb{E} \left( w_\tau^\top x_{\tau, \text{query}} H_\tau \right) u \right]_{(d+1)^2} &= \left[ \begin{pmatrix} \mathbf{0}_{d(d+1) \times d(d+1)} & A \\ A^\top & \mathbf{0}_{(d+1) \times (d+1)} \end{pmatrix} \cdot u \right]_{(d+1)^2} \quad (\text{from (31)}) \\ &= \left( V_1 + V_1^\top \quad \dots \quad V_d + V_d^\top \quad \mathbf{0}_{(d+1) \times (d+1)} \right)_{d+1} \cdot U \\ &\quad (\text{definition of } A \text{ in (32)}) \\ &= \sum_{i,j=1}^d \Lambda_i^\top \Lambda_j U_{ji} = \operatorname{tr} \left( \Lambda (U_{11})^\top \Lambda \right). \end{aligned}$$

688 Then, from (29), we have the dynamics of  $u_{-1}$  is

$$\frac{d}{dt} u_{-1} = - \operatorname{tr} \left[ u_{-1} \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - \Lambda^2 (U_{11})^\top \right]. \quad (38)$$

689

□

### 690 D.2.3 Proof of Lemma D.4

691 Lemma D.4 gives the form of global minima of an equivalent loss function. First, we prove that  
692 gradient flow on  $L$  defined in (16) from the initial values satisfying Assumption C.1 is equivalent to  
693 gradient flow on another loss function  $\tilde{\ell}$  defined below. Then, we derive an expression for the global  
694 minima of this loss function.

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

697 **Lemma D.6 (Loss Function).** *Consider gradient flow over  $L$  in (22) with respect to  $u$  starting from  
698 an initial value satisfying Assumption C.1. This is equivalent to doing gradient flow with respect to  
699  $U_{11}$  and  $u_{-1}$  on the loss function*

$$\tilde{\ell}(U_{11}, u_{-1}) = \operatorname{tr} \left[ \frac{1}{2} u_{-1}^2 \Gamma \Lambda U_{11} \Lambda (U_{11})^\top - u_{-1} \Lambda^2 (U_{11})^\top \right]. \quad (39)$$

700 *Proof.* The proof is simply by taking gradient of the loss function in (39). For techniques in matrix  
701 derivatives, see Lemma H.1. We take the gradient of  $\tilde{\ell}$  on  $U_{11}$  to obtain

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

702 since  $\Gamma$  and  $\Lambda$  are commutable. We take derivatives w.r.t.  $u_{-1}$  to get

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

703 Combining this with Lemma D.3, we have

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

704

□

705

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

708 **Corollary D.7** (Minimum of Loss Function). *The loss function  $\tilde{\ell}$  in Lemma D.6 satisfies*

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

709 *and*

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

710 *Proof.* First, we claim that

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

711 To calculate this, we just need to expand the terms in the brackets and notice that  $\Gamma$  and  $\Lambda$  are  
712 commutable:

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

713 Equations (i) and (ii) use that  $\Gamma$  and  $\Lambda$  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 H.4 that

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

714 which implies

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

715 Equality holds when

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

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

718 Lemma D.4 is an immediate consequence of Corollary D.7, since the loss will keep the same when  
719 we replace  $(U_{11}, u_{-1})$  by  $(cU_{11}, c^{-1}u_{-1})$  for any non-zero constant  $c$ .

## 720 D.2.4 Proof of Lemma D.5

721 In this section, we prove that the dynamical system in Lemma D.3 satisfies a PL inequality. Then, the  
722 PL inequality naturally leads to the global convergence of this dynamical system. First, we prove  
723 a simple lemma, which says the parameters in the LSA model will keep 'balanced' in the whole  
724 trajectory. From the proof of this lemma, we can understand why we assume a balanced parameter at  
725 the initial time.

726 **Lemma D.8** (Balanced Parameters). *Consider gradient flow over  $L$  in (22) with respect to  $u$  starting  
727 from an initial value satisfying Assumption C.1. For any  $t \geq 0$ , it holds that*

$$u_{-1}^2 = \operatorname{tr} [U_{11}(U_{11})^\top]. \quad (40)$$

728 *Proof.* From Lemma D.3, we multiply the first equation in (23) by  $(U_{11})^\top$  from the right to get

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

729 Also we multiply the second equation in Lemma D.3 by  $u_{-1}$  to obtain

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

730 Therefore, we have

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

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

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

732 Notice that from Assumption C.1, we know that at  $t = 0$ ,

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

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

734

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

738 **Lemma D.9.** Consider gradient flow over  $L$  in (22) with respect to  $u$  starting from an initial value  
739 satisfying Assumption C.1. If the initial scale satisfies

$$0 < \sigma < \sqrt{\frac{2}{\sqrt{d}\|\Gamma\|_{op}}}, \quad (41)$$

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

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

741 *Proof.* From Lemma D.6, we are actually doing gradient flow on the loss  $\tilde{\ell}$ . The loss function is  
742 non-increasing, because

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

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

746 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  
747 Assumption C.1, we can calculate the loss function at the initial time:

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

748 From the property of trace, we know

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

749 From Von-Neumann's trace inequality (Lemma H.3) and the fact that  $\|\Theta\Theta^\top\|_F = 1$ , we know

$$\text{tr}[\Gamma\Lambda\Theta\Theta^\top\Lambda\Theta\Theta^\top] \leq \sqrt{d}\|\Lambda\Theta\Theta^\top\Lambda\Theta\Theta^\top\|_F \cdot \|\Gamma\|_{op} \leq \sqrt{d}\|\Lambda\Theta\|_F^2 \|\Theta\Theta^\top\|_F \|\Gamma\|_{op} = \sqrt{d}\|\Lambda\Theta\|_F^2 \|\Gamma\|_{op}.$$



750 Therefore, we have

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

751 From Assumption C.1, we know  $\|\Lambda\Theta\|_F \neq 0$ . From (35), we know  $\|\Gamma\|_{op} > 0$ . Therefore, when

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

752 we have

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

753 □

754

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

758 **Lemma D.10.** *Consider gradient flow over  $L$  in (22) with respect to  $u$  starting from an initial value  
759 satisfying Assumption C.1 with initial scale  $0 < \sigma < \sqrt{\frac{2}{\sqrt{d}\|\Gamma\|_{op}}}$ . For any  $t \geq 0$ , it holds that*

$$u_{-1} \geq \sqrt{\frac{\sigma^2}{2\sqrt{d}\|\Lambda\|_{op}^2} \|\Lambda\Theta\|_F^2 \left[ 2 - \sqrt{d}\sigma^2 \|\Gamma\|_{op} \right]} > 0. \quad (42)$$

760 *Proof.* We prove by contradiction. Suppose the claim does not hold. From Lemma D.8, we know  
761  $u_{-1}^2 = \text{tr} [U_{11}(U_{11})^\top] = \|U_{11}\|_F^2$ . From Lemma D.9, we know  $u_{-1} = \|U_{11}\|_F$ . Recall the  
762 definition of loss function:

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

763 Since  $\Gamma \succeq 0, \Lambda \succeq 0$ , and they commute, we know from Lemma H.4 that  $\Gamma \Lambda \succeq 0$ . Again, since  
764  $U_{11} \Lambda (U_{11})^\top = \left( U_{11} \Lambda^{\frac{1}{2}} \right) \left( U_{11} \Lambda^{\frac{1}{2}} \right)^\top \succeq 0$ , from Lemma H.4 we have  $\text{tr} \left[ \frac{1}{2} u_{-1}^2 \Gamma \Lambda U_{11} \Lambda (U_{11})^\top \right] \geq$   
765 0. So

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

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

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

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

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

768 Here, the last inequality comes from the proof of Lemma D.9. This contradicts the non-increasing  
769 property of the loss function in gradient flow. □

770

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

773 **Lemma D.5.** *Suppose the initialization of gradient flow satisfies Assumption C.1 with initialization  
774 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*

$$\mu := \frac{\sigma^2}{\sqrt{d}\|\Lambda\|_{op}^2 \text{tr}(\Gamma^{-1}\Lambda^{-1}) \text{tr}(\Lambda^{-1})} \|\Lambda\Theta\|_F^2 \left[ 2 - \sqrt{d}\sigma^2 \|\Gamma\|_{op} \right] > 0, \quad (26)$$

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

$$\left\| \nabla \tilde{\ell}(U_{11}(t), u_{-1}(t)) \right\|_2^2 := \left\| \frac{\partial \tilde{\ell}}{\partial U_{11}} \right\|_F^2 + \left| \frac{\partial \tilde{\ell}}{\partial u_{-1}} \right|^2 \geq \mu \left( \tilde{\ell}(U_{11}(t), u_{-1}(t)) - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \tilde{\ell}(U_{11}, u_{-1}) \right). \quad (27)$$

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

$$\lim_{t \rightarrow \infty} u_{-1}(t) = \|\Gamma^{-1}\|_F^{-\frac{1}{2}} \text{ and } \lim_{t \rightarrow \infty} U_{11}(t) = \|\Gamma^{-1}\|_F^{-\frac{1}{2}} \Gamma^{-1}. \quad (28)$$

778 *Proof.* From the definition and Lemma D.10, we have

$$\begin{aligned} \|\nabla \ell(U_{11}, u_{-1})\|_2^2 &\geq \left\| \frac{\partial \ell}{\partial U_{11}} \right\|_F^2 = \|u_{-1}^2 \Gamma \Lambda U_{11} \Lambda - u_{-1} \Lambda^2\|_F^2 \\ &= u_{-1}^2 \left\| \Gamma \Lambda^{\frac{1}{2}} \left( u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \Lambda^{\frac{1}{2}} \right\|_F^2 \\ &\geq \frac{\sigma^2}{2\sqrt{d} \|\Lambda\|_{op}^2} \|\Lambda \Theta\|_F^2 \left[ 2 - \sqrt{d} \sigma^2 \|\Gamma\|_{op} \right] \left\| \Gamma \Lambda^{\frac{1}{2}} \left( u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \Lambda^{\frac{1}{2}} \right\|_F^2. \end{aligned} \quad (43)$$

779 To see why the second line is true, recall that  $u_{-1} \in \mathbb{R}$  and  $\Gamma$  and  $\Lambda$  commute. The last line comes  
780 from the lower bound of  $u_{-1}$  in Lemma D.10. From Corollary D.7, we know

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

781 Therefore, we know that

$$\begin{aligned} \ell - \min_{U_{11} \in \mathbb{R}^{d \times d}, u_{-1} \in \mathbb{R}} \ell(U_{11}, u_{-1}) &\leq \frac{1}{2} \left\| \Gamma \Lambda^{\frac{1}{2}} \left( u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \Lambda^{\frac{1}{2}} \right\|_F^2 \cdot \left\| \Gamma^{-\frac{1}{2}} \Lambda^{-\frac{1}{2}} \right\|_F^2 \left\| \Lambda^{-\frac{1}{2}} \right\|_F^2 \\ &= \frac{1}{2} \left\| \Gamma \Lambda^{\frac{1}{2}} \left( u_{-1} \Lambda^{\frac{1}{2}} U_{11} \Lambda^{\frac{1}{2}} - \Lambda \Gamma^{-1} \right) \Lambda^{\frac{1}{2}} \right\|_F^2 \cdot \text{tr}(\Gamma^{-1} \Lambda^{-1}) \text{tr}(\Lambda^{-1}) \end{aligned} \quad (44)$$

782 We compare (43) and (44) to obtain that in order to make the PL condition hold, one needs to let

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

783 Once we set this  $\mu$ , we get the PL inequality. The  $\mu$  is positive due to the assumption for  $\sigma$  in the  
784 lemma.

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

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

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

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

787 which implies

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

788 From Corollary D.7, we know this is

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

789 Since  $\Gamma$  and  $\Lambda$  are non-singular and positive definite, and they commute, we know

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

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

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

791 Therefore, we know

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

792

□

793 **E Theorem 3.2 and the proof**

794 **E.1 Formal statement and discussion**

795 First, we formally state the Theorem 3.2 and provide some discussion about the convergence rate of  
796 generalization risk.

797 **Theorem E.1.** *Let  $\mathcal{D}$  be a distribution over  $(x, y) \in \mathbb{R}^d \times \mathbb{R}$ , whose marginal distribution on  $x$  is  
798  $\mathcal{D}_x = \mathcal{N}(0, \Lambda)$ . 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  
799 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  $f_{\text{LSA}}^*$  be the  
800 LSA model with parameters  $W_*^{PV}$  and  $W_*^{KQ}$  in (19), and  $\hat{y}_{\text{query}}$  is the prediction for  $x_{\text{query}}$  given the  
801 prompt. If we define*

$$a := \Lambda^{-1} \mathbb{E}_{(x,y) \sim \mathcal{D}} [xy], \quad \Sigma := \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ (xy - \mathbb{E}(xy))(xy - \mathbb{E}(xy))^\top \right], \quad (45)$$

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

$$\begin{aligned} \mathbb{E}(\hat{y}_{\text{query}} - y_{\text{query}})^2 &= \underbrace{\min_{w \in \mathbb{R}^d} \mathbb{E}(\langle w, x_{\text{query}} \rangle - y_{\text{query}})^2}_{\text{Error of best linear predictor}} \\ &\quad + \frac{1}{M} \text{tr}[\Sigma \Gamma^{-2} \Lambda] + \frac{1}{N^2} \left[ \|a\|_{\Gamma^{-2}\Lambda^3}^2 + 2 \text{tr}(\Lambda) \|a\|_{\Gamma^{-2}\Lambda^2}^2 + \text{tr}(\Lambda)^2 \|a\|_{\Gamma^{-2}\Lambda}^2 \right], \end{aligned} \quad (46)$$

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

804 Now we make a few remarks on the above theorem before considering particular instances of  $\mathcal{D}$   
805 where we may provide more explicit bounds on the prediction error.

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

813 Let us now consider when  $\mathcal{D}$  corresponds to noiseless linear models, so that for some  $w \in \mathbb{R}^d$ ,  
814 we have  $(x, y) = (x, \langle w, x \rangle)$ , in which case the prediction of the trained transformer is given  
815 by (7). Moreover, a simple calculation shows that the  $\Sigma$  from Theorem E.1 takes the form  
816  $\Sigma = \|w\|_\Lambda^2 \Lambda + \Lambda w w^\top \Lambda$ . Hence Theorem E.1 implies the prediction error for the prompt  
817  $P = (x_1, \langle w, x_1 \rangle, \dots, x_M, \langle w, x_M \rangle, x_{\text{query}})$  is

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

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

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

820 From the upper bound above, we can see the rate w.r.t  $M$  and  $N$  are still at most  $O(1/M)$  and  
821  $O(1/N^2)$  respectively. Moreover, the generalization risk also scales with dimension  $d$ ,  $\text{tr}(\Lambda)$  and  
822 the condition number  $\kappa$ . This suggests that for in-context examples involving covariates of greater  
823 variance, or a more ill-conditioned covariance matrix, the generalization risk will be higher for  
824 the same lengths of training and testing prompts. Putting the above together with Theorem E.1,  
825 Definition 2.1 and Definition 2.2, we get the following corollary.

826 **Corollary E.2.** A transformer with a single linear self-attention layer trained on in-context ex-  
827 amples of functions in  $\{x \mapsto \langle w, x \rangle\}$  w.r.t.  $w \sim \mathbf{N}(0, I_d)$  and  $\mathcal{D}_x = \mathbf{N}(0, \Lambda)$  with gradient flow  
828 on the population loss (16) for initializations satisfying Assumption C.1 converges to the model  
829  $f_{\text{LSA}}(\cdot; W_*^{KQ}, W_*^{PV})$ . This model takes a prompt  $P = (x_1, y_1, \dots, x_M, y_M, x_{\text{query}})$  and returns a  
830 prediction  $\hat{y}_{\text{query}}$  for  $x_{\text{query}}$  given by

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

831 Moreover, the model  $f_{\text{LSA}}(\cdot; W_*^{KQ}, W_*^{PV})$  in-context learns the class of linear models  $\{x \mapsto \langle w, x \rangle\}$   
832 with respect to  $w \sim \mathbf{N}(0, I_d)$  and  $\mathcal{D}_x = \mathbf{N}(0, \Lambda)$ , provided  $M \geq 2(d+1) \text{tr}(\Lambda) \varepsilon^{-1}$  and the prompts  
833 seem during training were of length at least  $N \geq \sqrt{2(1+2d+d^2\kappa) \text{tr}(\Lambda) \varepsilon^{-1/2}}$ , where  $\kappa$  is the  
834 condition number of  $\Lambda$ .

## 835 E.2 Proof of Theorem E.1

836 *Proof.* Unless otherwise specified, we denote  $\mathbb{E}$  as the expectation over  $(x_i, y_i), (x_{\text{query}}, y_{\text{query}})$  <sup>i.i.d.</sup>  
837  $\mathcal{D}$ . 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  
838  $\mathbb{E}(\langle w, x_{\text{query}} \rangle - y_{\text{query}})^2$  exists for each  $w \in \mathbb{R}^d$ . We denote

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

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

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

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

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

844 We also define a similar vector for ease of computation:

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

845 Therefore, we can decompose the risk as

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

846 The term I is the first term on the right hand side of (46). So it suffices to calculate II to VI.

847

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

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

849 since

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

850

851 Similarly, for IV, we have

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

852

853 For VI, we have

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

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

856

857 For II, from the definition we have

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

858 The last line comes from the definition of  $\Sigma$ .

859

860 For III, we have

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

861 Combining all terms above, we conclude. □

## 862 F Transformers trained on prompts with random covariate distributions

### 863 F.1 Main theorem for the random covariance case

864 In this section, we will consider a variant of training on in-context examples (in the sense of  
 865 Definition 2.1) where the distribution  $\mathcal{D}_x$  is itself sampled randomly from a distribution, and training  
 866 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  $h \sim \mathcal{D}_{\mathcal{H}}$ .  
 867 More formally, we can generalize Definition 2.1 as follows.

868 **Definition F.1** (Trained on in-context examples with random covariate distributions). *Let  $\Delta$  be  
 869 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  
 870  $\mathcal{X} \rightarrow \mathcal{Y}$ , 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  
 871  $\mathcal{S} = \cup_{n \in \mathbb{N}} \{(x_1, y_1, \dots, x_n, y_n) : x_i \in \mathcal{X}, y_i \in \mathcal{Y}\}$  be the set of finite-length sequences of  $(x, y)$   
 872 pairs and let*

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

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

$$\theta^* \in \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}_{P=(x_1, h(x_1), \dots, x_N, h(x_N), x_{\text{query}})} [\ell(f_{\theta}(P), h(x_{\text{query}}))], \quad (49)$$

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

877 We recover the previous definition of training on in-context examples by taking  $\Delta$  to be concentrated  
 878 on a singleton,  $\operatorname{supp}(\Delta) = \{\mathcal{D}_x\}$ . The natural question is then, if a model  $f$  is trained on in-context  
 879 examples from a function class  $\mathcal{H}$  w.r.t.  $\mathcal{D}_{\mathcal{H}}$  and a distribution  $\Delta$  over covariate distributions, and if  
 880 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  
 881 that  $\mathcal{D}_x$  (cf. Definition 2.2)? Since  $\mathcal{D}_x$  is random, we can hope that this may hold in expectation or  
 882 with high probability over the sampling of the covariate distribution. In the remainder of this section,  
 883 we will explore this question for transformers with a linear self-attention layer trained by gradient  
 884 flow on the population loss.

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

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

892 Thus the diagonal entries of  $\Lambda_{\tau}$  are independent but could have different distributions, and  $\Lambda_{\tau}$  is  
 893 identically distributed for  $\tau = 1, \dots, B$ . Then, conditional on  $\Lambda_{\tau}$ , we sample independent and  
 894 identically distributed  $x_{\tau, 1}, \dots, x_{\tau, N}, x_{\tau, \text{query}} \sim \mathcal{N}(0, \Lambda_{\tau})$ . A training prompt is then given by  
 895  $P_{\tau} = (x_{\tau, 1}, \langle w_{\tau}, x_{\tau, 1} \rangle, \dots, x_{\tau, N}, \langle w_{\tau}, x_{\tau, N} \rangle, x_{\tau, \text{query}})$  Notice that here,  $x_{\tau, i}, x_{\tau, \text{query}}$  are condition-  
 896 ally independent given the covariance matrix  $\Lambda_{\tau}$ , but not independent in general. We consider the  
 897 same token embedding matrix as (3) and linear self-attention network, which forms the prediction  
 898  $\hat{y}_{\text{query}, \tau}$  as in (14). The empirical risk is the same as before (see (15)), and as in (16), we then take  
 899  $B \rightarrow \infty$  and consider the gradient flow on the population loss. The population loss now includes an  
 900 expectation over the distribution of the covariance matrices in addition to the task weight  $w_{\tau}$  and  
 901 covariate distributions, and is given by

$$L(\theta) = \frac{1}{2} \mathbb{E}_{w_{\tau}, \Lambda_{\tau}, x_{\tau, 1}, \dots, x_{\tau, N}, x_{\tau, \text{query}}} [(\hat{y}_{\tau, \text{query}} - \langle w_{\tau}, x_{\tau, \text{query}} \rangle)^2]. \quad (50)$$

902

903 In the main result for this section, we show that gradient flow with a suitable initialization converges to  
 904 a global minimum, and we characterize the limiting solution. The proof will be deferred to Appendix  
 905 F.2.

906 **Theorem F.2** (Global convergence in random covariance case). *Consider gradient flow of the linear*  
 907 *self-attention network  $f_{\text{LSA}}$  defined in (3) over the population loss (50), where  $\Lambda_\tau$  are diagonal with*  
 908 *independent diagonal entries which are strictly positive a.s. and have finite third moments. Suppose*  
 909 *the initialization satisfies Assumption C.1,  $\|\mathbb{E}\Lambda_\tau\Theta\|_F \neq 0$ , with initialization scale  $\sigma > 0$  satisfying*

$$\sigma^2 < \frac{2 \|\mathbb{E}\Lambda_\tau\Theta\|_F^2}{\sqrt{d} \left[ \mathbb{E} \|\Gamma_\tau\|_{\text{op}} \|\Lambda_\tau\|_F^2 \right]}. \quad (51)$$

910 *Then gradient flow converges to a global minimum of the population loss (50). Moreover,  $W^{PV}$  and*  
 911  *$W^{KQ}$  converge to  $W_*^{PV}$  and  $W_*^{KQ}$  respectively, where*

$$W_*^{KQ} = \left\| [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} \mathbb{E}[\Lambda_\tau^2] \right\|_F^{-\frac{1}{2}} \cdot \begin{pmatrix} [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} [\mathbb{E}\Lambda_\tau^2] & 0_d \\ 0_d^\top & 0 \end{pmatrix},$$

$$W_*^{PV} = \left\| [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} \mathbb{E}[\Lambda_\tau^2] \right\|_F^{\frac{1}{2}} \cdot \begin{pmatrix} 0_{d \times d} & 0_d \\ 0_d^\top & 1 \end{pmatrix}, \quad (52)$$

912 *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*  
 913  *$\Lambda_\tau$ .*

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

$$\begin{aligned} & \widehat{y}_{\text{query}} \quad (53) \\ &= \begin{pmatrix} 0_d^\top & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{M} \sum_{i=1}^M x_i x_i^\top + \frac{1}{M} x_{\text{query}} x_{\text{query}}^\top & \frac{1}{M} \sum_{i=1}^M x_i y_i \\ \frac{1}{M} \sum_{i=1}^M x_i^\top y_i & \frac{1}{M} \sum_{i=1}^M y_i^2 \end{pmatrix} \begin{pmatrix} [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} [\mathbb{E}\Lambda_\tau^2] & 0_d \\ 0_d^\top & 0 \end{pmatrix} \begin{pmatrix} x_{\text{query}} \\ 0 \end{pmatrix} \\ &= x_{\text{query}}^\top \cdot [\mathbb{E}\Lambda_\tau^2] [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} \cdot \left[ \frac{1}{M} \sum_{i=1}^M x_i x_i^\top \right] w \\ &\rightarrow x_{\text{query}}^\top \cdot [\mathbb{E}\Lambda_\tau^2] [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} \cdot \Lambda_{\text{new}} w \quad \text{almost surely when } M \rightarrow \infty. \quad (54) \end{aligned}$$

920 The last line comes from the strong law of large numbers. Thus, in order for the prediction on the  
 921 query example to be close to the ground-truth  $x_{\text{query}}^\top w$ , we need  $[\mathbb{E}\Lambda_\tau^2] [\mathbb{E}\Gamma_\tau\Lambda_\tau^2]^{-1} \cdot \Lambda_{\text{new}}$  to be  
 922 close to the identity. When  $\Lambda_\tau \equiv \Lambda_{\text{new}}$  is deterministic, this indeed is the case as we know from  
 923 Theorem E.1. However, this clearly does not hold in general when  $\Lambda_\tau$  is random.

924 To make things concrete, let us assume for simplicity that  $M, N \rightarrow \infty$  so that  $\Gamma_\tau \rightarrow \Lambda_\tau$  and the  
 925 identity (54) holds (conditionally on  $\Lambda_{\text{new}}$ ). Then, taking expectation over  $\Lambda_{\text{new}}$  in (54), we obtain

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

926 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  $\mathbb{E}[\Lambda_\tau^3] = 6I_d$ ,  
 927 we get

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

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

## 931 F.2 Proof of Theorem F.2

932 The proof of Theorem F.2 is very similar to that of Theorem D.1. The first step is to explicitly write  
 933 out the dynamical system. In order to do so, we notice that the Lemma D.2 does not depend on



934 the training data and data-generating distribution and hence, it still holds in the case of a random  
 935 covariance matrix. Therefore, we know when we input the embedding matrix  $E_\tau$  to the linear  
 936 self-attention layer with parameter  $\theta = (W^{KQ}, W^{PV})$ , the prediction will be

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

937 where the matrix  $H_\tau$  is defined as,

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

938 and

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

939 where  $U_{11} = W_{11}^{KQ} \in \mathbb{R}^{d \times d}$ ,  $u_{12} = w_{21}^{PV} \in \mathbb{R}^{d \times 1}$ ,  $u_{21} = w_{21}^{KQ} \in \mathbb{R}^{d \times 1}$ ,  $u_{-1} = w_{22}^{PV} \in \mathbb{R}$   
 940 correspond to particular components of  $W^{PV}$  and  $W^{KQ}$ , defined in (13).

941

## 942 F.2.1 Dynamical system

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

947 **Lemma F.3.** *Consider gradient flow on (50) with respect to  $u$  starting from an initial value that*  
 948 *satisfies Assumption C.1. We assume the covariance matrices  $\Lambda_\tau$  are sampled from some distribution*  
 949 *with finite third moment and  $\Lambda_\tau$  are positive definite almost surely. We denote  $u = \text{Vec}(U) :=$*

950  $\text{Vec} \begin{pmatrix} U_{11} & u_{12} \\ (u_{21})^\top & u_{-1} \end{pmatrix}$  *and define*

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

951 *Then the dynamics of  $U$  follows*

$$\begin{aligned} \frac{d}{dt} U_{11}(t) &= -u_{-1}^2 \mathbb{E} [\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau] + u_{-1} \mathbb{E} [\Lambda_\tau^2] \\ \frac{d}{dt} u_{-1}(t) &= -u_{-1} \text{tr} \mathbb{E} [\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top] + \text{tr} (\mathbb{E} [\Lambda_\tau^2] (U_{11})^\top), \end{aligned} \quad (55)$$

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

953 *Proof.* This lemma is a natural corollary of Lemma D.3. Notice that Lemma D.3 holds for any fixed  
 954 positive definite  $\Lambda_\tau$ . So when  $\Lambda_\tau$  is random, if we condition on  $\Lambda_\tau$ , the dynamical system will be

$$\begin{aligned} \frac{d}{dt} U_{11}(t) &= -u_{-1}^2 [\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau] + u_{-1} [\Lambda_\tau^2] \\ \frac{d}{dt} u_{-1}(t) &= -u_{-1} \text{tr} [\Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top] + \text{tr} ([\Lambda_\tau^2] (U_{11})^\top), \end{aligned} \quad (56)$$

955 and  $u_{12}(t) = 0_d$ ,  $u_{21}(t) = 0_d$  for all  $t \geq 0$ . Then, we conclude by simply taking expectation over  
 956  $\Lambda_\tau$ .  $\square$

957

958 The lemma above gives the dynamical system with general random covariance matrix. When  $\Lambda_\tau$  are  
 959 diagonal almost surely, we can actually simplify the dynamical system above. In this case, we have  
 960 the following corollary.

961 **Corollary F.4.** Under the assumptions of Lemma F.3, we further assume the covariance matrix  $\Lambda_\tau$  to  
 962 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

$$\begin{aligned}\gamma_i &= \mathbb{E} \left[ \frac{N+1}{N} \lambda_{\tau,i}^3 + \frac{1}{N} \lambda_{\tau,i}^2 \cdot \sum_{j=1}^d \lambda_{\tau,j} \right], \\ \xi_i &= \mathbb{E} [\lambda_{\tau,i}^2], \\ \zeta_{ij} &= \mathbb{E} \left[ \frac{N+1}{N} \lambda_{\tau,i}^2 \lambda_{\tau,j} + \frac{1}{N} \lambda_{\tau,i} \lambda_{\tau,j} \cdot \sum_{k=1}^d \lambda_{\tau,k} \right]\end{aligned}\tag{57}$$

963 for  $i, j \in [d]$ , where the expectation is over the distribution of  $\Lambda_\tau$ . Then, the dynamical system (55) is  
 964 equivalent to

$$\begin{aligned}\frac{d}{dt} u_{ii}(t) &= -\gamma_i u_{-1}^2 u_{ii} + \xi_i u_{-1} \quad \forall i \in [d], \\ \frac{d}{dt} u_{ij}(t) &= -\zeta_{ij} u_{-1}^2 u_{ij} \quad \forall i \neq j \in [d], \\ \frac{d}{dt} u_{-1}(t) &= -\sum_{i=1}^d [\gamma_i u_{-1} u_{ii}^2] - \sum_{i \neq j} \zeta_{ij} u_{-1} u_{ij}^2 + \sum_{i=1}^d [\xi_i u_{ii}].\end{aligned}\tag{58}$$

965 *Proof.* This is directly obtained by rewriting the equation for each entry of  $U_{11}$  and recalling the  
 966 assumption that  $\Lambda_\tau$  (and hence  $\Gamma_\tau$ ) is diagonal almost surely.  $\square$

## 967 F.2.2 Loss function and global minima

968 As in the proof of Theorem D.1, we can actually recover the loss function in the random covariance  
 969 case, up to a constant.

970 **Lemma F.5.** The differential equations in (58) are equivalent to gradient flow on the loss function

$$\begin{aligned}\ell_{\text{rdm}}(U_{11}, u_{-1}) &= \mathbb{E} \text{tr} \left[ \frac{1}{2} u_{-1}^2 \Gamma_\tau \Lambda_\tau U_{11} \Lambda_\tau (U_{11})^\top - u_{-1} \Lambda_\tau^2 (U_{11})^\top \right] \\ &= \frac{1}{2} \sum_{i=1}^d [\gamma_i u_{-1}^2 u_{ii}^2] + \frac{1}{2} \sum_{i \neq j} \zeta_{ij} u_{-1}^2 u_{ij}^2 - \sum_{i=1}^d [\xi_i u_{ii} u_{-1}]\end{aligned}\tag{59}$$

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

972 *Proof.* This can be verified by simply taking gradient of  $\ell_{\text{rdm}}$  to show that

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

973  $\square$

974

975 Next, we solve for the minimum of  $\ell_{\text{rdm}}$  and give the expression for all global minima.

976 **Lemma F.6.** Let  $\ell_{\text{rdm}}$  be the loss function in (59). We denote

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

977 Then, we have

$$\min \ell_{\text{rdm}} = -\frac{1}{2} \sum_{i=1}^d \frac{\xi_i^2}{\gamma_i}\tag{60}$$

978 and

$$\ell_{\text{rdm}}(U_{11}, u_{-1}) - \min \ell_{\text{rdm}} = \frac{1}{2} \sum_{i=1}^d \gamma_i \left( u_{ii} u_{-1} - \frac{\xi_i}{\gamma_i} \right)^2 + \frac{1}{2} \sum_{i \neq j} \zeta_{ij} u_{-1}^2 u_{ij}^2.\tag{61}$$

979 Moreover, denoting  $u_{ij}$  as the  $(i, j)$ -entry of  $U_{11}$ , all global minima of  $\ell_{\text{rdm}}$  satisfy

$$u_{-1} \cdot u_{ij} = \mathbb{I}(i = j) \cdot \frac{\xi_i}{\gamma_i}. \quad (62)$$

980 *Proof.* From the definition of  $\ell_{\text{rdm}}$ , we have

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

981 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  
 982 achieved by simply letting  $u_{-1} = 1$  and  $u_{ii} = \frac{\xi_i}{\gamma_i}$  for  $i \in [d]$ . Of course, when we replace  $(u_{-1}, u_{ii})$   
 983 with  $(cu_{-1}, c^{-1}u_{ii})$  for any constant  $c \neq 0$ , we can also achieve this global minimum.  $\square$

### 984 F.2.3 PL Inequality and global convergence

985 Finally, to end the proof, we prove a Polyak-Łojasiewicz Inequality on the loss function  $\ell_{\text{rdm}}$ , and  
 986 then prove global convergence. Before that, let's first prove the balanced condition of parameters will  
 987 hold during the whole trajectory.

988 **Lemma F.7** (Balanced condition). *Under the assumptions of Lemma F.3, for any  $t \geq 0$ , it holds that*

$$u_{-1}^2 = \text{tr} [U_{11}(U_{11})^\top]. \quad (63)$$

989 *Proof.* The proof is similar to the proof of Lemma D.8. From Lemma D.3, we multiply the first  
 990 equation in (55) by  $(U_{11})^\top$  from the right to get

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

991 Also we multiply the second equation in Lemma 55 by  $u_{-1}$  to obtain

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

992 Therefore, we have

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

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

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

994 Notice that from Assumption C.1, we know that

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

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

996

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

1000 **Lemma F.8.** *We do gradient flow on  $\ell_{\text{rdm}}$  with respect to  $u_{i,j}$  ( $\forall i, j \in [d]$ ) and  $u_{-1}$ . Suppose the  
 1001 initialization satisfies Assumption C.1 with initial scale*

$$0 < \sigma < \sqrt{\frac{2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2}{\sqrt{d} \left[ \mathbb{E} \|\Gamma_\tau\|_{\text{op}} \|\Lambda_\tau\|_F^2 \right]}}, \quad (64)$$

1002 then for any  $t \geq 0$ , it holds that

$$u_{-1}(t) > 0. \quad (65)$$

1003 *Proof.* From the dynamics of gradient flow, we know the loss function  $\ell_{\text{rdm}}$  is non-increasing:

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

1004 Since we assume  $U_{11}(0) = \Theta\Theta^\top$ , we know the loss function at  $t = 0$  is

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

1005 From the property of trace, we know

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

1006 From Von-Neumann's trace inequality and the assumption that  $\|\Theta\Theta^\top\|_F = 1$ , we know

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

1007 From the assumptions on  $\Theta$  and  $\Lambda_\tau$  we know  $\mathbb{E} \Lambda_\tau \Theta \neq 0_{d \times d}$  and  $\mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 > 0$ . Therefore,  
 1008 comparing the two displays above, we know when (64) holds, we must have  $\ell_{\text{rdm}}(0) < 0$ . So from  
 1009 the non-increasing property of the loss function, we know  $\ell_{\text{rdm}}(t) < 0$  for any time  $t \geq 0$ . Notice that  
 1010 when  $u_{-1} = 0$ , the loss function is also zero, which suggests that  $u_{-1}(t) \neq 0$  for any time  $t \geq 0$ .  
 1011 Since  $u_{-1}(0) > 0$  and the trajectory of  $u_{-1}$  must be continuous, we know that it stays positive at all  
 1012 times.  $\square$

1013

1014 **Lemma F.9.** *We do gradient flow on  $\ell_{\text{rdm}}$  with respect to  $u_{i,j}$  ( $\forall i, j \in [d]$ ) and  $u_{-1}$ . Suppose the*  
 1015 *initialization satisfies Assumption C.1 and the initial scale satisfies (64). Then, for any  $t \geq 0$ , it holds*  
 1016 *that*

$$u_{-1}(t) \geq \sqrt{\frac{\sigma^2}{2\sqrt{d} \|\mathbb{E} \Lambda_\tau^2\|_{op}}} \left[ 2 \|\mathbb{E} \Lambda_\tau \Theta\|_F^2 - \sqrt{d} \sigma^2 \left[ \mathbb{E} \|\Gamma_\tau\|_{op} \|\Lambda_\tau\|_F^2 \right] \right] > 0. \quad (66)$$

1017 *Proof.* From the dynamics of gradient flow, we know  $\ell_{\text{rdm}}$  is non-increasing (see the proof of Lemma  
 1018 F.8). Recall the definition of the loss function:

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

1019 Since  $\Lambda_\tau$  commutes with  $\Gamma_\tau$  and they are both positive definite almost surely, we know that  $\Gamma_\tau \Lambda_\tau \succeq$   
 1020  $0_{d \times d}$  almost surely from Lemma H.1. Again, since  $U_{11} \Lambda_\tau (U_{11})^\top \succeq 0_{d \times d}$  almost surely, from  
 1021 Lemma H.1 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

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

1022 From Von Neumann's trace inequality (Lemma H.3) and the fact that  $u_{-1}(t) > 0$  for any  $t \geq 0$   
 1023 (Lemma F.8), 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 F.7,  
 1024 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  
 1025  $u_{-1}(t) = \|U_{11}(t)\|_F$ . So we have

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

1026 From the proof of Lemma F.8, we know

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

1027 Combine the two preceding displays above, we have

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

1028 The last inequality comes from Lemma F.8.  $\square$

1029

1030 Finally, we prove the PL Inequality, which naturally leads to the global convergence.

1031 **Lemma F.10.** *We do gradient flow on  $\ell_{\text{rdm}}$  with respect to  $u_{i,j}$  ( $\forall i, j \in [d]$ ) and  $u_{-1}$ . Suppose the*  
 1032 *initialization satisfies Assumption C.1 and the initial scale satisfies (64). If we denote*

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

1033 and

$$\nu := \frac{\eta \cdot \sigma^2}{2\sqrt{d} \|\mathbb{E}\Lambda_\tau^2\|_{\text{op}}} \left[ 2 \|\mathbb{E}\Lambda_\tau \Theta\|_F^2 - \sqrt{d} \sigma^2 \left[ \mathbb{E} \|\Gamma_\tau\|_{\text{op}} \|\Lambda_\tau\|_F^2 \right] \right] > 0, \quad (67)$$

1034 then for any  $t \geq 0$ , it holds that

$$\|\nabla \ell_{\text{rdm}}(U_{11}, u_{-1})\|_2^2 := \sum_{i,j=1}^d \left| \frac{\partial \ell_{\text{rdm}}}{\partial u_{ij}} \right|^2 + \left| \frac{\partial \ell_{\text{rdm}}}{\partial u_{-1}} \right|^2 \geq \nu (\ell_{\text{rdm}} - \min \ell_{\text{rdm}}). \quad (68)$$

1035 Additionally,  $\ell_{\text{rdm}}$  converges to the global minimal value,  $u_{ij}$  and  $u_{-1}$  converge to the following  
 1036 limits,

$$\lim_{t \rightarrow \infty} u_{ij}(t) = \mathbb{I}(i=j) \cdot \left[ \sum_{i=1}^d \frac{\xi_i^2}{\gamma_i^2} \right]^{-\frac{1}{4}} \cdot \frac{\xi_i}{\gamma_i} \quad \forall i \in [d], \quad \lim_{t \rightarrow \infty} u_{-1}(t) = \left[ \sum_{i=1}^d \frac{\xi_i}{\gamma_i} \right]^{\frac{1}{4}}. \quad (69)$$

1037 Translating back to the original parameterization, we have this is equivalent to

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

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

1038 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  $\Lambda_\tau$ .

1039 *Proof.* First, we prove the PL Inequality. From Lemma F.6, we know

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

1040 where  $\xi_i, \zeta_{ij}, \gamma_i$  are defined in (57). Meanwhile, we calculate the square norm of the gradient of  $\ell_{\text{rdm}}$ :

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

1041 Comparing the two displays above, we know in order to achieve  $\|\nabla \ell_{\text{rdm}}\|_2^2 \geq \nu (\ell_{\text{rdm}} - \min \ell_{\text{rdm}})$ ,  
 1042 it suffices to make

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

1043 We define  $\eta := \min \{ \gamma_i, \zeta_{ij}, i \neq j \in [d] \}$ , then it is sufficient to make

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

1044 From Lemma F.9, we know that we can actually lower bound  $u_{-1}$  from below by a positive constant.  
 1045 Then, the inequality holds if we take

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

1046 Therefore, as long as we take  $\nu$  as above, a PL inequality holds for  $\ell_{\text{rdm}}$ .

1047 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  
 1048 gradient flow and the PL Inequality ((68)), we know

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

1049 which by Grönwall's inequality implies

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

1050 when  $t \rightarrow \infty$ . From Lemma F.6, we know

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

1051 This implies

$$\begin{aligned} u_{ii} u_{-1} &\rightarrow \frac{\xi_i}{\gamma_i} \quad \forall i \in [d], \\ u_{ij} u_{-1} &\rightarrow 0 \quad \forall i \neq j \in [d]. \end{aligned} \tag{70}$$

1052 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  
 1053 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 F.7, we know for any  $t \geq 0$ ,  $u_{-1}(t)^2 =$   
 1054  $\text{tr}(U_{11}(U_{11})^\top) = \sum_{i,j=1}^d u_{ij}^2$ . So we have

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

1055 which implies

$$u_{-1}(t) \rightarrow \left[ \sum_{i=1}^d \frac{\xi_i^2}{\gamma_i^2} \right]^{\frac{1}{4}} \tag{71}$$

1056 when  $t \rightarrow \infty$ . Combining (70) and (71), we conclude

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

1057 □

## 1058 G Experiments with large, nonlinear transformers

1059 We have shown that even when trained on prompts with random covariance matrices, transformers  
 1060 with a single linear self-attention layer fail to in-context learn linear models with random covariance  
 1061 matrices. We now investigate the behavior of more complex transformer architectures that are  
 1062 trained on in-context examples of linear models, both in the fixed-covariance case and in the random-  
 1063 covariance case.

1064 We examine the performance of transformers with a GPT2 architecture [Radford et al., 2019] that are  
 1065 trained on linear regression tasks with mean-zero Gaussian features with either a fixed covariance  
 1066 matrix or random covariance matrices. For the fixed covariance case, the covariance matrix is fixed  
 1067 to the identity matrix across prompts. For the random covariance case, covariates are drawn from  
 1068  $x \sim \mathcal{N}(0, c\Lambda)$  where  $\Lambda$  is diagonal with  $\lambda_i \stackrel{\text{i.i.d.}}{\sim} \text{Exponential}(1)$  and  $c > 0$  is a scaling factor. We  
 1069 set  $c = 1$  during training and vary this value at test time. The transformer is trained using the  
 1070 procedure of Garg et al. [2022] (see Appendix G for more details). We consider linear models in  
 1071  $d = 20$  dimensions and we train on prompt lengths of  $N = 40, 70, 100$  with either fixed or random  
 1072 covariance matrices. The performance of these trained models, when tested on new data with fixed  
 1073 covariance or random covariance matrices ( $c = 1, 4, 9$ ), is represented in six curves in Figure 2.  
 1074 Using the calculation (54), we can compare the prediction error for the linear self-attention networks  
 1075 in the  $M \rightarrow \infty, N \rightarrow \infty$  limit (the black dash line) to those of GPT2 architectures. We additionally  
 1076 compare these models to the ordinary least-squares solution which is optimal for this task.

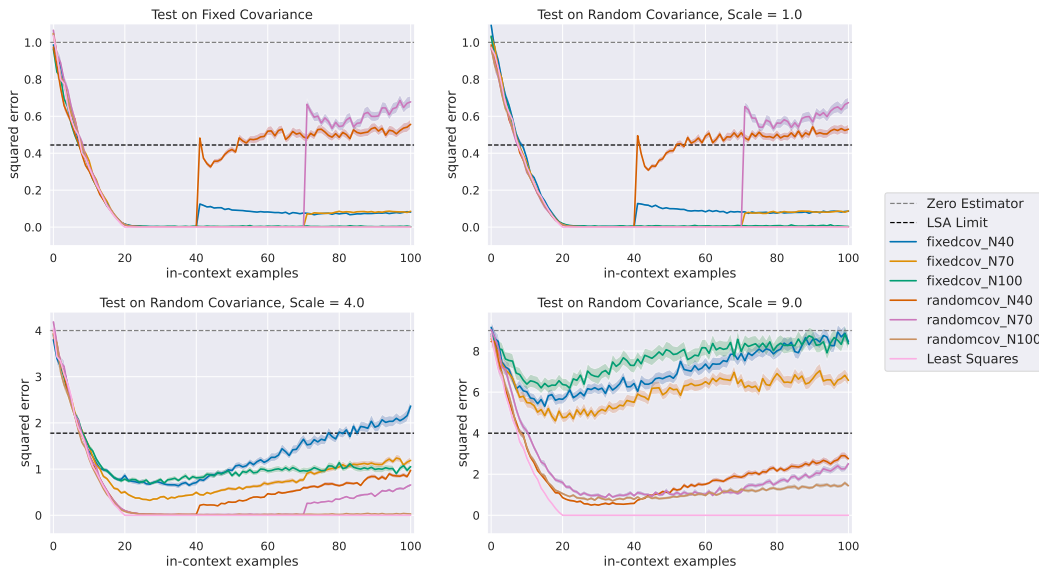


Figure 2: Normalized prediction error for transformers with GPT2 architectures as a function of the number of in-context test examples  $M$  when trained on in-context examples of linear models in  $d = 20$  dimensions. Colored lines correspond to different training context lengths ( $N \in \{40, 70, 100\}$ ) 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.

### 1077 G.1 Experiment details

1078 Here we provide more details for the experiment in Figure 2. Our experimental setup is based on  
 1079 the codebase provided by Garg et al. [2022], with a modification that allows for the possibility that  
 1080 the covariate distribution changes across prompts. We use the standard GPT2 architecture with 256  
 1081 embedding size, 12 layers and 8 heads [Radford et al., 2018] as implemented by HuggingFace [Wolf  
 1082 et al., 2020]. For the GPT2 models, we use the embedding method proposed by Garg et al. [2022],  
 1083 where instead of concatenating  $x$  and  $y$  into a single token, they are treated as separate tokens. It

1084 is also worth noting that the training objective function for the GPT2 model is different than those  
 1085 we consider for the linear self-attention network: for the GPT2 model, the objective function is the  
 1086 average over the full length of the context sequence (predictions for each  $x_i$  using  $(x_k, y_k)_{k < i}$ ), while  
 1087 in our setting the objective function is only for the final query point. However, in the figure, for both  
 1088 GPT2 and the linear self-attention model the error plotted corresponds to the error for predicting the  
 1089 final query point.

1090 In all experiments, covariates are sampled from a mean-zero Gaussian in  $d = 20$  dimensions with  
 1091 either fixed or random covariance matrix. For the fixed covariance case, we fix the covariance matrix  
 1092 to be identity; for the random case, the covariance matrices are restricted to be diagonal and all  
 1093 diagonal entries are i.i.d. sampled from the standard exponential distribution. The linear weights  
 1094 in all tasks are i.i.d. sampled from standard Gaussian distribution and also independently from all  
 1095 covariates. We trained the model for 500000 steps using Adam [Kingma and Ba, 2014] with a batch  
 1096 size of 64 and learning rate of 0.0001. We use the same curriculum strategy of Garg et al. [2022] for  
 1097 acceleration.

1098 For testing the trained model, we used ordinary least squares as a baseline which is optimal for  
 1099 noiseless linear regression tasks. For prompts at test time, covariates are sampled i.i.d. from a mean-  
 1100 zero Gaussian distribution. For the fixed-covariance evaluation, the covariance is the identity matrix.  
 1101 In the random-covariance evaluation, the covariance is a random diagonal matrix with diagonal entries  
 1102 sampled from the standard exponential distribution, multiplied by a scaling coefficient  $c \in \{1, 4, 9\}$ ,  
 1103 i.e. for each task  $\tau$ , the covariance matrix in the random case is

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

1104 where  $\lambda_{\tau,i} \stackrel{\text{i.i.d.}}{\sim} \text{Exponential}(1)$  for any  $\tau$  and  $i \in [d]$ . The plots in Figure 2 show the error averaged  
 1105 over  $64^2$  prompts, where we sample 64 covariance matrices for each curve and 64 prompts for each  
 1106 covariance matrix. We compute 90% confidence interval over 1000 bootstrap trials for each test.

1107 From the figure, we can see that the GPT2 model trained on fixed covariance succeeds in the random  
 1108 covariance setting if the variance is not too large, which shows that the larger nonlinear model is  
 1109 able to generalize better than the model with a single linear self-attention layer. However, when the  
 1110 variance is large ( $c = 4, 9$  for the bottom two figures), the GPT2 model trained with fixed covariance  
 1111 is unsuccessful. When trained on random covariance, the model performs better for test prompts from  
 1112 higher-variance random covariance matrices, but still fails to match least squares when the scaling is  
 1113 largest ( $c = 9$ ).

1114 Furthermore, we notice some surprising behaviors when the test prompt length exceeds the training  
 1115 prompt length (i.e.,  $M > N$ ): there is an evident spike in prediction error, regardless of whether  
 1116 training and testing were performed on fixed or random covariance, and the spike appears to decrease  
 1117 when evaluated on prompts with higher variance. Although we are unsure of why the spike should  
 1118 decrease with higher-variance prompts, the failure of large language models to generalize to larger  
 1119 contexts than they were trained on is a well-known problem [Dai et al., 2019, Anil et al., 2022]. In  
 1120 our setting, we conjecture that this spike in error comes from the absolute positional encodings in the  
 1121 GPT2 architecture. The positional encodings are randomly-initialized and are learnable parameters  
 1122 but the encoding for position  $i$  is only updated if the transformer encounters a prompt which has a  
 1123 context of length  $i$ . Thus, when evaluating on prompts of length  $M > N$ , the model is relying upon  
 1124 random positional encodings for  $M - N$  samples. We note that a concurrent work has explored  
 1125 the performance of transformers with GPT2 architectures for in-context learning of linear models  
 1126 and found that removing positional encoders improves performance when evaluating on larger  
 1127 contexts [Ahuja et al., 2023]. We leave further investigation of this behavior for future work.



## 1128 H Technical lemmas

1129 **Lemma H.1** (Matrix Derivatives, Kronecker Product and Vectorization, [Petersen et al., 2008]). We  
 1130 denote  $\mathbf{A}, \mathbf{B}, \mathbf{X}$  as matrices and  $\mathbf{x}$  as vectors. Then, we have

- 1131 •  $\frac{\partial \mathbf{x}^\top \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^\top) \mathbf{x}$ .
- 1132 •  $\text{Vec}(\mathbf{A} \mathbf{X} \mathbf{B}) = (\mathbf{B}^\top \otimes \mathbf{A}) \text{Vec}(\mathbf{X})$ .
- 1133 •  $\text{tr}(\mathbf{A}^\top \mathbf{B}) = \text{Vec}(\mathbf{A})^\top \text{Vec}(\mathbf{B})$ .
- 1134 •  $\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{X} \mathbf{B} \mathbf{X}^\top) = \mathbf{X} \mathbf{B}^\top + \mathbf{X} \mathbf{B}$ .
- 1135 •  $\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{A} \mathbf{X}^\top) = \mathbf{A}$ .
- 1136 •  $\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{A} \mathbf{X} \mathbf{B} \mathbf{X}^\top \mathbf{C}) = \mathbf{A}^\top \mathbf{C}^\top \mathbf{X} \mathbf{B}^\top + \mathbf{C} \mathbf{A} \mathbf{X} \mathbf{B}$ .

1137  
 1138 **Lemma H.2.** If  $X$  is Gaussian random vector of  $d$  dimension, mean zero and covariance matrix  $\Lambda$ ,  
 1139 and  $A \in \mathbb{R}^{d \times d}$  is a fixed matrix. Then

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

1140 *Proof.* We denote  $X = (X_1, \dots, X_d)^\top$ . Then,

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

1141 So we know  $(X X^\top A X X^\top)_{k,l} = \left( \sum_{i,j=1}^d A_{ij} X_i X_j \right) X_k X_l$ . From Isserlis' Theorem in probability  
 1142 theory (Theorem 1.1 in Michalowicz et al. [2009], originally proposed in Wick [1950]), we know for  
 1143 any  $i, j, k, l \in [d]$ , it holds that

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

1144 Then, we have for any fixed  $k, l \in [d]$ ,

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

1145 Therefore, we know

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

1146

□

1147

**Lemma H.3** (Von-Neumann's Trace Inequality). Let  $U, V \in \mathbb{R}^{d \times n}$  with  $d \leq n$ . We have

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

1148 where  $\sigma_1(X) \geq \sigma_2(X) \geq \dots \geq \sigma_d(X)$  are the ordered singular values of  $X \in \mathbb{R}^{d \times n}$ .

1149

1150 **Lemma H.4** ([Meenakshi and Rajian, 1999]). For any two positive semi-definitive matrices  $A, B \in$   
 1151  $\mathbb{R}^{d \times d}$ , we have

- 1152 •  $\text{tr}[AB] \geq 0$ .
- 1153 •  $AB \succeq 0$  if and only if  $A$  and  $B$  commute.