## TOWARDS BETTER UNDERSTANDING OF IN-CONTEXT LEARNING ABILITY FROM IN-CONTEXT UNCERTAINTY QUANTIFICATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Predicting simple function classes has been widely used as a testbed for developing theory and understanding of the trained Transformer's in-context learning (ICL) ability. In this paper, we revisit the training of Transformers on linear regression tasks, and different from the existing literature, we consider a bi-objective prediction task of predicting both the conditional expectation  $\mathbb{E}[Y|X]$  and the conditional variance Var(Y|X). This additional uncertainty quantification objective provides a handle to (i) better design out-of-distribution experiments to distinguish ICL from in-weight learning (IWL) and (ii) make a better separation between the algorithms with and without using the prior information of the training distribution. Theoretically, we show that the trained Transformer reaches near Bayes optimum, suggesting the usage of the information of the training distribution. Our method can be extended to other cases. Specifically, with the Transformer's context window S, we prove a new generalization bound of  $\mathcal{O}(\sqrt{\min\{S,T\}/(nT)})$  on n tasks with sequences of length T, providing sharper analysis compared to previous results of  $\mathcal{O}(\sqrt{1/n})$ . Empirically, we illustrate that while the trained Transformer behaves as the Bayes-optimal solution as a natural consequence of supervised training in distribution, it does not necessarily perform a Bayesian inference when facing task shifts, in contrast to the *equivalence* between these two proposed in many existing literature. We also demonstrate the trained Transformer's ICL ability over covariates shift and prompt-length shift and interpret them as a generalization over a meta distribution.

032 033 034

035

000

001

002

004

006

012 013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

#### 1 INTRODUCTION

A particularly remarkable characteristic of Large Language Models (LLMs) is their ability to perform in-context learning (ICL) (Brown et al., 2020). Once pre-trained on a vast corpus of data, LLMs can solve newly encountered tasks when provided with just a few training examples, without any updates to LLMs' parameters. ICL has significantly advanced the technique known as prompt engineering (Ekin, 2023), which has achieved widespread success in various aspects of daily life (Oppenlaender et al., 2023; Heston & Khun, 2023; Li et al., 2023a). Behind the empirical success of ICL, this method has captured the attention of the theoretical machine learning community, leading to considerable efforts into understanding ICL from different theoretical perspectives (Xie et al., 2021; Akyürek et al., 2022; Von Oswald et al., 2023; Zhang et al., 2023a).

This work aims to enhance the theoretical understanding of ICL by examining the Transformer's context window and showing its effects on the approximation-estimation tradeoff. Although we obtain the results for the case of uncertainty quantification where the model is asked to predict both the mean value and the uncertainty of its prediction, our analysis is applicable across various ICL tasks and provides sharper bounds compared to previous works. In addition to developing theories, we empirically demonstrate the effectiveness of Transformers to in-context predicting the mean and quantifying the variance of regression tasks. We design a series of out-of-distribution (OOD) experiments, which have generated significant interest within the community (Garg et al. (2022); Raventós et al. (2024); Singh et al. (2024)). These experiments provide insights in designing the pre-training process and understanding the ICL capabilities of transformers.

### Our contributions are as follows:

 $\begin{array}{rcl} & & - \mbox{ We study the problem of in-context uncertainty quantification for pre-trained Transformers which aims to predict both the mean and the variance of the conditional distribution of the target variable Y given the feature X. It is a task both of independent interest and providing insights into the in-context learning ability of the transformers. We derive the Bayes-optimal learner for the task and show that the transformer has the ability to fulfill this task of in-context uncertainty quantification. \\ \end{array}$ 

- We theoretically analyze the problem of in-context uncertainty quantification. We consider the case 061 when Transformers can only process the contexts within a context window capacity S and derive 062 a generalization bound of  $\mathcal{O}(\sqrt{\min\{S,T\}}/(nT))$  for pre-training over n tasks with sequences of 063 length T (Theorem 3.2). Our result can be easily extended to other cases under the assumption of 064 almost surely bounded and Lipschitz loss functions. As far as we know, our generalization bound is 065 the first of its kind and provides a tighter bound compared to the existing analyses (Li et al., 2023b; 066 Zhang et al., 2023b) when S < T. In particular, we use the context-window structure to establish 067 a Markov chain over the prompt sequence and construct an upper bound for its mixing time. We 068 also examine the extra approximation error term due to a finite context window S (Section B.2). 069 Combining those discussions together, we quantify the convergence of the trained Transformer's risk to the *Bayes-optimal* risk. Moreover, we note that all the theoretical results only show that the 071 trained Transformer achieves a near-optimal in-distribution risk compared to that of the Bayes-optimal 072 predictor. It is incorrect to draw (from the theory or the in-distribution numerical results) either of the conclusions that (i) the Transformer that achieves the near-optimal risk exhibits a similar structure as 073 the Bayes-optimal predictor by performing Bayesian inference (Zhang et al., 2023b; Panwar et al., 074 2023) or (ii) the Transformer performs as the Bayes-optimal predictor for out-of-distribution tasks. 075

076 - Numerically, we provide a comprehensive study of the in-context learning ability of the trained 077 Transformer (through the lens of uncertainty quantification) under three scenarios of distribution shifts: task shift (Section 4.1), covariates shift (Section 4.2), and prompt length shift (Section 4.3). 078 We find that transformers are capable of in-context learning of both mean and uncertainty predictions, 079 even under a moderate amount of task distribution shift, provided that the task diversity in the training data is relatively large. Additionally, we find that increasing the task diversity with a meta-learning 081 approach helps the transformer learn in-context robustly under covariates shift. Lastly, we observe 082 that removing positional encoding from the embedding vector massively helps the generalization 083 ability, enabling it to better learn tasks in-context with unseen prompt length. 084

085 We defer more discussions on the related literature to Section A.

086 087

880

094

#### 2 PROBLEM SETUP

Consider training a Transformer for some regression task  $f: \mathcal{X} \to \mathcal{Y}$  from a function class  $\mathcal{F}$ . The covariates  $x \in \mathcal{X} \subset \mathbb{R}^d$  are generated from a distribution  $\mathcal{P}_{\mathcal{X}}$ , and the output variable  $y = f(x) + \sigma \cdot \epsilon$ for some function  $f \in \mathcal{F}$ , noise level  $\sigma$ , and some random noise  $\epsilon$  with  $\mathbb{E}[\epsilon] = 0$  and  $\operatorname{Var}(\epsilon) = 1$ . The Transformer performs a sequential prediction task over the following sequence

$$x_1, y_1, \dots, x_T, y_T$$

where T is the total number of (in-context) samples. For a Transformer model with parameters  $\theta \in \Theta$ , we denote it as  $\text{TF}_{\theta}$ . At each time t = 1, ..., T, the model  $\text{TF}_{\theta}$  observes  $H_t \coloneqq (x_1, y_1, ..., x_{t-1}, y_{t-1}, x_t)$  (which is called *history* or *prompt*) and makes a bi-objective prediction of  $y_t$  to both predict the mean with  $\hat{y}_{\theta}(H_t)$  and quantify the uncertainty of the prediction with  $\hat{\sigma}_{\theta}(H_t)$ . With a slight abuse of notations, we denote the output of the model by  $\text{TF}_{\theta}(H_t) \coloneqq (\hat{y}_{\theta}(H_t), \hat{\sigma}_{\theta}(H_t))$ . The pre-training dataset consists of n sample sequences

$$\mathcal{D} \coloneqq \left\{ \left( x_1^{(i)}, y_1^{(i)}, x_2^{(i)}, y_2^{(i)}, \dots, x_T^{(i)}, y_T^{(i)} \right) \right\}_{i=1}^n$$

To generate each sample sequence in  $\mathcal{D}$ , a function  $f_i$  is sampled from a distribution  $\mathcal{P}_{\mathcal{F}}$  supported on  $\mathcal{F}$  and a noise level  $\sigma_i$  is sampled from a distribution  $\mathcal{P}_{\sigma}$  supported on  $[0, \bar{\sigma}] \subset \mathbb{R}$ . Then each  $x_t^{(i)}$ and  $y_t^{(i)}$  is generated pairwise by

107

where  $\epsilon_t^{(i)}$ 's are i.i.d. noise of mean zero and unit variance.

110 The Transformer is trained by minimizing the following empirical loss

$$\hat{\theta}^{\text{ERM}} \coloneqq \underset{\theta \in \Theta}{\operatorname{arg\,min}} \frac{1}{nT} \sum_{i=1}^{n} \sum_{t=1}^{T} \ell\left( \operatorname{TF}_{\theta}\left(H_{t}^{(i)}\right), y_{t}^{(i)} \right) \tag{1}$$

where  $H_t^{(i)} = (x_1^{(i)}, y_1^{(i)}, ..., x_{t-1}^{(i)}, y_{t-1}^{(i)}, x_t^{(i)})$  and  $l((\cdot, \cdot), \cdot) : (\mathbb{R} \times \mathbb{R}^+) \times \mathbb{R} \to \mathbb{R}$  denotes the loss function. We use  $x_t^{(i)}, y_t^{(i)}, H_t^{(i)}$  to denote the samples in the training dataset and  $x_t, y_t, H_t$  to denote an arbitrary feature, label, and history. Throughout this paper, we assume that each probability distribution is continuous and has a probability density function (p.d.f.), and we also assume the conditional distribution of  $y_t$  on observing  $H_t$  exists almost surely.

120 The loss function is accordingly defined by

$$\ell\left((\hat{y},\hat{\sigma}),y\right) := \log \hat{\sigma} + \frac{(y-\hat{y})^2}{2\hat{\sigma}^2}$$

**Definition 2.1** (Bayes-optimal predictor). The Bayes-optimal predictor under the distributions  $\mathcal{P}_{\mathcal{F}}$ ,  $\mathcal{P}_{\chi}$ ,  $\mathcal{P}_{\sigma}$  and  $\mathcal{P}_{\epsilon}$  is defined by

$$(y_t^*(\cdot), \sigma_t^*(\cdot)) \coloneqq \underset{(y(\cdot), \sigma(\cdot)) \in \mathcal{G}_t \times \mathcal{G}_t}{\operatorname{arg\,min}} \mathbb{E}\left[\ell\Big(\big(y(H_t), \sigma(H_t)\big), y_t\Big)\right]$$
(2)

where  $\mathcal{G}_t$  is the class of all measurable functions of  $H_t \in \mathcal{H}_t$ . The expectation is taken with respect to the following dynamics:  $x_t \sim \mathcal{P}_{\mathcal{X}}$ ,  $\epsilon_t \sim \mathcal{P}_{\epsilon}$ ,  $f \sim \mathcal{P}_{\mathcal{F}}$ ,  $\sigma \sim \mathcal{P}_{\sigma}$ ,  $y_t = f(x_t) + \sigma \cdot \epsilon_t$  and  $H_t = (x_1, y_1, \dots, x_t)$ .

The loss on the right-hand-side of equation 2 is the expectation of the empirical loss equation 1. With a rich enough function class and an infinite amount of training samples, the trained Transformer TF<sub> $\hat{\theta}_{ERM}$ </sub> converges to  $(y_t^*, \sigma_t^*)$  as will be shown in Theorem 3.2.

### 136 2.1 MOTIVATION FOR THE UNCERTAINTY QUANTIFICATION OBJECTIVE

We begin with a semi-formal definition of in-context learning and in-weight learning that are considered in this paper. *In-context learning* (ICL) refers to the ability of the model to accomplish some task regardless of the prior knowledge received during training. On the contrary, *In-weight learning* (IWL) represents the knowledge acquired during the training phase (that is stored in the model's weights).

However, it is not easy to separate the ICL ability completely from the IWL ability due to two major obstacles: the first is that whatever ability the model obtains from the training, the training distribution inevitably incorporates the prior knowledge into the model, which makes it difficult to test the pure ICL ability of a model. The second obstacle is that although the common practice to examine the ICL ability is to twist the distribution at the test phase, the criterion of setting the test distribution is still vague. If the performance of the model is bad for a twisted test distribution, is it due to the failure to attain an ICL ability or because the inputs are out-of-domain during the test phase?

148 We first present an abstract training/test framework to attack the second obstacle. We define the set 149 of test distributions capable of testing the ICL ability yet not intriguing the out-of-domain issues 150 by setting a meta-distribution and constraining any input sequences that are "common" in the test 151 distribution must also be "common" in the training distribution. More rigorously, at the training 152 phase, a meta distribution  $\Xi_{\text{train}}$  is chosen. For each sequence, we first draw a pattern  $\mathcal{P}$  (e.g. the weight vector and the noise level) from the meta-distribution  $\mathcal{P} \sim \Xi_{\text{train}}$ . Then a sequence S is drawn 153 from this pattern such that  $S \sim \mathcal{P}$ . At the test phase, another meta distribution  $\Xi_{\text{test}}$  is chosen. After 154 observing S, the model tries to identify the pattern behind the sequence by making a prediction 155  $\mathcal{P}(S)$  (e.g. trying to identify the noise level or the weight vector). The model feeds on sequences 156 S's as inputs while the performance is evaluated by how close the predicted pattern  $\hat{\mathcal{P}}$  is to the true 157 pattern  $\mathcal{P}$ . Therefore, we can solve the out-of-domain issue at the *sequence* level while keeping the 158 distribution twisted at the *pattern* level. Denote the marginal distribution of the sequence S under the 159 meta-distribution  $\Xi$  as 160

161

111 112 113

121 122 123

124

125 126 127

$$\mathbb{P}(S \in \mathcal{S}|\Xi) = \int \mathbb{P}(S \in \mathcal{S}|S \sim \mathcal{P}) d\mathbb{P}(\mathcal{P}|\mathcal{P} \sim \Xi).$$

Then to avoid the out-of-domain input problems, one just needs to impose a constraint that any sequence that "commonly appears" at the case  $\Xi = \Xi_{\text{test}}$  also "commonly appears" at the training phase  $\Xi = \Xi_{\text{train}}$ :

190 191 192

193

199

200 201

202

 $\exists C > 0, \quad \text{s.t.} \ \frac{\mathbb{P}(S \in \mathcal{S} | \Xi_{\text{train}})}{\mathbb{P}(S \in \mathcal{S} | \Xi_{\text{test}})} \ge C, \quad \forall \mathcal{S} \subset \text{supp}(\mathbb{P}(\cdot | \Xi_{\text{test}})).$ (3)

168 Denote the set of all those  $\Xi_{\text{test}}$ 's as  $\mathcal{T}(\Xi_{\text{train}})$ .

We provide an example here to help understand the concept. For any *d*-dimensional linear regression task sequence  $(X_1, Y_1, \ldots, X_t, Y_t)$  (say t < d), one cannot tell whether it comes from a noiseless pattern (by solving the linear equations) or from a noisy pattern. One can only guess which pattern is more likely based on their knowledge.

This constraint also explains why doing ICL v.s. IWL detection under a single-target regression task 174 (i.e. predicting the mean only) is hard. On the one hand, current observations (Garg et al., 2022) 175 show that under certain distribution-shift circumstances, the model performance is poor. We can 176 conclude those failed cases by violating the constraint equation 3. On the other hand, the typical 177 way of twisting the test distribution is to associate the pattern with the inputs and flip the connection 178 during the test phase (Wei et al., 2023; Singh et al., 2024) to distinguish in-context learning from 179 in-weight learning. In the mean-prediction-only case, by changing the weight vector's association 180 with the feature X, the domain of Y will also be entirely changed (e.g. flipped or multiplied). But if 181 we are considering the bi-objective task, by only changing the noise level, one can barely tell whether 182 it is a perturbation in the weight vector or a change in the noise level.

With the constraint equation 3 on the test distributions, we now discuss the way to deal with the first obstacle that the training will inevitably incorporate prior knowledge. Intuitively, when the prior knowledge is helpful (e.g.  $\Xi_{train} = \Xi_{test}$ ), IWL and ICL cooperate to achieve an "ICL + IWL" result; when the prior knowledge is a fault, IWL can hurt the test performance, making the final performance behave like "ICL - IWL". Then we can implicitly measure the influence caused by the IWL of the model (denoted by  $\mathcal{M}$ ) by defining

$$\mathrm{IWL}(\mathcal{M}) \coloneqq \frac{1}{2} \left( \sup_{\Xi_{\mathrm{test}} \in \mathcal{T}(\Xi_{\mathrm{train}})} \mathbb{E} \left[ \ell(\hat{\mathcal{P}}_{\mathcal{M}}(S), \mathcal{P}) \right] - \inf_{\Xi_{\mathrm{test}} \in \mathcal{T}(\Xi_{\mathrm{train}})} \mathbb{E} \left[ \ell(\hat{\mathcal{P}}_{\mathcal{M}}(S), \mathcal{P}) \right] \right)$$

We can use the average performance between the "best" and the "worst" case as a proxy of the pure ICL performance:

$$\mathrm{ICL}(\mathcal{M}) \coloneqq \frac{1}{2} \left( \sup_{\Xi_{\mathrm{test}} \in \mathcal{T}(\Xi_{\mathrm{train}})} \mathbb{E} \left[ \ell(\hat{\mathcal{P}}_{\mathcal{M}}(S), \mathcal{P}) \right] + \inf_{\Xi_{\mathrm{test}} \in \mathcal{T}(\Xi_{\mathrm{train}})} \mathbb{E} \left[ \ell(\hat{\mathcal{P}}_{\mathcal{M}}(S), \mathcal{P}) \right] \right).$$

Therefore, we can approximate the ICL and the IWL abilities by approximating the best and the worst model performances for all the distributions under the constraint equation 3.

#### 3 IN-CONTEXT LEARNING WHEN IN-DISTRIBUTION

In this section, we focus on the in-distribution property of the trained Transformer. We provide a finite-sample analysis of how trained Transformers reach near Bayes-optimum. While our analysis is made on the case of uncertainty quantification, it can be easily adapted to other loss functions such as mean squared error. To proceed, we first provide the exact form of the Bayes-optimal predictor defined in equation 2 for the mean and uncertainty prediction.

Proposition 3.1 (Bayes-optimal predictor for mean and uncertainty prediction). The Bayes-optimal predictor of the step-wise population risk defined in equation 2 is given by

211 
$$y_t^*(H_t) = \mathbb{E}[y_t|H_t], \ \sigma_t^{*2}(H_t) = \mathbb{E}[(y_t - y_t^*(H_t))^2|H_t] = \mathbb{E}[(f(x_t) - y_t^*(H_t))^2|H_t] + \mathbb{E}[\sigma^2|H_t].$$
212

The optimal mean predictor shares the same form as the Bayes-optimal predictor for a single-objective mean prediction task. The additional uncertainty prediction task does not change the nature of the mean prediction part. The two terms in the optimal uncertainty predictor can be interpreted as follows. The first term is epistemic uncertainty, which indicates the uncertainty (of identifying the f that



Figure 1: Transformer behaves close to the Bayes-optimal predictor for in-distribution tasks. Details of the 231 distributions in data generation are given in Section G.1. The numbers 4096 and 65536 refer to the number of 232 tasks (configurations of  $(w_i, \sigma_i)$ ) used in the training, which is formally defined in Section G.2. The Bayes-233 optimal predictor is stated in Proposition 3.1 and calculated analytically in Section G.3. For the left panel, the 234 y-axis gives the mean squared error in predicting  $y_t$ . For the right panel, the y-axis gives the average of the 235 predicted uncertainty over all the test samples (average of  $\hat{\sigma}(H_t)$  or  $\sigma^*(H_t)$  on test samples). In particular, we note that ridge regression and linear regression (ordinary least squares) do not naturally produce a measurement 236 of uncertainty, so we use the sum of residuals on the in-context samples as their estimates of uncertainty. More 237 visualizations are deferred to Section C.1. 238

239

243

249

250

251

253

254

255

256

257

258

263

240 governs the history  $H_t$ ) due to lack of information. The term decreases as the samples accumulate, 241 i.e., as the number of in-context samples t increases. The second term is aleatoric uncertainty also 242 known as intrinsic uncertainty.

Recall that the empirical risk estimator is defined by equation 1. Now we define the population risk as

$$R(\mathrm{TF}_{\theta}) \coloneqq \frac{1}{T} \mathbb{E}_{H_t} \left[ \sum_{t=1}^T \ell \big( \mathrm{TF}_{\theta}(H_t), y_t \big) \right],$$

where  $H_t$  is another sampled sequence that is independently and identically distributed as  $H_t^{(i)}$ 's in the training data. We denote the population risk minimizer as  $\theta^*$ :

$$\theta^* \in \operatorname*{arg\,min}_{\theta \in \Theta} R(\mathrm{TF}_{\theta}). \tag{4}$$

Now we present our main theoretical result.

**Theorem 3.2.** Let  $\hat{\theta}^{ERM}$  denote the ERM estimator as defined in equation 1 over the function class of the L-layer, M-heads Transformer models. Suppose that at each time t, the Transformer has a context window of making predictions based on  $x_t$  and previous S pairs of  $(x_s, y_s)$  for  $s = \max\{1, t - S\}, \ldots, t - 1$ . Then under some boundedness assumptions of the Transformer's parameters (Assumption B.5 and B.6), we have with probability at least  $1 - \delta$ ,

$$R(\mathrm{TF}_{\hat{\theta}^{\mathrm{ERM}}}) - R(\mathrm{TF}_{\theta^*}) \leq \tilde{\mathcal{O}}\left(\sqrt{\min\{S,T\}/(nT)}\right).$$

where  $\mathcal{O}$  omits poly-logarithmic terms that depend on  $n, T, 1/\delta$  and boundedness parameters.

**Proof sketch.** First, we prove that (a slightly redefined version of) the truncated history forms up a Markov chain conditioned on observing the full hidden information  $f^{(i)}$  and  $\sigma^{(i)}$ , and upper bound the mixing time by min{S, T} to enable the concentration arguments. Second, we prove that the loss function is almost surely bounded (Lemma E.3) in preparation for McDiarmid-type concentration inequalities (Lemma F.2, (Paulin, 2015)). Third, we show that the loss is almost surely Lipschitz to control the difference between loss functions with respect to the change of the parameter (Lemma E.7). Fourth, we prove that there exist two distributions  $\rho_{\hat{\theta}\text{ERM}}$  and  $\pi$  over parameter space  $\Theta$ , satisfying a number of properties as constructed in Lemma E.11. Lastly, we use standard PAC-Bayes arguments over  $\rho_{\hat{\rho}_{ERM}}$  and  $\pi$  and conclude the proof. The detailed proofs are deferred to Section D.2.

272 **Comparison with previous results.** There are also other theoretical results that characterize the 273 outcomes of the (pre-)training on Transformer models (Zhang et al., 2023a; Wu et al., 2023; Xie 274 et al., 2021; Li et al., 2023b; Bai et al., 2024; Zhang et al., 2023b; Lin et al., 2023a). Our analysis 275 differs from theirs in terms of both the conclusion and the techniques. One stream of results examines 276 the property of the gradient flow (or gradient descent) over the loss function for linear regression 277 problems. The exact quantification of the gradient flow entails a simplification of the Transformer's 278 architecture to the case of a single-layer attention mechanism under linear activation or even simpler 279 settings (Zhang et al., 2023a; Wu et al., 2023). While their analyses provide insights into the learning 280 dynamics of Transformer models, the learning of the single-layer attention Transformer can be very different from multiple-layer Transformers (Olsson et al., 2022; Reddy, 2023). Another major line 281 of research uses statistical learning arguments (Xie et al., 2021; Li et al., 2023b; Bai et al., 2024; 282 Zhang et al., 2023b; Lin et al., 2023a) such as algorithm stability, chaining, or PAC-Bayes arguments. 283 Bai et al. (2024) focus on making predictions after observing a fixed length of variables under the 284 i.i.d. setting (which is more aligned with the standard supervised learning setting), which differs 285 from the more practical setting of making predictions at every position as in Theorem 3.2. Xie et al. 286 (2021) prove the convergence between the Bayesian inference and the true underlying distribution 287 rather than the trained model and the Bayesian inference. Lin et al. (2023a) consider a sequential 288 decision-making problem and use covering arguments to derive generalization bounds, while their 289 analysis does not adopt the concentration arguments inside each sequence, resulting in an  $\tilde{O}(\sqrt{1/n})$ 290 upper bound for the average regret. The most related works to ours are Li et al. (2023b); Zhang et al. 291 (2023b). The major difference is that they all consider the only case of S > T. Li et al. (2023b) use 292 the algorithm stability arguments to give a generalization bound over |R - r| of order  $O(\sqrt{1/(nT)})$ . 293 They prove the loss difference caused by perturbing one input pair over a history of length t is 294 controlled by  $\mathcal{O}(1/t)$ . Averaging those differences leads to a  $\mathcal{O}(\log(T)/T) = \mathcal{O}(1/T)$  inside each sequence (see their equation (15) in their Appendix C), which appears in the Azuma-Hoeffding 295 argument to prove that the loss per sequence is  $ilde{\mathcal{O}}(T^{-1/2})$ -sub-Gaussian. However, in the case of 296  $S \ll T$ , the algorithm stability term is of  $\mathcal{O}(1/S)$ . Averaging these terms inside each sequence 297 leads to a difference of order  $\mathcal{O}(1/S)$ . If we stick to the original Azuma-Hoeffding arguments, 298 the sum of squares of these terms is of  $\mathcal{O}(T/S^2)$ , leading to a far worse sub-Gaussian norm of 299  $\mathcal{O}(T^{1/2}S^{-1})$ , resulting in a final generalization bound of order  $\tilde{O}(T^{1/2}S^{-1}n^{-1/2})$  that is clearly 300 suboptimal compared to our  $\tilde{O}(\sqrt{S/(nT)})$ . Besides, such a bound also grows with T, which is 301 undesirable. Similar to ours, Zhang et al. (2023b) also use a concentration argument for Markov 302 chains. However, their Theorem 5.3 has two limitations: The first is that their result is of the order 303  $\mathcal{O}(\sqrt{\tau_{\min}}/(nT))$  but they do not specify  $\tau_{\min}$ . Since they do not consider the truncated history 304 but the full history, the Markov chain (which is not verified by them) will never mix inside each 305 task sequence (see our discussions in Section D.2). Thus, the term  $\tau_{\min}$  in their result is actually 306 T, leading to an order of  $\tilde{\mathcal{O}}(\sqrt{1/n})$ , which is suboptimal compared to our  $\tilde{\mathcal{O}}(\sqrt{S/(nT)})$  when the 307 context window  $S \ll T$ . The second limitation is that their error decomposition is not tight: their 308 excessive risk bound (measured by the total variation distance between the distribution induced by  $\theta$ 309 and that by  $\theta^*$ ) has a term  $D_{kl}(\mathcal{P}_{true}, \mathcal{P}_{\theta^*}) - TV(\mathcal{P}_{true}, \mathcal{P}_{\theta^*})$ , which means their result has an extra 310 term of the approximation error since the Kullback-Leibler divergence is stronger than the total 311 variation distance (Polyanskiy & Wu, 2024). Our work is the first theoretical analysis showing the 312 effects of the context window S on the performance of the Transformer up to our knowledge. The 313 construction of the truncated history serves two-fold: not only does the truncation fit the practical 314 model of finite context window but it also gives an upper bound on the mixing time. Concentration 315 inside each sequence makes it possible to analyze the training dynamics broader than fixed-length sequences and prove the convergence to near Bayes-optimum. The context window S also captures a 316 novel dimension of the approximation-estimation tradeoff in the Transformer model. 317

**Extension of Theorem 3.2 to other problems.** We remark that the result and its derivation do not pertain to the uncertainty quantification setting, but hold for more general loss functions and are of independent interests. In particular, our analysis still holds as long as the loss function is almost surely bounded and Lipschitz with respect to the change of parameter  $\theta$ , as we can see from the proof sketch. We note here that to enable the Markov chain's concentration arguments, the almost surely **223**  bounded loss requirement cannot be relaxed to other tail properties such as sub-Gaussian (see the counter example in Theorem 4 of Fan et al. (2021)). 

We defer discussions on the approximation error to Section B.2.

#### **IN-CONTEXT LEARNING UNDER DISTRIBUTION SHIFTS**

In Section 2, we describe in-context learning ability as algorithm-like that predicts based on the learning from in-context samples, and such an ability should be generalizable to an out-of-distribution (OOD) environment. In this section, we differentiate the OOD scenarios into task shift, covariate shift and length shift, and examine the Transformer's in-context learning ability in each scenario. As far as we know, we provide the first comprehensive group of numerical experiments (for the linear regression task) that demonstrates the Transformer's ability to handle these three types of distribution shifts. We provide preliminary theoretical discussions for such abilities and hope this points directions for future theoretical research. 

4.1 TASK SHIFT

When the trained Transformer performs well on the OOD data, it means that the Transformer gains an algorithmic ability that learns to make predictions based on the in-context samples, because such an ability is not restricted to the distribution of the inputs. Comparatively, the mere observation that the Transformer works well on the in-distribution data does not demonstrate its in-context learning ability as a traditional supervised learning model also has such ability and generalization performance over in-distribution data. 

In the previous section, when we show the in-distribution performance of the Transformer, the variance parameter  $\sigma^2$  is generated by the prior of the inverse-Gamma distribution  $\sigma^2 \sim \text{Inv-Gamma}(\tau, \bar{\tau})$ with parameters  $\tau$  and  $\bar{\tau}$ . The details of the other generation distributions are deferred to Section G.1. For the in-distribution setting, we set  $\tau = \overline{\tau} = 20$  which leads to a prior mean around 1. Now we consider three out-of-distribution (OOD) settings for the

- S-OOD (small OOD):  $\tau = 80$ ,  $\overline{\tau} = 20$ . The prior mean of  $\sigma$  is around 0.5.
- M-OOD (medium OOD):  $\tau = 100, \bar{\tau} = 400$ . The prior mean of  $\sigma$  is around 2.
- L-OOD (large OOD):  $\underline{\tau} = 100, \overline{\tau} = 1600$ . The prior mean of  $\sigma$  is around 4.





378 We make following observations based on Figure 2: First, the Bayes-optimal predictor predicts well. 379 We note that the Bayes-optimal is computed based on the in-distribution prior distribution (with 380 respect to  $\sigma^2$ ). Thus when the Bayes-optimal predictor is tested under the OOD environment as in 381 Figure 2, the prior used by the Bayes-optimal predictor is wrong. But we note from Figure 2 that the 382 Bayes-optimal predictor has the OOD ability to correct the prior as the in-context samples accumulate (noting that the three Bayes-optimal curves converging to the correct mean of 0.5, 2, and 4). This is also known as the washing out of priors in Bayesian statistics. Second, Transformers deviate from the 384 Bayes-optimal on these OOD tasks. For both plots in Figure 2, we note that the predicted values from 385 the Transformers deviate from those of the Bayes-optimal predictor when the OOD intensity is large. 386 This tells that the trained Transformer does not conduct Bayesian inference under task shift. In other 387 words, it is incorrect to conclude that the trained Transformer behaves as the Bayes-optimal predictor 388 just from the matching in-distribution loss (as Figure 1). Moreover, the Transformer achieves a 389 near-optimal loss for in-distribution tasks (as Figure 1) but it does so via a different avenue than the 390 Bayes-optimal predictor (as Figure 2). This is in contrast with the findings/claims in the previous 391 papers (Zhang et al., 2023b; Panwar et al., 2023). Third, the deviation of the trained Transformer 392 from the Bayes-optimal is smaller when the task diversity is large or the OOD intensity is small. This 393 is aligned with the findings in (Raventós et al., 2024) for in-distribution performance, while the OOD setting is not studied therein. 394

The theoretical evidence only states that the trained Transformer has a near-optimal in-distribution 396 loss as the Bayes-optimal predictor. But it does not give any evidence that these two have a structural 397 similarity that persists for OOD tasks. In particular, we note that the trained Transformer may take 398 statistical shortcuts: When evaluated under in-distribution tasks or some simple task shifts (e.g. 399 scaling the weights vectors or changing the signal-noise ratio), Zhang et al. (2023a); Wu et al. (2023) show that Transformer will construct shortcuts using the statistical property of the training distribution. 400 More specifically, Transformers (can, and will) encode the information of the covariance matrix into 401 their model parameters to reach near-optimal in-distribution performance. Such statistical shortcuts 402 are beneficial to the in-distribution performance but can hurt its OOD ability. Increasing the training 403 task diversity, such as a larger training pool size, may remove some of these statistical shortcuts to 404 obtain near-optimal empirical loss, and thus better enable its in-context learning ability. 405

406 We defer more discussions and visualizations on this OOD experiment to Section C.2.

408 4.2 COVARIATES SHIFT

For all the numerical experiments so far, the covariates are generated from  $\mathcal{N}(0, I_d)$ . This follows the standard setup of the existing literature (Akyürek et al., 2022; Von Oswald et al., 2023; Li et al., 2023b; Raventós et al., 2024). It is also noted from the literature (Garg et al., 2022; Zhang et al., 2023a) that the trained Transformer in this way lacks in-context learning ability under covariates shift. In this subsection, we propose a meta-training procedure that effectively improves the trained Transformer's ability to handle covariates shifts. Specifically, we consider generating the covariates in the training data as follows:

For each training sequence (say, the *i*-th), we first sample a vector (λ<sub>1</sub>,...,λ<sub>d</sub>) where each λ<sub>j</sub> is i.i.d. Uniform[0,2]. Then all the X<sub>t</sub><sup>(i)</sup>'s for t = 1,...,T are sampled from N(0, diag((λ<sub>1</sub>,...,λ<sub>d</sub>))). In this sense, the covariance matrix of X<sub>t</sub><sup>(i)</sup>'s is also a random variable, and the X<sub>t</sub><sup>(i)</sup>'s can be viewed as being sampled in a hierarchical manner from a meta-distribution.

We examine the performance of such a training procedure under four OOD test settings. In other words, the  $X_t^{(i)}$ 's in the test data is generated from the following four distributions where d = 8.

425 426 427

407

409

417

418 419

420

421

422 423

424

428 429

- Large covariance (L-cov):  $X_t^{(i)}$ 's are sampled from  $\mathcal{N}(0, 4I_d)$ .
- Decreasing diagonal (Dec.):  $X_t^{(i)}$ 's are sampled from  $\mathcal{N}(0, \operatorname{diag}([d/i]_{i=1}^d))$ .
- Shrinking diagonal (Shr.):  $X_t^{(i)}$ 's are sampled from  $\mathcal{N}(0, \operatorname{diag}([d/i^2]_{i=1}^d))$ .
- Rotation (Rot.):  $X_t^{(i)}$ 's are sampled from  $\mathcal{N}(0, U_i \operatorname{diag}([d/i]_{i=1}^d) U_i^{\top})$  where  $U_i$  is an orthogonal matrix independently generated for each sequence.

Figure 9 gives the evaluation result under the 4 OOD settings. We note that the meta-distribution used is still significantly different from the four OOD test environments. Thus the results show the effectiveness of the meta-training approach.

435 436

437

### 4.3 LENGTH SHIFT AND POSITIONAL EMBEDDING

438 Existing work (Dai et al., 2019; Anil et al., 2022; Zhang et al., 2023a) have pointed out the failure 439 of Transformers to generalize to longer contexts than the ones they have seen during training. It is worth mentioning that the code implementations of some previous works (Zhang et al., 2023a; 440 Garg et al., 2022) are based on the "transformers" package of Hugging Face. Although these works 441 have not included positional embedding explicitly, the GPT2 module imported from this package 442 adds a built-in positional encoding implicitly. We suspect that some unexpected behaviors (like the 443 "unexpected spikes of prediction error" mentioned in Zhang et al. (2023a)) are due to that the built-in 444 positional encoding is not disabled. In this subsection, we investigate the length generalization ability 445 of the trained Transformer on the uncertainty quantification task. Specifically, we control the prompt 446 lengths that the model is trained on. Previous experiments train the model on prompts with lengths 447 (number of in-context samples) ranging from 1 to 100. In this experiment, we control the training 448 prompts such that the lengths are either shorter than 44 or longer than 45 (the choice of 45 as the 449 cutoff point is not essential). We specify these two configurations below.

450 451

452 453

454

- Trained on  $\leq 44$ : the model is trained on prompts with length ranging from 1 to 44, and is evaluated with prompt length from 1 to 100
- Trained on  $\geq 45$ : the model is trained on prompts with length ranging from 45 to 100, and is evaluated with prompt length from 1 to 100

455 We regard this difference in prompt length between training and testing as length shift. We evaluate 456 the effect of removing positional encoding under this prompt length generalization task. If positional 457 encoding is added to the embedding, samples at unseen positions will be associated with an unseen 458 positional encoding vector in the embedding space. This requires the model to handle not just 459 an unseen number of in-context samples, but also a possibly unseen embedding distribution, and 460 generalization ability will likely deteriorate. As mentioned previously, the built-in positional encoding 461 of GPT2 model use a positional encoding which is set to be  $(t, 0, \dots, 0)^{\top}$  for the t-th token, and the 462 encoding will then be concatenated to the embedding vector. We validate the above intuitions with 463 the following 4 training configurations.

- 464
- 465 466

467

468

469

470

471

472

473

474

- No positional encoding (w/o Pos.): the model is trained without positional encoding.
- Add positional encoding (w/ Pos.): the model is trained with GPT2's built-in positional encodings.
  - Add segment encoding (w/ S-Pos.): the positional encoding is added with a random amount offset. For the *i*-th training sequence, a random offset t<sub>i</sub> is first uniformly sampled from {0, 1, ..., 22}. Next, for each token in this prompt at position t, the positional encoding is set to (t + t<sub>i</sub>, 0, ..., 0)<sup>T</sup>.
  - Add full range encoding (w/ F-Pos.): similar to the S-Pos. configuration, the positional encoding is added with a random amount offset. But here the offset is uniformly sampled from {0, 1, ..., 100}.
- For the model trained with the "w/o Pos." configuration, it is also tested without positional encodings. For the models trained with the rest configurations, they are all tested with the "w/ Pos." way of encoding.
- The results are shown in Figure 3. The models in the left figure are trained on prompts shorter than
  44, and the models in the right figure are trained on prompts longer than 45. We make the following
  observations. The pre-trained transformer in general can generalize to prompts with unseen length,
  under the condition of using/removing the positional embedding properly. The "w/o Pos." curve in
  the left figure shows that even at positional encoding hurts the generalization ability. From the "w/
  Pos." curve in the left figure, we find that the model's performance drops significantly at positions
  larger than 44. The main cause of the failure of length generalization is due to the distribution shift



Figure 3: The effect of removing positional encoding on prompt length generalization. The *y*-axis records the average error of uncertainty prediction, which is the difference between the uncertainty predicted by the transformer and the Bayes-optimal estimator. (a) For models trained with prompt lengths  $\leq$  44, the figure on the left shows that positional encoding has the worst generalization capacity with a larger length, and removing positional encoding could effectively enhance the length generalization power. (b) For models trained with prompt lengths  $\geq$  45, removing positional encoding can help generalize to smaller lengths, although the generalization ability for smaller lengths is generally weaker compared to that for larger lengths.

507 508

in the positional embedding space. As given in the "w/ S-Pos." and "w/ F-Pos." curves in the left
figure, if the model has seen the *positional encodings* for a certain position during training, then its
performance at this position is significantly improved, even if the *corresponding prompt length* is
never seen. The length generalization ability is not unrestrictively strong, and such generalization
ability for smaller lengths is generally weaker compared to that for larger lengths. The right figure
shows that even for the "w/o Pos." configuration, its performance still degrades when the prompt
length is shorter than 20.

Theoretically, Wu et al. (2023)'s Theorem 5.3 points out that under the case of the single-layer linear-attention-only Transformer model on a linear regression task with Gaussian priors, if we train the model to only predict one single label after observing T context exemplars, the optimally trained model under  $T = T_1$  also performs well at the case  $T = T_2$  (compared to the Bayes-optimal predictor for  $T = T_2$ ) if  $T_1$  and  $T_2$  are close. This result implies the possibility of context length generalization by a simplified Transformer model due to shared structures in the attention matrices.

522 523

#### 5 CONCLUSION

524 525 526

In this paper, we study the in-context learning ability of the trained Transformer through the lens of 527 uncertainty quantification. In particular, we train the Transformer for a bi-objective task of mean 528 prediction and uncertainty prediction. We develop new results both theoretically and numerically. 529 The takeaway messages are: First, the Transformer can perform in-context uncertainty quantification. 530 Second, the trained Transformer is only guaranteed to achieve a near-optimal in-distribution risk 531 against the Bayes-optimal predictor. This does not imply that the Transformer behaves as the 532 Bayes-optimal predictor either in-distribution or out-of-distribution. Third, the Transformer has the 533 in-context ability for out-of-distribution tasks, but this in-context ability is contingent on a proper 534 training method such as sufficient task diversity, meta-training for covariates shift, and effective removal of the positional encoding. Two important future directions are as follows. First, we believe 536 our method for deriving the generalization bound has implications for a scope much larger than 537 uncertainty quantification and can be used to improve the existing bounds for various tasks using Transformers. Second, all the numerical experiments in the paper are conducted for the linear 538 functions  $f_i$ 's. We believe the same results still hold for nonlinear functions as well; and such results can further consolidate the in-context ability for uncertainty quantification of the Transformer.

### 540 REFERENCES

- Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information fusion*, 76:243–297, 2021.
- Gustaf Ahdritz, Tian Qin, Nikhil Vyas, Boaz Barak, and Benjamin L Edelman. Distinguishing the knowable from the unknowable with language models. *arXiv preprint arXiv:2402.03563*, 2024.
- 549 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement
   550 preconditioned gradient descent for in-context learning. Advances in Neural Information Processing
   551 Systems, 36, 2024.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556, 2022.
- Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:
   Provable in-context learning with in-context algorithm selection. *Advances in neural information processing systems*, 36, 2024.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
  Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
  few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 567 Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. Inside:
  568 Llms' internal states retain the power of hallucination detection. *arXiv preprint arXiv:2402.03744*,
  569 2024.
- Herman Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, pp. 493–507, 1952.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdi nov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Monroe D Donsker and SR Srinivasa Varadhan. Asymptotic evaluation of certain markov process expectations for large time. iv. *Communications on pure and applied mathematics*, 36(2):183–212, 1983.
- Sabit Ekin. Prompt engineering for chatgpt: a quick guide to techniques, tips, and best practices.
   *Authorea Preprints*, 2023.
- Fabian Falck, Ziyu Wang, and Christopher C Holmes. Are large language models bayesian? a martingale perspective on in-context learning. In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*, 2024.
- Jianqing Fan, Bai Jiang, and Qiang Sun. Hoeffding's inequality for general markov chains and its applications to statistical learning. *Journal of Machine Learning Research*, 22(139):1–35, 2021.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Jakob Gawlikowski, Cedrique Rovile Njieutcheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt,
   Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, et al. A survey of
   uncertainty in deep neural networks. *Artificial Intelligence Review*, 56(Suppl 1):1513–1589, 2023.

| 594<br>595<br>596        | Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai. How do transformers learn in-context beyond simple functions? a case study on learning with representations. <i>arXiv preprint arXiv:2310.10616</i> , 2023.  |
|--------------------------|---|
| 597<br>598<br>599        | Thomas F Heston and Charya Khun. Prompt engineering in medical education. <i>International Medical Education</i> , 2(3):198–205, 2023.  |
| 600<br>601               | Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. <i>arXiv preprint arXiv:2302.09664</i> , 2023.  |
| 602<br>603<br>604        | Beibin Li, Konstantina Mellou, Bo Zhang, Jeevan Pathuri, and Ishai Menache. Large language models for supply chain optimization. <i>arXiv preprint arXiv:2307.03875</i> , 2023a.  |
| 605<br>606<br>607        | Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In <i>International Conference on Machine Learning</i> , pp. 19565–19594. PMLR, 2023b.   |
| 608<br>609<br>610        | Licong Lin, Yu Bai, and Song Mei. Transformers as decision makers: Provable in-context reinforce-<br>ment learning via supervised pretraining. <i>arXiv preprint arXiv:2310.08566</i> , 2023a.  |
| 611<br>612               | Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantifica-<br>tion for black-box large language models. <i>arXiv preprint arXiv:2305.19187</i> , 2023b.   |
| 614<br>615<br>616        | Potsawee Manakul, Adian Liusie, and Mark JF Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. <i>arXiv preprint arXiv:2303.08896</i> , 2023.   |
| 617<br>618               | Colin McDiarmid et al. On the method of bounded differences. <i>Surveys in combinatorics</i> , 141(1): 148–188, 1989.   |
| 619<br>620<br>621<br>622 | Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. <i>arXiv preprint arXiv:2209.11895</i> , 2022.   |
| 623<br>624               | Jonas Oppenlaender, Rhema Linder, and Johanna Silvennoinen. Prompting ai art: An investigation into the creative skill of prompt engineering. <i>arXiv preprint arXiv:2303.13534</i> , 2023.  |
| 625<br>626<br>627        | Madhur Panwar, Kabir Ahuja, and Navin Goyal. In-context learning through the bayesian prism. In <i>The Twelfth International Conference on Learning Representations</i> , 2023.   |
| 628<br>629<br>630        | Daniel Paulin. Concentration inequalities for Markov chains by Marton couplings and spectral methods. <i>Electronic Journal of Probability</i> , 20(none):1 – 32, 2015. doi: 10.1214/EJP.v20-4039. URL https://doi.org/10.1214/EJP.v20-4039.  |
| 631<br>632               | Yury Polyanskiy and Yihong Wu. Information theory: From coding to learning. 2024.   |
| 633<br>634               | Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.   |
| 635<br>636<br>637<br>638 | Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the emergence of non-bayesian in-context learning for regression. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.  |
| 639<br>640               | Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context classification task. <i>arXiv preprint arXiv:2312.03002</i> , 2023.   |
| 641<br>642<br>643<br>644 | Aaditya Singh, Stephanie Chan, Ted Moskovitz, Erin Grant, Andrew Saxe, and Felix Hill. The transient nature of emergent in-context learning in transformers. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.   |
| 645<br>646<br>647        | Aviv Slobodkin, Omer Goldman, Avi Caciularu, Ido Dagan, and Shauli Ravfogel. The curious case of hallucinatory (un) answerability: Finding truths in the hidden states of over-confident large language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 3607–3625, 2023. |

648

| 648<br>649 | Ralph C Smith. Uncertainty quantification: theory, implementation, and applications. SIAM, 2013.   |
|------------|--|
| 650        | Timothy John Sullivan. Introduction to uncertainty quantification, volume 63. Springer, 2015.  |
| 651        | Alexander Tsigler and Peter L Bartlett. Benign overfitting in ridge regression. Journal of Machine   |
| 652<br>653 | Learning Research, 24(123):1–76, 2023.   |
| 654        | Johannes Von Oswald, Evvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,  |
| 655        | Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In  |
| 656        | International Conference on Machine Learning, pp. 35151–35174. PMLR, 2023.   |
| 657        | Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu  |
| 658        | Da Huang, Denny Zhou, et al. Larger language models do in-context learning differently. <i>arXiv</i>   |
| 659        | preprint arXiv:2303.03846, 2023.   |
| 661        | Jingteng Wu Difan Zou, Ziviang Chen, Vladimir Brayerman, Auanguan Gu, and Peter I. Bartlett  |
| 662        | How many pretraining tasks are needed for in-context learning of linear regression? <i>arXiv preprint</i>  |
| 663        | arXiv:2310.08391, 2023.  |
| 664        | Sana Mishael Vie Aditi Dashurathan Dana Linna and Tanawa Ma An anglanatian af in contact   |
| 665        | learning as implicit havesian inference arXiv preprint arXiv:2111.02080, 2021  |
| 666        | ioanning as implicit obyesian inference. <i>arxiv preprint arxiv.2111.</i> 02000, 2021.  |
| 667        | Fengzhuo Zhang, Boyi Liu, Kaixin Wang, Vincent Tan, Zhuoran Yang, and Zhaoran Wang. Relational   |
| 668        | reasoning via set transformers: Provable efficiency and applications to marl. Advances in Neural   |
| 670        | Information 1 rocessing Systems, 55.55825–55858, 2022.   |
| 671        | Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context.  |
| 672        | <i>arXiv preprint arXiv:2306.09927</i> , 2023a.  |
| 673        | Yufeng Zhang, Fengzhuo Zhang, Zhuoran Yang, and Zhaoran Wang. What and how does in-context   |
| 674        | learning learn? bayesian model averaging, parameterization, and generalization. arXiv preprint   |
| 675        | <i>arXiv:2305.19420</i> , 2023b.   |
| 676        |  |
| 677        | A RELATED WORKS  |
| 678        |  |
| 680        | Theoretical Understanding of In-Context Learning. There are two streams of research in the   |
| 681        | theoretical understanding of ICL: the first tries to give sharp approximation error bounds on different  |
| 682        | tasks, while the second focuses on how the trained Transformer approaches the potential optimum.   |
| 683        | For the approximation error, following the pioneering empirical investigations on simple function classes (Garg et al. 2022). Von Oswald et al. (2023): A kyiirek et al. (2022) conjecture that the        |
| 684        | Transformer is doing ICL via gradient descent, and verify it both empirically and theoretically.   |
| 685        | Based on the mechanism of layer-wise gradient descent construction, Bai et al. (2024) show that  |
| 686        | Transformers are able to behave (approximately) as well as some well-known algorithms on some  |
| 687        | statistical problems. Some following works generalize the layer-wise gradient descent construction to  |
| 000<br>080 | outer settings such as decision-making (Lin et al., $2025a$ ) and integression under representations (Guo et al. 2023). Apart from the layer-wise gradient descent, some other works consider the ope-step |
| 600        | gradient descent reached by a single-layer linear-activated Transformer (Zhang et al., 2023a; Wu   |
| 691        | et al., 2023) and curve the excessive population risk of the optimal model compared to oracle or the   |
| 692        | Bayes-optimal predictor. Ahn et al. (2024) give a set of global optima for some specific one-layer or  |
| 693        | two-layer attention-only models with linear or ReLU activation. Aside from characterizing where the  |
| 694        | Iransformer <i>can</i> reach, another group of works is making efforts towards understanding where the   |
| 695        | et al. (2023a) start the analysis of the training dynamics of the gradient flow over the population  |
| 696        | risk on the linear regression task and show that a single-layer linear-attention-only model converges  |
| 697        | to some specific sets with suitable initialization. Wu et al. (2023) keep the same spirit and give a   |
| 698        | sample complexity bound based on a certain gradient descent scheme. For general Transformer  |
| 699        | models, technical tools from the statistical learning theory are applied. As for the task of predicting  |
| 700        | the next token in natural language tasks, Ale et al. (2021) provide a viewpoint from the Hidden  |
| 701        | Markov Model (HMM) and prove the asymptotic consistency under the regularity condition. Rejet al   |

(2024); Lin et al. (2023a) use chaining arguments with covering numbers for generalization, where

702 Bai et al. (2024) consider the training under fixed length and Lin et al. (2023a) consider the problem 703 of sequential decision-making. Li et al. (2023b) adopt algorithm stability arguments obtaining a 704 bound of  $O(1/\sqrt{nT})$ . As is discussed in the main text (see discussions after Theorem 3.2), their 705 analysis will result in a suboptimal  $O(S/\sqrt{nT})$  for the case S < T. Zhang et al. (2023b) adopt a 706 similar concentration inequality for Markov chains to get a bound of  $\mathcal{O}(\sqrt{\tau_{\text{mix}}/(nT)})$ . Since they 707 do not consider the limit of context window S, their derivation ends up with  $\tau_{mix} \ge T$ , which is 708 suboptimal compared to our case. In short, our paper is the first theoretical analysis on the limit of 709 context window S and gets a tighter generalization bound than previous works on the generalization 710 bound when S < T.

711

712 Bayesian Behavior of In-context Learning. Due to the complex structure of transformers, showing 713 the theoretical properties of ICL without proper assumptions is challenging. There has been growing interest in developing experiments to test various properties of ICL, leading to new observations and 714 insights. Some of the earliest works that show transformers behave like Bayesian estimator can be 715 found in Akyürek et al. (2022); Garg et al. (2022), and this argument is supported in follow-up works 716 including Li et al. (2023b); Wu et al. (2023); Bai et al. (2024). However, there is also increasing 717 empirical evidence demonstrating transformers' non-Bayesian behavior. Singh et al. (2024) design 718 flipped experiment and show transformers' Bayesian behavior could be transient. Raventós et al. 719 (2024); Panwar et al. (2023) demonstrate that the Bayesian behavior of transformers is dependent 720 on the task diversity in the pre-training dataset, and transformers could deviate from the Bayesian 721 predictor if number of different training tasks is large. Falck et al. (2024) design experiments based 722 on the martingale property, a necessary condition of Bayesian behavior, and provide evidence that 723 transformers exhibit non-Bayesian behavior from a statistical perspective.

724

735

736

725 Transformers for Uncertainty Quantification. Uncertainty quantification has seen significant development within the general machine learning and deep learning domains (Abdar et al. (2021); Gaw-726 likowski et al. (2023)), generating considerable interests within communities working on transformer-727 based large language models (LLMs). See Kuhn et al. (2023); Manakul et al. (2023); Lin et al. 728 (2023b) for uncertanty quantification using black-box LLMs, and Slobodkin et al. (2023); Chen et al. 729 (2024); Ahdritz et al. (2024) for that of white-box LLMs. Most of these works focus on natural 730 language processing tasks which have less statistical properties. Indeed, uncertainty quantification 731 has traditionally been developed from a more statistical and probabilistic perspective (Smith (2013); 732 Sullivan (2015)). By adopting transformer models to study more statistics-related problems, our work 733 aims to bridge and contribute to both fields. 734

### **B** TRANSFORMER MODEL

737 Following Radford et al. (2019), we consider a decoder-only L-layer Transformer model that pro-738 cesses the input sequence  $H_t$  by applying multi-head attention (MHA) of M heads and multi-layer 739 perceptron (MLP) layer-wise. Without loss of generality, we assume  $x_t \in \mathbb{R}^d$  for some  $d \geq 2$ . 740 We concatenate each  $y_t$  with d-1 zeros so that it matches the format of each  $x_t$ , while we still 741 denote the concatenated vector by  $y_t$  with a slight abuse of notations. We denote  $H_t$  by a matrix 742 in  $\mathbb{R}^{d \times (2t-1)}$  for  $t = 1, \ldots, T$ , where  $H_t = [x_1, y_1, \cdots, x_t]$ . We may also refer to  $x_t$  by  $h_{2t-1}$ 743 and  $y_t$  by  $h_{2t}$ . In practice, the attention mechanism has a maximum dependence length, and there-744 fore the Transformer model can only produce an output based on the most recent tokens up to a 745 context window size S. Hence we assume that at each time step t, the Transformer model has a 746 maximum capacity of making predictions based on  $x_t$  and previous S pairs of  $(x_s, y_s)$  observations for  $s = t - S, \dots, t - 1$ . In other words, the Transformer has a maximum capacity of processing 747 2S + 1 tokens, and it is making predictions  $TF_{\theta}(H_t) = TF_{\theta}(H_t^S)$  based on the truncated history 748  $H_t^S$ , where  $H_t^S := (x_{\max\{1,t-S\}}, y_{\max\{1,t-S\}}, \dots, x_t)$ . In the following, we formally describe the 749 architecture of the Transformer used in this paper. 750

**Definition B.1** (Multi-Head Attention). A multi-head attention layer with M heads and activation function  $act(\cdot)$  can be defined as a function  $MHA_W(\cdot)$  for any sequence  $Z_t \in \mathbb{R}^{d \times (2t-1)}$  and  $t = 1, \dots, S+1$ ,

$$\mathsf{MHA}_W(Z_t) = Z_t + \sum_{m=1}^M (W_V^m Z_t) \mathsf{act}\left((W_K^m Z_t)^\top (W_Q^m Z_t)\right),$$

756 where  $W = \{(W_Q^m, W_K^m, W_V^m)\}_{m=1}^M$  denotes all the parameters,  $W_Q^m, W_K^m \in \mathbb{R}^{d_m \times d}, W_V^m \in \mathbb{R}^{d \times d}$  for each  $m = 1, \ldots, M$ , and  $\operatorname{act} : \mathbb{R}^{(2t-1) \times (2t-1)} \to \mathbb{R}^{(2t-1) \times (2t-1)}$  is the activation function. 757 758

759 Here we merge the residual connection into the multi-head layer and skip the layer normalization to 760 ease the notations and simplify the analysis. The activation function is usually set to be columns-wise softmax in practice: for each vector  $z \in \mathbb{R}^{2t-1}$ , 761

$$\operatorname{softmax}(z) \coloneqq \left(\frac{\exp(z_1)}{\sum_{i=1}^{2t-1}\exp(z_i)}, \dots, \frac{\exp(z_{2t-1})}{\sum_{i=1}^{2t-1}\exp(z_i)}\right)^\top.$$

765 Some theoretical results also consider alternative choices for act. For example, Akyürek et al. 766 (2022); Ahn et al. (2024); Zhang et al. (2023a) consider the linear activation (that is, to entry-wise 767 divide by the sequence length 2t - 1). Bai et al. (2024); Guo et al. (2023) also examine the ReLU 768 activation (that is, to entry-wise apply a ReLU function ReLU(z) = max{0, z} and later divide by 769 the sequence length 2t - 1).

770 Definition B.2 (Multi-Layer Perceptron). A multi-layer perceptron layer with hidden dimension 771  $d_h$  can be defined as a (token-wise) function MLP<sub>A</sub>(·) for any sequence  $Z_t \in \mathbb{R}^{d \times (2t-1)}$  and 772  $t = 1, \ldots, S + 1,$ 773

$$\mathsf{MLP}_A(Z_t) = Z_t + A_2 \mathsf{ReLU}(A_1 Z_t),$$

where  $A = (A_1, A_2)$  denotes all the parameters,  $A_1 \in \mathbb{R}^{d_h \times d}$ ,  $A_2 \in \mathbb{R}^{d \times d_h}$ , and ReLU is the 774 775 entry-wise ReLU function. 776

We merge the residual connection into the multi-layer perceptron layer and omit the layer normaliza-777 tion to simplify the theoretical development. 778

Definition B.3 (Transformer). A Transformer model with L layers can be defined as a function 779  $TF_{\theta}(\cdot)$  for any sequence  $Z_t \in \mathbb{R}^{d \times (2t-1)}$  and  $t = 1, \ldots, S+1$ . For the *l*-th layer, the model receives 780  $Z_t^{(l-1)}$  as the input and processes it by an MHA block and an MLP block, such that 781

$$Z_t^{(l)} = \mathrm{MLP}_{A^{(l)}}(\mathrm{MHA}_{W^{(l)}}(Z_t^{l-1})), \quad \forall l = 1, \dots, L,$$

where  $Z_t^{(0)} = Z_t$ . After the *L*-th layer, the model linearly maps the  $Z_t^{(L)} \in \mathbb{R}^{d \times (2t-1)}$  onto  $\mathbb{R}^{2 \times (2t-1)}$  via a matrix  $P \in \mathbb{R}^{2 \times d}$ , and we process the second dimension by a softplus function to get the 784 785 final prediction as 786

 $\hat{y}_{\theta}(Z_t) = (PZ_t^{(L)})_{1,2t-1},$ 

and

782 783

787 788

789

798

762 763 764

$$\hat{\sigma}_{\theta}(Z_t) = \text{softplus}\left((PZ_t^{(L)})_{2,2t-1}\right)$$

790 Here  $\theta = (\{(W^{(l)}, A^{(l)})\}_{l=1}^{L}, P)$  encapsulates all the parameters and the function softplus(z) =791  $\log(1 + \exp(z))$  is introduced to avoid negative output. The output is summarized as  $\text{TF}_{\theta}(Z_t) \coloneqq$ 792  $(\hat{y}_{\theta}(Z_t), \hat{\sigma}_{\theta}(Z_t)).$ 793

Remark B.4. To enable parallel training, the decoder-only Transformer receives a full sequence in 794 the training phase. The model has a masking component that prevents the model from seeing into the 795 "future". However, such masking is unnecessary in our setting as the Transformer model receives 796 exactly what it should "see" at each time t, and the full dynamics are identical to those in the masked 797 setting.

**Miscellaneous notations.** Denote the set  $\{1, \ldots, K\}$  by [K]. Denote the consecutive sequence 799  $\{i, i+1, \ldots, j\}$  by i : j. Denote the matrix A's entry at the *i*-th row and the *j*-th column by 800  $A_{i,j}$ . Denote the vector x's *i*-th element by  $(x)_i$ . Define the d-dimensional vector x's p-norm as 801  $(\sum_{i=1}^{d} (x)_{i}^{p})^{1/p}$  for  $p \in [1, \infty]$ , where  $||x||_{\infty} = \max_{1 \le i \le d} (x)_{i}$ . Define the  $m \times n$ -sized matrix A's (p, q)-norm as  $(\sum_{j=1}^{n} ||A_{:,j}||_{p}^{q})^{1/q}$ . Denote the d-dimensional diagonal matrix by  $\operatorname{diag}\{\lambda_{1}, \ldots, \lambda_{d}\}$ . 802 803 804 Denote the *d*-dimensional identity matrix by  $I_d$ . Denote the total variation distance between two 805 probability distributions P and Q by TV(P, Q). Denote the Kullback-Leibler divergence between 806 two probability distributions such that  $P \ll Q$  by  $D_{kl}(P||Q)$ . Denote the product measure of P and Q by  $P \times Q$  or  $P \otimes Q$ . Denote the Cartesian product of two spaces  $\mathcal{X}$  and  $\mathcal{Y}$  by  $\mathcal{X} \times \mathcal{Y}$ . Denote the 807 tensor-product  $\sigma$ -algebra of two  $\sigma$ -algebras  $\Sigma_1$  and  $\Sigma_2$  by  $\Sigma_1 \otimes \Sigma_2$ . Denote the limiting behavior of 808 being upper (lower, both upper and lower, respectively) bounded by up to some constant(s) by  $\mathcal{O}(\Omega)$ , 809  $\Theta$ , respectively). Denote  $\hat{O}$  to be the O but omitting some poly-logarithmic terms.

# 810 B.1 ASSUMPTIONS

Based on this setup of the Transformer model, we introduce the following bounded assumptions used
for Theorem 3.2. Such assumptions are common in the analyses of the Transformer model (Bai et al.,
2024; Zhang et al., 2023b) by either assuming an extra clipping operator or explicit upper bounds.

Assumption B.5. Assume  $\Theta = \mathcal{B}(0, B_{\text{TF}})$ , where the norm is defined as

816 817

819 820 821

822

823

824

825 826

827 828

829

830

835

836

837 838

839 840

841 842

843 844

845 846

$$\|\theta\| \coloneqq \max\{\|W^{(l)}\|, \|A^{(l)}\|, \|P\|: l = 1, \dots, L\}.$$

<sup>818</sup> The corresponding norms are defined as

$$||W|| \coloneqq \max\{||W_V^m||_{2,2}, ||W_K^m||_{2,2}, ||W_Q^m||_{2,2} : m = 1, \dots, M\},\$$

 $||A|| \coloneqq \max\{||A_1||_{2,2}, ||A_2||_{2,2}\}, \quad ||P|| \coloneqq ||P^+||_{2,\infty},$ 

where we omit some superscripts/subscripts of the layer number (l) for simplicity.

**Assumption B.6.** Assume  $||H_t||_{2,\infty}$  is bounded by  $B_H$  almost surely. Such a regularization is equivalent to assuming  $||x_t||_2 \leq B_H$  and  $|y_t| \leq B_H$  almost surely.

#### B.2 APPROXIMATION ERROR

In Section 3, we provide the generation bound in Theorem 3.2. Now we give an analysis for the approximation error. We define the Bayes-optimal risk obtained by the Bayes-optimal predictor in Proposition 3.1: for each t = 1, ..., T,

$$R_t^* := \mathbb{E}\bigg[\ell\Big(\big(y_t^*(H_t), \sigma_t^*(H_t)\big), y_t\Big)\bigg].$$
(5)

However, Transformers only have access to the truncated history  $H_t^S$ , which prevents them from reaching  $R_t^*$ . By using Proposition 3.1 for the  $H_t^S$ , we denote the *truncated* Bayes optimum for each t:

$$y_t^{S*}(H_t^S) \coloneqq \mathbb{E}[y_t | H_t^S],$$

and

$$\left(\sigma_t^{S*}(H_t^S)\right)^2 \coloneqq \mathbb{E}[(f(x_t) - y_t^{S*}(H_t^S))^2 | H_t^S] + \mathbb{E}[\sigma^2 | H_t^S].$$

We denote the truncated Bayes-optimal risk as

$$R_t^{S*} \coloneqq \mathbb{E}\bigg[\ell\Big(\big(y_t^{S*}(H_t^S), \sigma_t^{S*}(H_t^S)\big), y_t\Big)\bigg].$$
(6)

It is straightforward to check that

$$R_t^{S*} = R_t^*, \quad \text{for any } t \le S. \tag{7}$$

However, the equality is generally not true for t > S. We give an example to illustrate the gap.

**Example B.7.** Consider the case where one has oracle access to the noise level  $\sigma$ . Note that the oracle knowledge only reduces the risk  $R_t^{S*}$ , since we use information that is not a measurable function of  $H_t^S$ . The problem is reduced to a regression problem.

Suppose the function f is linear and its weight vector has a prior distribution of  $\mathcal{N}(0, \sigma^2 I_d)$ , and the noise  $\epsilon \sim \mathcal{N}(0, 1)$ . Suppose  $x_t \sim \mathcal{N}(0, I_d)$ . Then the optimal estimator is an optimally tuned Ridge regression.

**Trigler & Bartlett** (2023) show that with high probability, the optimal Ridge regression estimator has an average risk of  $\frac{1}{2} + \Theta(1/S)$ , where the term  $\frac{1}{2}$  is due to  $\mathbb{E}[(y - f(x))^2]/(2\sigma^2)$ . But as the length t approaches infinity, the average risk of the optimal Ridge regression over the full sequence  $H_t$  will converge to  $\frac{1}{2}$  with high probability, meaning that the estimated  $\hat{f}$  will converge to true f for every sequence. Hence one can always construct an uncertainty estimation by averaging all the residuals, and such an estimation  $\hat{\sigma}$  will converge to the true  $\sigma$ . Thus, we have  $R_t^*$  approaching  $\frac{1}{2}$  as t grows to infinity, leading to the conclusion that

$$R_t^{S*} - R_t^* \ge \Omega(1/S),$$

for sufficiently large t.

Example B.7, together with Theorem 3.2, shows the approximation-estimation tradeoff in selecting the context window S of Transformer models. Previous works (Wu et al., 2023; Bai et al., 2024; Zhang et al., 2023b; Guo et al., 2023) consider the case where  $t \leq S$ , and establish the upper bounds for the approximation error. In other words, these existing results are all made with respect to the gap between  $R_t^{S*}$  and  $R(TF_{\theta^*})$ . To our knowledge, we are the first work to point out the extra term of approximation error due to truncation.

#### C MORE NUMERICAL RESULTS AND DISCUSSIONS

#### C.1 IN-DISTRIBUTION PERFORMANCE

In Figure 1, we provide a comparison of the in-distribution performance of the trained Transformer v.s. the Bayes-optimal predictor. A subtle point is that for the uncertainty prediction, we only plot the average predicted uncertainty, this does not fully imply that the Transformer gives a similar prediction as the Bayes-optimal predictor. To this end, Figure 4 plots the difference between each of the models and the Bayes-optimal predictor in terms of uncertainty estimation.



Figure 4: In-distribution performance of the uncertainty prediction against the Bayes-optimal predictor. The y-axis gives an estimate of  $\mathbb{E}\left[-\log |\hat{\sigma}(H_t) - \sigma^*(H_t)|\right]$  where the expectation is taken with respect to  $H_t$ . Here  $\hat{\sigma}(H_t)$  is the uncertainty estimate produced by an algorithm (ridge regression, linear regression, or transformer), and  $\sigma^*(H_t)$  is the Bayes-optimal predictor given in Proposition 3.1 and calculated by Section G.3. The figure shows that the Transformer and the Bayes-optimal predictor produce similar uncertainty predictions. In addition, the Transformer trained on a larger pool of tasks (larger N) produces a better approximation of the Bayes-optimal predictor.

901 902

903

904

870 871

872 873

874 875

876

877

878

879

880

882

883

885

886

888

889

890 891

892

893

894

### C.2 OUT-OF-DISTRIBUTION PERFOMANCE

In Figure 2, we plot the Bayes-optimal predictor under three OOD settings, and we note that though 905 the Bayes-optimal predictor uses a wrong prior, it has the ability to work as an algorithm to correct 906 the prediction with the in-context samples. Now in Figure 5 (a), we compare the Bayes-optimal 907 predictor that uses the wrong prior with the Bayes-optimal predictor that uses the correct prior (which 908 replaces the in-distribution prior with the correct OOD prior of  $\sigma^2$ ). The figure is based on the large 909 OOD setting. We observe that the Bayes-optimal predictors with the ID prior or the OOD prior both 910 converge to the true uncertainty level. For Figure 5 (b), we plot the performances under the same 911 large OOD setting. As a reference line, we copy and paste the Bayes-optimal predictor's curve in 912 Figure 1 (b) here. We note this reference line is computed based on the in-distribution (ID) data and is 913 not comparable at all to the predicted uncertainty level on the OOD data. Yet, we note that when the 914 number of tasks is small when training the Transformer (say N = 4096), it tends to make predictions 915 on the OOD data by treating the OOD data just as ID data, and this means the trained Transformer is doing in-weight learning and has no in-context learning ability. As the number of tasks increases, the 916 Transformer gradually gains the in-context ability and moves towards the Bayes-optimal predictor on 917 the OOD data.



Figure 5: Performance under L-OOD setting. For both (a) and (b), the y-axis gives the average of the predicted uncertainty over all the test samples (average of  $\hat{\sigma}(H_t)$  or  $\sigma^*(H_t)$  on test samples), and ideally the curves should converge to the true uncertainty level of 4 as the number of in-context samples increases. In (a), we compare the Bayes-optimal predictor that uses the wrong prior with the Bayes-optimal predictor that uses the correct prior (which replaces the in-distribution prior with the correct OOD prior of  $\sigma^2$ ). Both work well in that the curves converge to the true mean uncertainty level of around 4. The Transformers deviate from both Bayes-optimal predictors due to the large OOD intensity. In (b), we observe that as the training task diversity increases. The transformer gradually moves from the ID reference line to the Bayes-optimal predictor.

939 940 941

942

963 964

933

934

935

936

937

938

#### C.3 TRAINING DYNAMICS AND TASK SHIFT OOD PERFORMANCE

Now we zoom into the training dynamics to further investigate the OOD performance under task shift. In the following example, we derive a theoretical result based on Theorem 4.1 in Zhang et al. (2023a). Specifically, R and R' (following the notations therein) denote the in-distribution and out-of-distribution expected risk. The result says that while the in-distribution risk continues decreasing over time, the out-of-distribution risk may keep increasing or may first decrease and then increase. Importantly, the out-of-distribution risk may depend on the initial point of the training procedure.

Reaching the Bayes optimum requires prior knowledge of the underlying distribution. We provide
a simple example of where the Transformer stores its prior knowledge and how it hurts the OOD
performance even under a mild distribution shift.

953 **Example C.1** (A corollary that can be derived based on Theorem 4.1 in Zhang et al. (2023a)). 954 Consider a one-layer attention-only Transformer model with linear activation and one attention 955 head on the linear regression task. We now concatenate each of the inputs to be  $[x_t^{\top}, y_t]^{\top} \in \mathbb{R}^{d+1}$ . Suppose we focus on the linear regression task on the (T + 1)-th sample after observing T context 956 exemplars, where each  $w^{(i)} \sim \mathcal{N}(0, I_d)$ , each  $x_t^{(i)} \sim \mathcal{N}(0, I_d)$ , each  $\epsilon_t^{(i)} \sim \mathcal{N}(0, 1)$ , and  $y_t^{(i)} = 0$ 957 958  $w^{(i)\top}y_t^{(i)} + \sigma_0 \cdot \epsilon_t^{(i)}$ . If we adopt the same training setup as Zhang et al. (2023a) (with details referred 959 to therein), then for any  $|\sigma'_0 - \sigma_0| \ge \Delta$  for some  $\Delta > 0$ , if we train on the distribution w.r.t.  $\sigma_0$  but 960 test on the distribution w.r.t.  $\sigma'_0$  (denoted by R'), then for  $C = d/(16(2 + \sigma_0))$  and any sequence 961  $0 < \delta_1 < \delta_2 < \cdots < C\Delta$ , there exists a non-decreasing sequence  $0 \le \tau(\delta_1) \le \tau(\delta_2) \le \ldots$ , such 962 that

$$R'(\tau(\delta_i)) - R'_{\theta^{*'}} \ge \delta_i, \quad \text{for each } i = 1, 2, \dots,$$

while the parameter  $(W_K^{\top}W_Q)_{1:d,1:d}(W_V)_{d+1,d+1}$  converges to  $1/(1+(2+\sigma_0)/T) \cdot I_d$  (which is the corresponding part of some  $\theta^*$ ). Here  $\theta^*$  and  $\theta^{*'}$  minimize the population risk R and R', accordingly.

We design an experiment to show that as training proceeds, the model's OOD performance is improved abruptly in the starting phase, but then degrades steadily after too many steps of training. We introduce the experiment settings below. A visualization of the setup is given in Figure 6.

Each linear task in our uncertainty quantification setting is characterized by parameters  $(w, \sigma)$ . We define two regions for w, denoted by  $W_1$  and  $W_2$ . And two regions for  $\sigma$ , denoted by  $\mathcal{G}_1$  and  $\mathcal{G}_2$ .



Figure 6: The settings of the OOD experiment. (Left) The in-distribution (ID) tasks are sampled from regions denoted by the green blocks, and the OOD tasks are sampled from the red blocks. (Right) In the starting phase, training improves both the ID and OOD performance. But if training for too many steps, the ID performance is only marginally improved, while the OOD performance steadily degrades.

992 When w is sampled from  $W_1$ , w follows the following distribution

 $w = |\beta|, \quad \beta \sim \mathcal{N}(0, I_8).$ 

996 When w is sampled from  $W_2$ , w follows the following distribution

$$w = -|\beta|, \quad \beta \sim \mathcal{N}(0, I_8).$$

1000 For  $\sigma$ , define  $\mathcal{G}_1 = [0.1, 0.3] \cup [0.5, 0.7]$  and  $\mathcal{G}_2 = [0.3, 0.5] \cup [0.7, 0.9]$ . Define  $\mathcal{G}_1$  and  $\mathcal{G}_2$  to be the "complementary" group of each other. We sample  $\sigma$  independently from w, and we always 1001 sample  $\sigma$  uniformly from either group  $\mathcal{G}_1$  or  $\mathcal{G}_2$ . As marked in Figure 6, the "ID" tasks sample its 1002 parameters from  $(w, \sigma) \in W_1 \otimes \mathcal{G}_1 \bigcup W_2 \otimes \mathcal{G}_2$  and the "OOD" tasks sample its parameters from 1003  $(w,\sigma) \in \mathcal{W}_1 \otimes \mathcal{G}_2 \bigcup \mathcal{W}_2 \otimes \mathcal{G}_1$ . The training is on "ID" tasks, and the trained model is tested on 1004 both "ID" tasks and "OOD" tasks. The metric we evaluate in this experiment is the "prediction 1005 accuracy" of uncertainty. The accuracy denotes the probability that the model predicts the  $\sigma$  into its "right" group. (for a prompt generated from  $(w, \sigma)$  with  $\sigma \in \mathcal{G}$ , we say that the model makes a "right" prediction if the predicted  $\sigma$  falls into  $\mathcal{G}$ ). 1008

Figure 7 presents the experiment result. The prediction accuracy on the ID dataset peaks after 20k
steps of training. At the same time, the prediction accuracy on the OOD dataset also increases to 80%. After that, the ID performance remains unchanged, but the OOD accuracy keeps dropping.

In order to verify that the degradation of OOD performance is due to the increasing confidence in 1012 the prior information of the training data, we check for the OOD distribution whether the model 1013 has predicted  $\sigma$  into the complementary group. The result is presented in Figure 8, which verifies 1014 that after training too many steps, the model tends to predict  $\sigma$  following the training prior. A more 1015 concrete way to explain it: consider an OOD sampled task  $(w, \sigma)$  where w > 0. According to the 1016 sampling rule of OOD tasks, it must have  $\sigma \in \mathcal{G}_2$ . If the model has the OOD ability, it should 1017 predict  $\sigma \in \mathcal{G}_2$ . But if it has too much confidence in its training prior, it will predict  $\sigma$  into  $\mathcal{G}_1$ , the 1018 complementary group of  $\mathcal{G}_2$ . Figure 8 shows that for the misclassified OOD tasks, the model has 1019 predicted them into complementary groups. 1020

1021

985

986

987

988 989

990 991

993

994 995

997 998

999

#### **1022** C.4 COVARIATE SHIFT EXPERIMENT

1023

The experiment result of Section 4.2 is given in Figure 9. We evaluate the prediction error of models ordinarily trained, and models trained by the meta-training process. For both mean and uncertainty, the models trained by the meta-training procedure have a smaller prediction error.



Figure 7: The accuracy denotes the probability that the model predicts the  $\sigma$  into the "right" group. For example, if the sampled tasks take  $\sigma$  from group  $\mathcal{G}_1$ , then accuracy denotes the probability that the model predicts  $\sigma$  into  $\mathcal{G}_1$ . The data is collected for the 100-th token in order to eliminate the epistemic uncertainty due to insufficient in-context samples. The x-axis denotes the training steps. This figure shows that when training too many steps (> 40k in this case), the generalization ability of the model steadily declines.



Figure 8: This figure validates that the decline of OOD ability is due to increasing confidence in the training prior. The blue bars correspond to the OOD accuracy, and the red bars give the probability that the model predicts uncertainty  $\sigma$  into the complementary group (i.e. the training distribution of  $\sigma$ ). As the training proceeds, most of the misclassified  $\sigma$  are predicted following the training prior.



Figure 9: The errors of the mean and uncertainty prediction where the error is measured by the absolute difference against the Bayes-optimal predictor. The *static\_x\_model* corresponds to models trained with the standard way in generating  $X_t$ 's, while the *meta\_x\_model* corresponds to the new approach of drawing  $X_t$ 's from the meta-training procedure. In all 4 OOD settings, meta-trained models have better performance.

1130

1131

1132

### D PROOFS OF THE RESULTS IN THE MAIN PAPER

# 1136 D.1 PROOF OF PROPOSITION 3.1

1138 *Proof.* Recall that the population risk is

1156

1159 1160

 $L(\hat{y}, \hat{\sigma}) \coloneqq \mathbb{E}_{f, x_{[t]}, \epsilon_{[t]}, \sigma} \left[ \log \hat{\sigma}(H_t) + \frac{(y_t - \hat{y}(H_t))^2}{2\hat{\sigma}^2(H_t)} \right].$ 

1142 We first prove that for any  $\hat{\sigma}(H_t)$ , the choice of  $\hat{y}_t = y^* = \mathbb{E}[y_t|H_t]$  minimizes the population risk. 1143 With any fixed  $\sigma_0 > 0$ , when  $\hat{\sigma}_t = \sigma_0$ , then minimizing the population risk reduces to minimizing 1144  $\mathbb{E}[(y_t - \hat{y}(H_t))^2]$ . Using Fubini's Theorem and the fact that the conditional distribution exists, we 1145 have

$$\mathbb{E}_{f,x_{[t]},\epsilon_{[t]},\sigma}\left[(y_t - \hat{y}(H_t))^2\right] = \mathbb{E}_{H_t}\left[\mathbb{E}\left[(y_t - \hat{y}(H_t))^2 \big| H_t\right]\right]$$
$$= \mathbb{E}_{H_t}\left[\mathbb{E}\left[(f(x_t) - \hat{y}(H_t))^2 \big| H_t\right]\right] + \mathbb{E}_{H_t}\left[\mathbb{E}[\sigma^2 | H_t]\right], \qquad (8)$$

where the last equality follows from the fact that  $\epsilon_t$  is independent of  $H_t$  and  $\sigma$  and is of zero mean and unit variance. Since the second term on the right-hand-side of equation 8 does not depend on  $\hat{y}$ , we only need to focus on the first term. For each realization of  $H_t$ , the prediction  $\hat{y}(H_t)$  is a single point; combining it with the fact that the squared loss is minimized with respect to one single point prediction if and only if that point is the expectation (in this case, the conditional expectation  $\mathbb{E}[f(x_t)|H_t]$ ), we prove that for any  $\sigma_0$ , the population risk's minimizer

$$y_t^*(\sigma_0) = \mathbb{E}[f(x_t)|H_t] = \mathbb{E}[f(x_t) + \sigma \cdot \epsilon_t|H_t] = \mathbb{E}[y_t|H_t]$$

where the second equality follows again from the fact that  $\epsilon_t$  is independent of  $H_t$  and  $\sigma$ , and  $\epsilon_t$  is of zero mean. Since this equality holds for an arbitrary  $\sigma_0$ , we can conclude that

$$y_t^* = \mathbb{E}[y_t|H_t].$$

Now we have confirmed the optimal choice of  $y_t^*$  regardless of whatever  $\hat{\sigma}$  is. We can thus find the optimal choice of  $\hat{\sigma}$  by fixing  $\hat{y} = y_t^*$  and minimizing the population risk. Similarly, we can change the integration order so that we only need to minimize  $\mathbb{E}[\log \hat{\sigma}(H_t) + \frac{(y_t - \hat{y}(H_t))^2}{2\hat{\sigma}^2(H_t)}|H_t]$  for any realization of  $H_t$ . Calculations show that

$$\begin{aligned} \frac{\partial \mathbb{E} \left[ \log \hat{\sigma}(H_t) + \frac{(y_t - \hat{y}(H_t))^2}{2\hat{\sigma}^2(H_t)} \mid |H_t \right]}{\partial \hat{\sigma}(H_t)} &= \frac{\partial \left( \log \hat{\sigma}(H_t) + \frac{\mathbb{E}[(y_t - \hat{y}(H_t))^2 \mid H_t]}{2\hat{\sigma}^2(H_t)} \right)}{\partial \hat{\sigma}(H_t)} \\ &= \frac{\hat{\sigma}^2(H_t) - \mathbb{E}[(y_t - \hat{y}(H_t))^2 \mid H_t]}{\hat{\sigma}^3(H_t)}, \end{aligned}$$

where the first equality follows from the fact that on observing  $H_t$ ,  $\hat{\sigma}(H_t)$  is a fixed value, and the second equality from the calculus. Thus, the risk is minimized if and only if  $\hat{\sigma}(H_t) = \mathbb{E}[(y_t - \hat{y}(H_t))^2 | H_t]$ . Substituting  $\hat{y}$  for  $y_t^*$ , we have

1178

1180

1183 1184

1187

$$\sigma_t^{*2}(H_t) = \mathbb{E}[(y_t - y_t^*(H_t))^2 | H_t] = \mathbb{E}[(f(x_t) - y_t^*(H_t))^2 | H_t] + \mathbb{E}[\sigma^2 | H_t],$$

where the last equality follows again from the fact that  $\epsilon_t$  is independent of  $H_t$  and is of zero mean and unit variance.

#### 1179 D.2 PROOF OF THEOREM 3.2

1181 Proof.  $\hat{\theta}^{\text{ERM}}$  Before we start the detailed proof, we define another flattened sequence  $(\tilde{x}_k, \tilde{y}_k)$  for  $k = 1, \dots, nT$ , where for k = iT + t we have

$$\left(\tilde{x}_{iT+t}, \tilde{y}_{iT+t}\right) \coloneqq \left(x_t^{(i)}, y_t^{(i)}\right). \tag{9}$$

Here, we merge all the sequences  $\{(x_t^{(i)}, y_t^{(i)})\}_{t=1}^T$  for i = 1, ..., n into one sequence  $(\tilde{x}_k, \tilde{y}_k)_{k=1}^{nT}$ . Similarly, we can define a flattened truncated history  $\tilde{H}_k^S$  as

$$\tilde{H}_{iT+t}^{S} \coloneqq (x_{\max\{t-S,1\}}^{(i)}, y_{\max\{t-S,1\}}^{(i)}, \dots, x_{t}^{(i)}, y_{t}^{(i)}.$$
(10)

1188 1189 1190 Note that  $\tilde{H}_{k,k=iT+t}^{S} = (H_t^{S(i)}, y_t^{(i)})$ , since we have added the target label  $y_t^{S(i)}$  into the flattened truncated history  $\tilde{H}_k^S$  for notation simplicity. With a slight abuse of notations, we have

$$\ell_{\theta}(\tilde{H}_{k,k=iT+t}^{S}) \coloneqq \ell(\mathbb{TF}_{\theta}(H_{t}^{S}), y_{t}) = \ell(\mathbb{TF}_{\theta}(H_{t}), y_{t}),$$
(11)

where the equality holds since we are making predictions based on at most S pairs of  $(x_t, y_t)$ . We can similarly replace the  $\ell$  function in the definition of empirical risk r and population risk R, obtaining

 $\begin{aligned} r(\mathrm{TF}(\theta)) &= \frac{1}{nT} \sum_{i=1}^{n} \sum_{t=1}^{T} \ell(\mathrm{TF}_{\theta}(H_t^S), y_t) \\ &= \frac{1}{nT} \sum_{i=1}^{nT} \ell_{\theta}(\tilde{H}_k^S), \end{aligned}$ 

1191 1192

1195

1196 1197 1198

1201 and 1202

 $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{H_t} \left[ \sum_{t=1}^T \ell(TF(\theta)(H_t), y_t) \right]$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{H_k^{S'}} \left[ \sum_{t=1}^T \ell_{\theta}(\tilde{H}_{k,k=iT+t}^{S'}) \right]$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \sum_{t=1}^T \ell_{\theta}(\tilde{H}_{k,k=iT+t}^{S'}) \right]$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \sum_{t=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$   $R(TF(\theta)) = \frac{1}{T} \mathbb{E}_{\tilde{H}_k^{S'}} \left[ \frac{1}{nT} \sum_{k=1}^T \ell_{\theta}(\tilde{H}_k^{S'}) \right],$ 

where  $\tilde{H}_t^{S'}$  is another flattened truncated history that is i.i.d. to  $\tilde{H}_t^S$ . For notation simplicity, we define

$$\tilde{\mathcal{H}}^S \coloneqq (\tilde{H}_1^S, \dots, \tilde{H}_{nT}^S).$$
(14)

(12)

1215 Then we simplify the notations as

$$r_{\theta}\left(\tilde{\mathcal{H}}^{S}\right) \coloneqq r(\operatorname{TF}(\theta)),$$
(15)

1218 and

1214

1216 1217

1219

1220

$$R_{\theta} \coloneqq \mathbb{E}_{\tilde{\mathcal{H}}^{S'}} \left[ r_{\theta} \left( \tilde{\mathcal{H}}^{S'} \right) \right] = R(\mathsf{TF}(\theta)).$$
(16)

To control the difference between  $R_{\theta}$  and  $r_{\theta}(\tilde{\mathcal{H}}^S)$  for any  $\theta$  (which could potentially depend on training data  $\mathcal{D}$ ), we use PAC-Bayes arguments for simplicity.

All the following arguments are made with the conditional distribution on knowing each  $f^{(i)}$  and  $\sigma^{(i)}$ , for each i = 1, ..., n. We omit the conditional dependencies in our notations only for simplicity.

By our definition of data generation, the flattened truncated history  $\tilde{H}_k^S$  naturally forms up a Markov chain on the space  $\otimes_{k=1}^{nT} \Omega_k$  (verified in Lemma E.13), since the newly generated  $(x_t, y_t)$  are conditionally independent of all previous observations. Here  $\Omega_{k,k=iT+t} := (\mathcal{X} \times \mathcal{Y})^{\otimes \min\{t,S\}}$ .

Fix a  $\theta$  that does not depend on the training data  $\mathcal{D}$ . We now bound the difference between  $R_{\theta}$ and  $r_{\theta}(\tilde{\mathcal{H}}^S)$  via concentration inequality for Markov chains. From Lemma F.2, we know that if the Markov chain's mixing time is small enough (which means it quickly converges to the stationary distribution), the concentration properties over the Markov chain would be good enough to enable the standard PAC-Bayes arguments. We also know from Lemma E.15 that the flattened truncated history has a mixing time no greater than min{S, T}, since all the histories S pairs before the current time would be truncated from the input, and the history  $H_t^S$  restarts every time a sequence reaches length T. With these observations, we start our detailed derivation.

Since the function  $\ell$  is almost surely bounded by  $C_2$  as is shown in Lemma E.3, we have almost surely for any  $\tilde{\mathcal{H}}^S$  and  $\tilde{\mathcal{H}}^{S'}$ ,

$$r_{\theta}(\tilde{\mathcal{H}}^S) - r_{\theta}(\tilde{\mathcal{H}}^{S'}) \le \sum_{k=1}^{nT} \frac{2C_2}{nT} \cdot \mathbb{1}\{\tilde{H}_k^S \neq \tilde{H}_k^{S'}\}.$$
(17)

We can use McDiarmid type's inequality for Markov chains (Lemma F.2, with the mixing time upper bound no greater than min{S, T} (specified in Lemma E.15), such that for any  $\lambda \in \mathbb{R}$ ,

$$\mathbb{E}_{\mathcal{D}}\left[\exp\left(\lambda(r_{\theta}(\tilde{\mathcal{H}}^{S}) - R_{\theta}(\tilde{\mathcal{H}}^{S}))\right)\right] \le \exp\left(\frac{2\lambda^{2}C_{2}^{2}\min\{S, T\}}{nT}\right).$$
(18)

<sup>1247</sup> Set  $\pi$  to be the distribution over  $\Theta$  defined in Lemma E.11. Since  $\pi$  is chosen independently from  $\mathcal{D}$ , <sup>1248</sup> we can integrate equation 18 with respect to  $\theta \sim \pi$  such that

$$\mathbb{E}_{\theta \sim \pi} \left[ \mathbb{E}_{\mathcal{D}} \left[ \exp\left( \lambda (r_{\theta}(\tilde{\mathcal{H}}^{S}) - R_{\theta}(\tilde{\mathcal{H}}^{S})) \right] \right] \leq \exp\left( \frac{2\lambda^{2}C_{2}^{2}\min\{S, T\}}{nT} \right).$$
(19)

Using Fubini's Theorem, we can exchange the order of integration, such that

$$\mathbb{E}_{\mathcal{D}}\left[\mathbb{E}_{\theta \sim \pi}\left[\exp\left(\lambda(r_{\theta}(\tilde{\mathcal{H}}^{S}) - R_{\theta}(\tilde{\mathcal{H}}^{S})\right)\right]\right] \leq \exp\left(\frac{2\lambda^{2}C_{2}^{2}\min\{S, T\}}{nT}\right).$$
(20)

By applying Donsker-Varadhan's formula (Lemma F.3), we derive from equation 20 that

$$\mathbb{E}_{\mathcal{D}}\bigg[\exp\bigg(\sup_{\rho\in\mathcal{P}(\Theta)}\big\{\mathbb{E}_{\theta\sim\rho}\big[\lambda(r_{\theta}(\tilde{\mathcal{H}}^{S})-R_{\theta}(\tilde{\mathcal{H}}^{S})\big]-D_{\mathrm{kl}}(\rho\|\pi)\big\}\bigg)\bigg]\leq\exp\bigg(\frac{2\lambda^{2}C_{2}^{2}\min\{S,T\}}{nT}\bigg).$$

1261 Rearranging terms, we have

$$\mathbb{E}_{\mathcal{D}}\left[\exp\left(\sup_{\rho\in\mathcal{P}(\Theta)}\left\{\mathbb{E}_{\theta\sim\rho}\left[\lambda(r_{\theta}(\tilde{\mathcal{H}}^{S})-R_{\theta}(\tilde{\mathcal{H}}^{S})\right]-D_{\mathrm{kl}}(\rho\|\pi)\right\}-\frac{2\lambda^{2}C_{2}^{2}\min\{S,T\}}{nT}\right)\right]\leq1.$$
 (21)

Using Chernoff's bound (Lemma F.4) with probability  $\delta/4$ , we have with probability at least  $1 - \frac{\delta}{4}$  w.r.t.  $\mathcal{D}$ ,

$$\sup_{\rho \in \mathcal{P}(\Theta)} \left\{ \mathbb{E}_{\theta \sim \rho} \left[ \lambda(r_{\theta}(\tilde{\mathcal{H}}^{S}) - R_{\theta}(\tilde{\mathcal{H}}^{S}) \right] - D_{\mathrm{kl}}(\rho \| \pi) \right\} - \frac{2\lambda^{2}C_{2}^{2} \min\{S, S\}}{nT} \leq \log(4/\delta).$$
(22)

Since this bound equation 22 holds for *any* distribution  $\rho$  over  $\Theta$ , we can set  $\rho$  to be  $\rho_{\hat{\theta}^{\text{ERM}}}$  as defined in Lemma E.11, resulting in a high-probability bound

$$\begin{array}{ll} & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}}\left[r_{\theta}(\tilde{\mathcal{H}}^{S}) - R_{\theta}(\tilde{\mathcal{H}}^{S})\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}}\left[r_{\theta}(\tilde{\mathcal{H}}^{S}) - R_{\theta}(\tilde{\mathcal{H}}^{S})\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}}}\left[\pi\right] \\ & \mathbb{E}_{\theta \sim \rho_{\hat{\theta}^{\text{ERM}}$$

1280 By Lemma E.12, the loss function is Lipschitz. Since for any  $\theta \in \operatorname{supp}(\rho_{\hat{\theta}^{\text{ERM}}}), \theta$  is up to  $\mathcal{O}(1/(nT))$ 1281 away from  $\hat{\theta}^{\text{ERM}}$ , we can control the difference between the risks of any  $\theta \in \operatorname{supp}(\rho_{\hat{\theta}^{\text{ERM}}})$  and  $\hat{\theta}^{\text{ERM}}$ 1283 as  $|m_{\theta}(\tilde{\mathcal{U}}^{S}) - m_{\theta}(\tilde{\mathcal{U}}^{S})| \leq \tilde{\mathcal{O}}(1/(nT))$ (24)

$$\left| r_{\theta}(\tilde{\mathcal{H}}^{S}) - r_{\hat{\theta}^{\text{ERM}}}(\tilde{\mathcal{H}}^{S}) \right| \le \tilde{\mathcal{O}}(1/(nT)), \tag{24}$$

$$\left|R_{\theta} - R_{\hat{\theta}^{\text{ERM}}}\right| \le \tilde{\mathcal{O}}(1/(nT)).$$
<sup>(25)</sup>

1286 Thus, we have 1287

$$r_{\hat{\theta}^{\text{ERM}}}(\tilde{\mathcal{H}}^S) - R_{\hat{\theta}^{\text{ERM}}} \le \tilde{\mathcal{O}}(\sqrt{\min\{S,T\}/(nT)}).$$
(26)

Applying the above arguments again for the negative of r, we have with probability at least  $1 - \delta/2$ ,

$$r_{\hat{\theta}^{\text{ERM}}}(\tilde{\mathcal{H}}^S) - R_{\hat{\theta}^{\text{ERM}}} \Big| \le \tilde{\mathcal{O}}(\sqrt{\min\{S,T\}/(nT)}).$$
(27)

1292 For  $\theta^*$ , we can repeat the above steps and get

$$\left| r_{\theta^*}(\tilde{\mathcal{H}}^S) - R_{\theta^*} \right| \le \tilde{\mathcal{O}}(\sqrt{\min\{S,T\}/(nT)}).$$
(28)

The probability that all these bounds hold simultaneously is at least  $1 - \delta$  w.r.t.  $\mathcal{D}$ .

 $\leq r_{\hat{\theta}\text{ERM}}(\tilde{\mathcal{H}}^S) - r_{\theta^*}(\tilde{\mathcal{H}}^S) + \tilde{\mathcal{O}}(\sqrt{\min\{S,T\}/(nT)})$ 

We now take the expectation over each  $f^{(i)}$  and  $\sigma^{(i)}$  to conclude the proof.

 $= r(\mathrm{TF}_{\hat{\theta}\mathrm{FRM}}) - r(\mathrm{TF}_{\theta^*}) + \tilde{\mathcal{O}}(\sqrt{1/n} + \sqrt{S/T})$ 

1296 Hence with probability at least  $1 - \delta$ , 1297  $R(\mathrm{TF}_{\hat{\theta}\mathrm{FRM}}) - R(\mathrm{TF}_{\theta^*})$ 

 $\leq \tilde{\mathcal{O}}(\sqrt{\min\{S,T\}/(nT)})$ 

 $= R_{\hat{a}_{\text{FRM}}} - R_{\theta^*}$ 

1298

1299 1300

1301 1302

- 1303
- 1304

1305

1306

Remark D.1 (Why truncation). Previous analysis (Zhang et al., 2023b) to derive a similar Bayesoptimal argument does not truncate the history and treats the whole history as an inhomogeneous Markov chain. Then they apply the concentration inequalities on Markov chains (for example, Lemma 1309 F.2) to control the difference between R and r. However, their arguments have two limitations: the 1310 first one is that their model is assumed to make decisions based on the full history, which clearly 1311 exceeds the Transformer's model's capacity. The second limitation is that such a concentration argument for Markov chains often relies on upper bounding the mixing time or lower bounding the 1313 spectral gap (for example, Fan et al. (2021)). But Zhang et al. (2023b) do not specify this the mixing 1314 time. Furthermore, in each sampled task sequence (assume we know the task  $f^{(i)}$ ), the mixing time 1315 of the (untruncated) history  $\dot{H}_t$  is infinity: if two sequences start with different initial pairs of  $(x_1, y_1)$ , 1316 then they will never become identical no longer what comes consecutively. Thus, their mixing time 1317 will be T, leading to an  $\tilde{O}(1/\sqrt{n})$  generalization, which is suboptimal if  $S \ll T$  compared to our 1318 result. 1319

(by definition in equation 16)

(by definition in equation 15)

(by equation 27 and equation 28

(by definition of ERM equation 1)

(29)

1320

#### E PROOFS OF LEMMAS 1321

1322

1326

1328

In this section, we prove these lemmas based on the choice of the activation function act =1323 1324 softmax. Similar results for other options act = ReLU can also be found in many existing literatures (for example, see Bai et al. (2024)). 1325

#### 1327 E.1 **BOUNDEDNESS OF TRANSFORMERS**

Lemma E.1 (Layer-wise boundedness). Suppose at the l-th layer of the Transformer, we have 1329  $\|W_V^{m,(l)}\|_{2,2} \leq B_V$  for any  $m = 1, \ldots, M$ ,  $\|A_1^{(l)}\|_{2,2}, \|A_2^{(l)}\|_{2,2} \leq B_A$ . Then for any input  $H^{(l-1)}$ , 1330 we have 1331

$$\|H^{(l)}\|_{2,\infty} \le (1+B_A^2)(1+MB_V)\|H^{(l-1)}\|_{2,\infty}.$$

1332 1333

1337 1338

1339

1340

1341 1342

1344 1345

1347 1348

1349

1334 *Proof of Lemma E.1.* For notation simplicity, we denote  $\operatorname{softmax}((W_K^{(l)}H^{(l-1)})^\top W_Q^{(l)}H^{(l-1)})$  as 1335  $S^m$ . Note that every column of  $S^m$  is of unit 1-norm. Denote each column of  $S^m$  by  $s_t^m$ . For any 1336 input H, we have

 $\| MHA_{W^{(l)}}(H) \|_{2,\infty}$  $\leq \|H\|_{2,\infty} + \sum_{-}^{M} \|W_V^{m,(l)}H\mathbf{S}\|_{2,\infty}$ (by triangle inequality)  $= \|H\|_{2,\infty} + \sum_{t=1}^{M} \max_{t} \|W_{V}^{m,(l)}Hs_{t}^{m}\|_{2}$ (by definition of  $\|\cdot\|_{2,\infty}$ )  $\leq \|H\|_{2,\infty} + \sum_{t=1}^{M} \max_{t} \|W_{V}^{m,(l)}H\|_{2,\infty} \|s_{t}^{m}\|_{1} \quad \text{(by Lemma F.5)}$  $= \|H\|_{2,\infty} + \sum_{l=1}^{M} \|W_{V}^{m,(l)}H\|_{2,\infty}$ (since  $s_t^m$  is of unit 1-norm)

1350  $\leq \|H\|_{2,\infty} + \sum^M \|W^{m,(l)}_V\|_{2,2} \|H\|_{2,\infty}$ 1351 (by Lemma F.6) 1352 1353  $\leq (1 + MB_V) \|H\|_{2,\infty}.$ (by assumption of bounded norm) (30)1354 For any input H, we have 1355  $\| MLP_{A(l)}(H) \|_{2,\infty}$ 1356  $\leq \|H\|_{2,\infty} + \|A_2^{(l)} \operatorname{ReLU}(A_1^{(l)}H)\|_{2,\infty}$ (by triangle inequality) 1358  $\leq \|H\|_{2,\infty} + \|A_2^{(l)}\|_{2,2} \|\operatorname{ReLU}(A_1^{(l)}H)\|_{2,\infty}$ (by Lemma F.6) 1359  $= \|H\|_{2,\infty} + \|A_2^{(l)}\|_{2,2} \max \|\operatorname{ReLU}(A_1^{(l)}H)_{:,t}\|_2$ (by definition of  $\|\cdot\|_{2,\infty}$ ) 1360 1361  $\leq \|H\|_{2,\infty} + \|A_2^{(l)}\|_{2,2} \max \|(A_1^{(l)}H)_{:,t}\|_2$ (since  $|\text{ReLU}(z)| \leq |z|$  for any  $z \in \mathbb{R}$ ) 1362 1363  $= \|H\|_{2,\infty} + \|A_2^{(l)}\|_{2,2} \|A_1^{(l)}H\|_{2,\infty}$ (by definition of  $\|\cdot\|_{2,\infty}$ ) 1364  $\leq \|H\|_{2,\infty} + \|A_2^{(l)}\|_{2,2} \|A_1^{(l)}\|_{2,2} \|H\|_{2,\infty}$ 1365 (by Lemma F.6) 1366  $\leq (1 + B_A^2) \|H\|_{2,\infty}.$ (by assumption of bounded norm) (31)1367 Combining equation 30 and equation 31 yields the conclusion. 1368 1369 **Lemma E.2** (Transformer's boundedness). Suppose  $||W_V^{m,(l)}||_{2,2} \leq B_V$  for any  $m = 1, \ldots, M$ , 1370  $||A_1^{(l)}||_{2,2}, ||A_2^{(l)}||_{2,2} \leq B_A$  for any  $l \in [L]$ . We further assume the projection matrix P is of bounded 1371 norm  $||P^{\top}||_{2,\infty} \leq B_P$ . Then the Transformer's outputs satisfy that 1372  $|\hat{y}(H)| \leq C_1 ||H||_{2,\infty}$ , and  $\exp(-C_1 ||H||_{2,\infty}) \leq \hat{\sigma}(H) \leq 1 + C_1 ||H||_{2,\infty}$ , 1373 where  $C_1 := B_P (1 + B_A^2)^L (1 + M B_V)^L$  is a specified constant. 1374 1375 Proof of Lemma E.2. By Lemma E.1 and a "peeling" argument, we can easily prove that 1376 1377  $\|H^{(L)}\|_{2,\infty} \le (1+B_A^2)^L (1+MB_V)^L \|H^{(0)}\|_{2,\infty}.$ 1378 Thus, 1379  $\|H_{:t}^{(L)}\|_{2} \leq \|H^{(L)}\|_{2,\infty} \leq (1+B_{A}^{2})^{L}(1+MB_{V})^{L}\|H^{(0)}\|_{2,\infty}.$ 1380 Denote P by  $P = [p_1, p_2]^{\top}$ , where  $p_1$  and  $p_2$  are vectors of dimension d. Then the first output 1381 1382  $\hat{y} = p_1^\top H_{\cdot t}^{(L)},$ 1383 where we have (by Cauchy-Schwarz inequality), 1384  $|\hat{y}| \le \|p_1\|_2 \|H_{:t}^{(L)}\|_2 \le B_P (1 + B_A^2)^L (1 + MB_V)^L \|H^{(0)}\|_{2,\infty}.$ 1385 1386 The other output  $\hat{\sigma}$  can be proved similarly as long as one notices 1387  $\log(1 + \exp(-x)) \ge \exp(-x)$ , and  $\log(1 + \exp(x)) \le 1 + x$ , 1388 for any  $x \ge 0$ . 1389 1390 **Lemma E.3** (Boundedness of loss). Under Assumption B.5 with  $\|\theta\| \leq B_{TF}$  and Assumption B.6 1391 with  $||H||_{2,\infty} \leq B_H$  almost surely, we have 1392  $|\ell(TF_{\theta}(H_t), y_t)| \le C_2$ 1393 almost surely, where  $C_2 \coloneqq (C_1 + 1)^2 B_H^2 \cdot \exp(2C_1 B_H) + \max\{C_1 B_H, 1 + \log(C_1 B_H)\}$  is a 1394 specified constant, and  $C_1$  is a constant defined in Lemma E.2. 1395 1396 *Proof of Lemma E.3.* By Lemma E.2, we have 1397  $\frac{(y_t - \hat{y}(H_t))^2}{2\hat{\sigma}^2(H_t)} \le (y_t - \hat{y}_t(H_t))^2 \cdot \frac{\exp(2C_1 B_H)}{2}$ 1398 1399 1400  $< (y_t^2 + \hat{y}_t (H_t)^2) \cdot \exp(2C_1 B_H)$ 1401  $< (C_1 + 1)^2 B_H^2 \cdot \exp(2C_1 B_H),$ 1402

where the second inequality follows from Cauchy's inequality. Combining with a triangle inequality, we have the desired result.  $\Box$ 

# 1404 E.2 LIPSCHITZNESS OF TRANSFORMERS

**Lemma E.4** (Lipschitzness of multi-head attention). Suppose we define the output's norm as  $\|\cdot\|_{2,\infty}$ , the norm of W as

 $||W|| \coloneqq \max\{||W_V^m||_{2,2}, ||W_K^m||_{2,2}, ||W_Q^m||_{2,2}: m = 1, \dots, M\},\$ 

1410and the input H's norm as  $\|\cdot\|_{2,\infty}$ . Suppose at the l-th layer of the Transformer, we have  $\|W^{m,(l)}\| \le B_W$  for any  $m = 1, \ldots, M$ , and  $\|H^{(l-1)}\|_{2,\infty} \le B_H^{(l-1)}$  almost surely. Then  $MHA_{W^{(l)}}(H^{(l-1)})$  is1413 $C_3^{(l)}$ -Lipschitz with respect to  $W^{(l)}$  and  $C_4$ -Lipschitz with respect to  $H^{(l-1)}$  almost surely. Here1414 $C_3^{(l)} \coloneqq 2B_W^2(B_H^{(l-1)})^3 + (B_H^{(l-1)})$  and  $C_4 \coloneqq 1 + MB_W$  are specified constants.

1416 *Proof of Lemma E.4.* We first prove the Lipschitzness result for W. To ease the notations, we omit 1417 the dependence on l and sometimes abbreviate  $W_K^{\top}W_Q$  as  $W_{KQ}$ . For any W and W', using triangle 1418 inequality twice, we have

1422 1423

1424 1425 1426

1427

1428 1429

1456 1457

1415

$$\begin{split} \left\| \mathsf{M}\mathsf{H}\mathsf{A}_{W}(H) - \mathsf{M}\mathsf{H}\mathsf{A}_{W'}(H) \right\|_{2,\infty} \\ &\leq \sum_{m=1}^{M} \left\| W_{V}^{m}H \operatorname{softmax}(H^{\top}W_{KQ}^{m}H) - W_{V}^{m'}H \operatorname{softmax}(H^{\top}W_{KQ}^{m'}H) \right\|_{2,\infty} \\ &\leq \sum_{m=1}^{M} \left\| W_{V}^{m}H \left( \operatorname{softmax}(H^{\top}W_{KQ}^{m}H) - \operatorname{softmax}(H^{\top}W_{KQ}^{m'}H) \right) \right\|_{2,\infty} \\ &+ \sum_{m=1}^{M} \left\| \left( W_{V}^{m} - W_{V}^{m'} \right) H \operatorname{softmax}(H^{\top}W_{KQ}^{m'}H) \right\|_{2,\infty}. \end{split}$$
(32)

We now deal with two terms in equation 32 separately. Since our conclusion will be made for arbitrary  $m \in [M]$ , we omit the dependence on m for notation simplicity from now on.

For the first term, we have

$$\begin{aligned} \|W_{V}H(\operatorname{softmax}(H^{\top}W_{KQ}H) - \operatorname{softmax}(H^{\top}W'_{KQ}H))\|_{2,\infty} \\ &= \max_{t} \|W_{V}H(\operatorname{softmax}(H^{\top}W_{KQ}h_{t}) - \operatorname{softmax}(H^{\top}W'_{KQ}h_{t}))\|_{2} \qquad (by \text{ definition of } \|\cdot\|_{2,\infty}) \\ &\leq \|W_{V}H\|_{2,\infty} \cdot \max_{t} \|(\operatorname{softmax}(H^{\top}W_{KQ}h_{t}) - \operatorname{softmax}(H^{\top}W'_{KQ}h_{t}))\|_{1} \qquad (by \text{ Lemma F.5}) \\ &\leq \|W_{V}\|_{2,2}\|H\|_{2,\infty} \cdot \max_{t} \|(\operatorname{softmax}(H^{\top}W_{KQ}h_{t}) - \operatorname{softmax}(H^{\top}W'_{KQ}h_{t}))\|_{1} \qquad (by \text{ Lemma F.6}) \\ &\leq \|W_{V}\|_{2,2}\|H\|_{2,\infty} \cdot \max_{t} \|H^{\top}W_{KQ}h_{t} - H^{\top}W'_{KQ}h_{t}\|_{\infty} \qquad (by \text{ Lemma F.7}) \\ &\leq 2\|W_{V}\|_{2,2}\|H\|_{2,\infty} \cdot \max_{t} \|H\|_{2,\infty} \|W_{KQ}h_{t} - W'_{KQ}h_{t}\|_{2} \qquad (by \text{ Lemma F.5}) \\ &\leq 2\|W_{V}\|_{2,2}\|H\|_{2,\infty}^{2} \cdot \max_{t} \|W_{KQ} - W'_{KQ}\|_{2,2}\|h_{t}\|_{2} \qquad (by \text{ Lemma F.5}) \\ &\leq 2\|W_{V}\|_{2,2}\|H\|_{2,\infty}^{2} \cdot \max_{t} \|W_{KQ} - W'_{KQ}\|_{2,2}\|h_{t}\|_{2} \qquad (by \text{ Lemma F.5}) \\ &\leq 2\|W_{V}\|_{2,2}\|H\|_{2,\infty}^{2} \cdot \|W_{KQ} - W'_{KQ}\|_{2,2} \|h_{t}\|_{2,\infty} \qquad (by \text{ definition of } \|\cdot\|_{2,\infty}) \\ &\leq 2\|W_{V}\|_{2,2}\|H\|_{2,\infty}^{2} \cdot (\|W_{K}W_{Q} - W'_{KQ}\|_{2,2} + \|W_{K}W'_{Q} - W'_{K}W'_{Q}\|_{2,2}) \qquad (by \text{ triangular inequality}) \\ &\leq 2\|W_{V}\|_{2,2}\|H\|_{2,\infty}^{3} \cdot (\|W_{K}\|_{2,2}\|W_{Q} - W'_{Q}\|_{2,2} + \|W_{K} - W'_{K}\|_{2,2}\|W'_{Q}\|_{2,2}). \qquad (by \text{ sub-multiplicativity of matrix no 1)} \\ &\leq 2B_{W}^{2}(B_{H}^{(1-1)})^{3} \cdot (\|W_{Q} - W'_{Q}\|_{2,2} + \|W_{K} - W'_{K}\|_{2,2}). \qquad (by \text{ sub-multiplicativity of matrix no 1)} \\ &\leq 2B_{W}^{2}(B_{H}^{(1-1)})^{3} \cdot (\|W_{Q} - W'_{Q}\|_{2,2} + \|W_{K} - W'_{K}\|_{2,2}). \qquad (by \text{ sub-multiplicativity of matrix no 1)} \\ &\leq 2B_{W}^{2}(B_{H}^{(1-1)})^{3} \cdot (\|W_{Q} - W'_{Q}\|_{2,2} + \|W_{K} - W'_{K}\|_{2,2}). \qquad (by \text{ sub-multiplicativity of matrix no 1)} \\ &\leq 2B_{W}^{2}(B_{H}^{(1-1)})^{3} \cdot (\|W_{Q} - W'_{Q}\|_{2,2} + \|W_{K} - W'_{K}\|_{2,2}). \qquad (by \text{ sub-multiplicativity of matrix no 1)} \\ &\leq 2B_{W}^{2}(B_{H}^{(1-1)})^{3} \cdot (\|W_{Q} - W'_{Q}\|_{2,2} + \|W_{K} - W'_{K}\|_{2,2}). \qquad (by \text{ sub-multiplicativity of matrix no 1)} \\ &\leq 2B_{W}^{2}(B_{H}^{(1-1)})^{3} \cdot (\|W_{Q} - W'_{Q}\|_{2,2} + \|W_{K} - W'_{K}\|_{2,2}). \qquad (by \text{ sub-multiplicativity of matrix n$$

For notation simplicity, we denote  $softmax(H^{\top}W_{KQ}^{m'}H)$  by S. Note that every column of S is of unit 1-norm. Denote each column of S by  $s_t$ . For the second term, we have

$$\begin{split} & \left\| \left( W_V - W'_V \right) H \texttt{softmax}(H^\top W_{KQ}^{m'} H) \right\|_{2,\infty} \\ &= \| (W_V - W'_V) H \texttt{S} \|_{2,\infty} \end{split}$$
 (by notation substitution)

1458  $= \max_{t} \| (W_V - W'_V) H s_t \|_2$ (by definition of  $\|\cdot\|_{2,\infty}$ ) 1459  $\leq \max_{t} \| (W_V - W'_V) H \|_{2,\infty} \| s_t \|_1$ 1460 (by Lemma F.5) 1461  $= \|(W_V - W'_V)H\|_{2\infty}$ (since  $s_t$  is of unit 1-norm) 1462  $\leq \|W_V - W'_V\|_{2,2} \|H\|_{2,\infty}$ (by Lemma F.6) 1463  $\leq B_{H}^{(l-1)} \| W_{V} - W_{V}' \|_{2.2}.$ 1464 (by assumption of bounded norm) (34)1465 Substituting equation 33 and equation 34 into equation 32, we can conclude that  $MHA_{W^{(l)}}(H^{(l-1)})$  is 1466  $C_3^{(l)}$ -Lipschitz with respect to  $W^{(l)}$  for  $C_3^{(l)} \coloneqq 2B_W^2(B_H^{(l-1)})^3 + (B_H^{(l-1)}).$ 1467 1468 As for the second Lipschitzness conclusion (the one w.r.t. H), it is straightforward if one replaces H1469 with H - H' in the proof of equation 30. 1470 Lemma E.5 (Lipschitzness of multi-layer perceptron). Suppose we define the output's norm as 1471  $\|\cdot\|_{2,\infty}$ , the norm of W as 1472  $||A|| \coloneqq \max\{||A_1||_{2,2}, ||A_2||_{2,2}\},\$ 1473 1474 and the input H's norm as  $\|\cdot\|_{2,\infty}$ . Suppose at the l-th layer of the Transformer, we have  $\|A^{(l)}\| \leq B_A$ 1475 and  $||H||_{2,\infty} \leq B'^{(l-1)}_H$  almost surely. Then  $MLP_{A^{(l)}}(H)$  is  $C_5^{(l)}$ -Lipschitz with respect to  $A^{(l)}$  and 1476  $C_6$ -Lipschitz with respect to H almost surely. Here  $C_5^{(l)} \coloneqq B_A B_H^{\prime (l-1)}$  and  $C_6 \coloneqq 1 + B_A^2$  are 1477 specified constants. 1478 1479 Proof of Lemma E.5. We first prove the Lipschitzness result for A. To ease the notations, we omit 1480 the dependence on l. For any A and A', we have 1481  $\|\operatorname{MLP}_A(H) - \operatorname{MLP}_{A'}(H)\|_{2,\infty}$ 1482  $\leq \left\| (A_2 - A_2') \mathrm{Relu}(A_1 H) \right\|_{2,\infty} + \left\| A_2' (\mathrm{Relu}(A_1 H) - \mathrm{Relu}(A_1')) \right\|_{2,\infty}$ 1483 (by triangle inequality) 1484  $\leq \|A_2 - A_2'\|_{2,2} \|\operatorname{ReLU}(A_1H)\|_{2,\infty}$ 1485 1486  $+ \|A'_2\|_{2,2} \|\text{Relu}(A_1H) - \text{Relu}(A'_1H)\|_{2\infty}$ (by Lemma F.6) 1487  $= \|A_2 - A_2'\|_{2,2} \max_t \|\text{Relu}(A_1H)_{:,t}\|_2$ 1488 1489 +  $||A'_2||_{2,2} \max_{t} ||\operatorname{ReLU}(A_1H)_{:,t} - \operatorname{ReLU}(A'_1H)_{:,t}||_2$ (by definition of  $\|\cdot\|_{2,\infty}$ ) 1490  $\leq \|A_2 - A_2'\|_{2,2} \max_{I} \|(A_1H)_{:,t}\|_2 + \|A_2'\|_{2,2} \max_{I} \|(A_1H)_{:,t} - (A_1'H)_{:,t}\|_2$ 1491 1492 (since  $|\text{ReLU}(z_1) - \text{ReLU}(z_2)| \le |z_1 - z_2|$  for any  $z_1, z_2 \in \mathbb{R}$ ) 1493  $= \|A_2 - A_2'\|_{2,2} \|A_1 H\|_{2,\infty} + \|A_2'\|_{2,2} \|A_1 H - A_1' H\|_{2,\infty}$ (by definition of  $\|\cdot\|_{2,\infty}$ ) 1494  $\leq \|A_2 - A_2'\|_{2,2} \|A_1\|_{2,2} \|H\|_{2,\infty} + \|A_2'\|_{2,2} \|A_1 - A_1'\|_{2,2} \|H\|_{2,\infty}$ 1495 (by Lemma F.6) 1496  $\leq B_A B'^{(l-1)}_H (\|A_1 - A'_1\|_{2,2} + \|A_2 - A'_2\|_{2,2})$ (by assumption of bounded norm) 1497 (35)1498 1499 As for the second Lipschitzness conclusion (the one w.r.t. H), it is straightforward if one replaces H1500 with H - H' in the proof of equation 31. 1501 **Lemma E.6** (Lipshitzness of Transformer). Suppose we define each output's norm as  $|\cdot|$  for  $\hat{y}$  and 1502  $\hat{\sigma}$ , the norm of  $\theta$  as 1503  $\|\theta\| \coloneqq \max\{\|W\|, \|A\|, \|P\|\},\$ 1504 where ||W|| is as defined in Lemma E.4, ||A|| is as defined in Lemma E.5, and  $||P|| := ||P^{\top}||_{2,\infty}$ , and the input H's norm as  $\|\cdot\|_{2,\infty}$ . Suppose we have  $\|\theta\| \le B_{TF}$ , and  $\|H\|_{2,\infty} \le B_H$  almost surely. Then  $\hat{y}_{\theta}(H)$  is  $C_7$ -Lipschitz with respect to  $\theta$ , and  $\hat{\sigma}_{\theta}(H)$  is  $C_8$ -Lipschitz with respect to  $\theta$ . 1506 1507 1508 *Proof of Lemma E.6.* First we quantify the constants  $B_H^{(l-1)}$  in Lemma E.4 and the constants  $B_H^{'(l-1)}$  in Lemma E.5 via Lemma E.1. As is shown in the proof of Lemma E.1, we can define 1509 1510 1511  $B_H^{(l-1)} \coloneqq (1 + M B_{\rm TF})^{l-1} (1 + B_{\rm TF}^2)^{l-1}, \quad l = 1, \dots, L,$ 

1512 and 1513

1514

1519 1520

1523 1524

1545 1546

1550

1559

1560

1563

$$B'_{H}^{(l-1)} \coloneqq (1 + MB_{\mathrm{TF}})^{l} (1 + B_{\mathrm{TF}}^{2})^{l-1}, \quad l = 1, \dots, L,$$

such that all requirements in Lemma E.4 and Lemma E.5 are met almost surely. Thus, we bound the gap between  $H^{(l)}$  (the output of  $TF_{\theta}$  after *l* layers) and  $H'^{(l)}$  (the output of  $TF_{\theta'}$  after *l* layers) by induction. We claim that if  $H^{(0)} = H'^{(0)}$ , then there exists a constant  $C_9^{(l)}$  for any l = 1, ..., L that do not depend on  $\theta$  or H, such that

$$||H^{(l)} - H'^{(l)}||_{2,\infty} \le C_9^l ||\theta - \theta'||_2$$

We prove it by induction. For l = 1, the case can be verified by calculation: by Lemma E.4, 1522

$$\|\operatorname{MHA}_{W^{(1)}}(H^{(0)}) - \operatorname{MHA}_{W'^{(1)}}(H^{(0)})\|_{2,\infty} \le C_3^{(1)} \|\theta - \theta'\|.$$

1525 Similarly, by Lemma E.5,

$$\begin{aligned} \|H^{(1)} - H^{\prime(1)}\|_{2,\infty} &= \|\mathsf{MLP}_{A^{(1)}}(\mathsf{MHA}_{W^{(1)}}(H^{(0)})) - \mathsf{MLP}_{A^{\prime(1)}}(\mathsf{MHA}_{W^{\prime(1)}}(H^{(0)}))\|_{2,\infty} \\ &\leq \|\mathsf{MLP}_{A^{(1)}}(\mathsf{MHA}_{W^{(1)}}(H^{(0)})) - \mathsf{MLP}_{A^{(1)}}(\mathsf{MHA}_{W^{\prime(1)}}(H^{(0)}))\|_{2,\infty} \\ &+ \|\mathsf{MLP}_{A^{(1)}}(\mathsf{MHA}_{W^{\prime(1)}}(H^{(0)})) - \mathsf{MLP}_{A^{\prime(1)}}(\mathsf{MHA}_{W^{\prime(1)}}(H^{(0)}))\|_{2,\infty} \\ &\leq C_{6}\|\mathsf{MHA}_{W^{(1)}}(H^{(0)}) - \mathsf{MHA}_{W^{\prime(1)}}(H^{(0)})\|_{2,\infty} + C_{5}^{(1)}\|\theta - \theta^{\prime}\| \\ &\leq (C_{6}C_{3}^{(1)} + C_{5}^{(1)})\|\theta - \theta^{\prime}\|, \end{aligned}$$
(36)

where we define  $C_9^1$  as  $C_9^1 \coloneqq C_6 C_3^{(1)} + C_5^{(1)}$ . Suppose our conclusion holds for any  $l \le l_0 - 1$ . Then for  $l = l_0$ , we have

$$\begin{aligned} & \text{IS36} \\ & \text{IMHA}_{W^{(l_0)}}(H^{(l_0-1)}) - \text{MHA}_{W'^{(l_0)}}(H'^{(l_0-1)}) \|_{2,\infty} \\ & \leq \|\text{MHA}_{W^{(l_0)}}(H^{(l_0-1)}) - \text{MHA}_{W^{(l_0)}}(H'^{(l_0-1)}) \|_{2,\infty} + \|\text{MHA}_{W^{(l_0)}}(H'^{(l_0-1)}) - \text{MHA}_{W'^{(l_0)}}(H'^{(l_0-1)}) \|_{2,\infty} \\ & \leq C_4 \|H^{(l_0-1)} - H'^{(l_0-1)}\|_{2,\infty} + C_3^{(l_0)} \|\theta - \theta'\| \\ & \leq (C_4 C_9^{(l_0-1)} + C_3^{(l_0-1)}) \|\theta - \theta'\|, \end{aligned}$$

by applying Lemma E.4. We can again compute the difference between  $H^{(l_0)}$  and  $H'^{(l_0)}$  similar to what we do in equation 36 as

$$||H^{(l_0)} - H'^{(l_0)}||_{2,\infty} \le \left(C_6(C_4C_9^{(l_0-1)} + C_3^{(l_0-1)}) + C_5^{(l_0)}\right)||\theta - \theta'||.$$

Hence the induction holds if we define  $C_9^{(l_0)} \coloneqq C_6(C_4C_9^{(l_0-1)} + C_3^{(l_0-1)}) + C_5^{(l_0)}$ . Now we have proved

$$\|H^{(L)} - H'^{(L)}\|_{2,\infty} \le C_9^{(L)} \|\theta - \theta'\|$$

<sup>1551</sup> We shall see from Cauchy-Schwarz inequality that

$$\begin{aligned} |\hat{y} - \hat{y}'| &\leq \|p_1 - p_1'\|_2 \|H^{(L)}\|_{2,\infty} - \|p_1'\|_2 \|H^{(L)} - H'^{(L)}\|_{2,\infty} \\ &\leq \|\theta - \theta'\| (1 + MB_{\mathrm{TF}})^L (1 + B_{\mathrm{TF}}^2)^L + B_{\mathrm{TF}} C_9^{(L)} \|\theta - \theta'\| \\ &= \left( (1 + MB_{\mathrm{TF}})^L (1 + B_{\mathrm{TF}}^2)^L + B_{\mathrm{TF}} C_9^{(L)} \right) \|\theta - \theta'\|, \end{aligned}$$

where the second inequality follows from the proof of Lemma E.6. We can now define

$$C_7 \coloneqq (1 + MB_{\rm TF})^L (1 + B_{\rm TF}^2)^L + B_{\rm TF} C_9^{(L)},$$

and conclude the proof for  $\hat{y}$ . As for  $\hat{\theta}$ , we can see from the fact  $\log(1 + \exp(\cdot))$  is 1-Lipschitz that the Lipschitzness also holds for  $C_8 \coloneqq C_7$ .

**Lemma E.7** (Lipschitzness of loss). Suppose we have  $\|\theta\| \le B_{TF}$  and  $\|H\|_{2,\infty} \le B_H$  almost surely, where the norm of  $\theta$  is the same as defined in Lemma E.6. Then  $\ell(TF_{\theta}(H), y)$  is  $C_{10}$ -Lipschitz with respect to  $\theta$  almost surely. <sup>1566</sup> *Proof of Lemma E.7.* Based on the Lipschitzness of the Transformer w.r.t.  $\theta$  (Lemma E.6), we only need to prove that both partial derivatives  $\frac{\partial \ell}{\partial \hat{y}}$  and  $\frac{\partial \ell}{\partial \hat{\sigma}}$  are bounded. For the first partial derivative, we have

1569 1570

$$\left| \frac{\partial \ell}{\partial \hat{y}} \right| = \left| (y - \hat{y}) \right| \cdot \frac{1}{\hat{\sigma}^2} \le (1 + C_1) B_H \exp(2C_1 B_H). \quad \text{(by Lemma E.2)}$$
(37)

1573 1574 For the second partial derivative, we have

1575 1576

1577

1581 1582

1584

1571 1572

$$\left| \frac{\partial \ell}{\partial \hat{\sigma}} \right| = \frac{\left| - (y - \hat{y})^2 + \sigma^2 \right|}{\hat{\sigma}^3} \le \left( (1 + C_1)^2 B_H^2 + (1 + C_1 B_H)^2 \right) \cdot \exp(3C_1 B_H). \quad \text{(by Lemma E.2)}$$
(38)

<sup>1579</sup> Combining inequalities equation 37 and equation 38 with the Lipschitzness of  $\hat{y}$  and  $\hat{\sigma}$  w.r.t.  $\theta$ , we conclude the result with

$$C_{10} \coloneqq (1+C_1)B_H \exp(2C_1B_H)C_7 + \left((1+C_1)^2B_H^2 + (1+C_1B_H)^2\right)\exp(3C_1B_H)C_8,$$

where  $C_7$  and  $C_8$  are constants that appear in Lemma E.6.

# 1585 E.3 CONSTRUCTING DISTRIBUTIONS OVER PARAMETER SPACE

1587 In this section, we formally define two distributions over the parameter space  $\Theta$ . The first distribution 1588  $\rho_{\hat{\theta}}$  may depend on the empirical distribution, while the second distribution  $\pi_{\theta}$  should be independent 1589 of the training dataset. We control the Kullback-Leibler divergence between  $\rho_{\hat{\theta}}$  and  $\pi_{\theta}$  in Lemma 1590 E.11. For notation simplicity, we may use some notations of different meanings from the main text.

For any dimension d, we denote the Lebesgue measure over  $\mathbb{R}^d$  by  $\lambda_d(\cdot)$ . Then we have the following lemma.

**Lemma E.8** (Upper bound for p.d.f.). Suppose  $\rho$  is the uniform distribution over  $\mathcal{B}(x_0, 3r) \cap \mathcal{B}(0, R)$ for some  $x_0 \in \mathcal{B}(0, R) \subset \mathbb{R}^d$ , where the Lebesgue measure is defined as  $\lambda_d(\cdot)$ , and R > 3r. Then the p.d.f.  $p_{\rho}(\cdot)$  exists and

$$p_{\rho}(x) \leq \frac{1}{\lambda_d (\mathcal{B}(0,r))}.$$

1598

*Proof of Lemma E.8.* Denote the set to be  $S := \mathcal{B}(x_0, 3r) \cap \mathcal{B}(0, R)$ . Since  $\rho$  is the uniform distribution, we just need to prove that

$$\lambda_d(S) \ge \lambda_d(\mathcal{B}(0,r)).$$

This is true because there exists some  $x' \in \mathbb{R}^d$  s.t.  $\mathcal{B}(x', r) \subset S$ . In fact, we can construct the small ball as

$$\mathcal{B}\Big(x_0 - rac{x_0}{\|x_0\|} \cdot 1.5r, r\Big) \subset S.$$

**Lemma E.9** (Upper bound for KL divergence). Suppose the probability space is defined on  $\mathcal{B}(0, R)$ . Suppose  $\rho$  is the uniform distribution over  $\mathcal{B}(x_0, 3r) \cap \mathcal{B}(0, R)$  for some  $x_0 \in \mathcal{B}(0, R) \subset \mathbb{R}^d$ , where the Lebesgue measure is defined as  $\lambda_d(\cdot)$ , and R > 3r. Suppose  $\pi$  is the uniform distribution over  $\mathcal{B}(0, R)$ . Then

$$D_{\mathrm{kl}}(\rho \| \pi) \le \mathcal{O}(C_d \cdot \log(R/r))$$

1614 where  $C_d \coloneqq \log(\lambda_d(\mathcal{B}(0,1)))$  is some constant related to d.

1615

1613

1607

1608

1616 *Proof of Lemma E.9.* Since  $\rho \ll \pi$ , we can define the Radon-Nikodym derivative as  $\frac{d\rho}{d\pi}$ . By Lemma E.8, we can upper bound the RN derivative by

1619 
$$\frac{\mathrm{d}\rho}{\mathrm{d}\pi}(x) = \frac{1/\lambda_d(\mathcal{B}(x_0, 3r) \cap \mathcal{B}(0, R))}{1/\lambda_d \mathcal{B}(0, R)} \le \mathcal{O}(C_d \cdot \log(R/r)).$$

1620 Hence, 1621  $D_{\rm kl}(\rho \| \pi) = \int_{\pi \in \mathcal{B}(0, R)} \log\left(\frac{\mathrm{d}\rho}{\mathrm{d}\pi}(x)\right) \mathrm{d}\rho(x)$ 1622 1623 1624  $\leq \int_{x \in \mathcal{B}(0,R)} \mathcal{O}(C_d \cdot \log(R/r)) \mathrm{d}\rho(x)$ 1625 1626  $= \mathcal{O}(C_d \cdot \log(R/r)).$ 1628 1629 *Remark* E.10. Note that  $C_d = \frac{\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2}+1)}$  is uniformly upper bounded. Here  $\pi$  denotes the ratio of a 1630 circle's circumference to its diameter, and  $\Gamma$  is the Gamma-function. 1631 1632 **Lemma E.11** (Upper bound for  $D_{kl}(\rho_{\hat{\theta}} || \pi_{\theta})$ ). Suppose we are considering probability measures over 1633 the space specified by Assumption **B.6** (that is,  $\Theta = \mathcal{B}(0, B_{TF})$ ). For each layer l and each m, suppose we define the norm over each  $W_Q, W_K \in \mathbb{R}^{d_m \times d}, W_V \in \mathbb{R}^{d \times d}, A_1 \in \mathbb{R}^{d_h \times d}, A_2 \in \mathbb{R}^{d \times d_h}$  to be the 1634 Frobenius norm (that is,  $\|\cdot\|_{2,2}$ ). Suppose  $P = [p_1, p_2]^{\top}$ , and we define the norm over  $p_1, p_2 \in \mathcal{R}^d$  to 1635 be the Euclidean norm. For each layer l and each m, suppose we have the probability measures  $\rho_{\hat{W}_{o}}$ , 1636  $\rho_{\hat{W}_K}$ ,  $\rho_{\hat{W}_V}$ ,  $\rho_{\hat{A}_1}$ ,  $\rho_{\hat{A}_2}$  as the uniform distribution over  $\mathcal{B}(0, 1/(nT)) \cap \mathcal{B}(0, B_{TF})$ ), and the probability 1637 measures  $\pi_{W_Q}$ ,  $\pi_{W_K}$ ,  $\pi_{W_V}$ ,  $\pi_{A_1}$ ,  $\pi_{A_2}$  as the uniform distribution over  $\mathcal{B}(0, B_{TF})$ ). Suppose we have 1638 the probability measures  $\rho_{\hat{p}_1}$ ,  $\rho_{\hat{p}_2}$  as the uniform distribution over  $\mathcal{B}(0, 1/(nT)) \cap \mathcal{B}(0, B_{TF})$ ), and 1639 the probability measures  $\pi_{p_1}$ ,  $\pi_{p_2}$  as the uniform distribution over  $\mathcal{B}(0, B_{TF})$ ). Suppose we define 1640 1641  $\rho_{\hat{\theta}} \coloneqq \left(\bigotimes_{m \ l} \rho_{\hat{W}_Q^{m,(l)}}\right) \otimes \left(\bigotimes_{m \ l} \rho_{\hat{W}_K^{m,(l)}}\right) \otimes \left(\bigotimes_{m \ l} \rho_{\hat{W}_K^{m,(l)}}\right)$ 1642 1643  $\otimes \left(\bigotimes \rho_{\hat{A}_1^{(l)}}\right) \otimes \left(\bigotimes \rho_{\hat{A}_2^{(l)}}\right)$ 1645 1646 (39) $\otimes \rho_{\hat{p}_1} \otimes \rho_{\hat{p}_2},$ 1647 and 1648  $\pi_{\theta} \coloneqq \left(\bigotimes_{m=l} \pi_{W_Q^{m,(l)}}\right) \otimes \left(\bigotimes_{m=l} \pi_{W_K^{m,(l)}}\right) \otimes \left(\bigotimes_{m=l} \pi_{W_V^{m,(l)}}\right)$ 1650 1651  $\otimes \left(\bigotimes_{I} \pi_{A_{1}^{\left(l\right)}}\right) \otimes \left(\bigotimes_{I} \pi_{A_{2}^{\left(l\right)}}\right)$ 1652 1653 1654 (40) $\otimes \pi_{p_1} \otimes \pi_{p_2},$ 1655 where  $\otimes$  represents the product of measures. Then we have 1656  $D_{\mathrm{kl}}(\rho_{\hat{\theta}} \| \pi_{\theta}) \leq \mathcal{O}(C_{11} \log(nTB_{TF})),$ 1657 1658 where  $C_{11}$  is some specified constant that depends polynomially on  $L, M, d, d_m, d_h$ . 1659 *Proof of Lemma* E.11. By setting  $r = \frac{1}{3nT}$  for Lemma E.9 and  $R = B_{\text{TF}}$ , we have for each m =1, ..., M and l = 1, ..., L, 1662  $D_{\mathrm{kl}}(\rho_{\hat{W}^{m,(l)}_{\mathcal{O}}} \| \pi_{W^{m,(l)}_{\mathcal{O}}}) \leq \mathcal{O}\big(C_{dd_m} \log(nTB_{\mathrm{TF}})\big),$ 1663 1664  $D_{\mathrm{kl}}(\rho_{\hat{W}^{m,(l)}_{\kappa}} \| \pi_{W^{m,(l)}_{\kappa}}) \leq \mathcal{O}\big(C_{dd_m} \log(nTB_{\mathrm{TF}})\big),$ 1665  $D_{\mathrm{kl}}(\rho_{\hat{W}^{m,(l)}} \| \pi_{W^{m,(l)}}) \leq \mathcal{O}\big(C_{d^2} \log(nTB_{\mathrm{TF}})\big),$  $D_{\mathrm{kl}}(\rho_{\hat{A}_{*}^{(l)}} \| \pi_{A_{*}^{(l)}}) \leq \mathcal{O}\big(C_{dd_{h}} \log(nTB_{\mathrm{TF}})\big),$  $D_{\mathrm{kl}}(\rho_{\hat{A}_{\alpha}^{(l)}} \| \pi_{A_{\alpha}^{(l)}}) \leq \mathcal{O}\big(C_{dd_{h}} \log(nTB_{\mathrm{TF}})\big),$ 1670  $D_{\mathrm{kl}}(\rho_{\hat{p}_1} \| \pi_{p_1}) \le \mathcal{O}\big(C_d \log(nTB_{\mathrm{TF}})\big),$ 1671 1672  $D_{\mathrm{kl}}(\rho_{\hat{p}_2} \| \pi_{p_2}) \le \mathcal{O}\big(C_d \log(nTB_{\mathrm{TF}})\big).$ 1673 By Lemma F.8, we can sum up the above inequalities and get the final result. 

1674 **Lemma E.12** (Bounded difference). For any  $\hat{\theta} \in \Theta$ , suppose we construct the distribution  $\rho_{\hat{\alpha}}$  as in 1675 equation 39. Then for any  $\theta \in \text{supp}(\rho_{\hat{a}})$ , under Assumption B.6 and Assumption B.5, we have 1676

1677 1678

1680

1683 1684

1686

1688

 $\left|\ell(\mathrm{TF}_{\theta}(H), y) - \ell(\mathrm{TF}_{\hat{\theta}}(H), y)\right| \leq \mathcal{O}\big(C_{10}/(nT)\big),$ 

1679 almost surely. Here  $C_{10}$  is the same as defined in Lemma E.7.

1681 *Proof of Lemma E.12.* By construction shown in equation 39, we can see that for any  $\theta \in \text{supp}(\rho_{\hat{a}})$ , 1682

$$\|\theta - \hat{\theta}\| < 1/(nT).$$

Then from the Lipschitzness of the loss function w.r.t.  $\theta$  (Lemma E.7), we conclude the proof. 1685

1687 E.4 MARKOV CHAIN'S PROPERTY

**Lemma E.13** ( $\hat{\mathcal{H}}^S$  is a Markov chain (conditioned on knowing f and  $\sigma$ )). Suppose we have  $\hat{\mathcal{H}}^S$ 1689 defined as equation 14. Then  $\tilde{\mathcal{H}}^S$  is a Markov chain conditioned on knowing each  $f^{(i)}$  and  $\sigma^{(i)}$  for 1690 *each* i = 1, ..., n*.* 

1693 *Proof of Lemma E.13.* By definition, the state of  $\tilde{\mathcal{H}}^S$  will restart and does not depend on all previous histories once  $\hat{H}_k^k$ 's index k reaches the point of k = iT + 1. Therefore, we only need to verify that 1695 inside each task's sequence, the state  $\tilde{H}_k^S$  is also Markovian. 1696

Suppose k = iT + t for some *i*, and we considering  $k = iT + 1, \dots, iT + T$  for each  $t = 1, \dots, T$ . 1697 We write  $\tilde{H}_k^S$  and  $(x_{\max\{1,t-S\}}, y_{\max\{1,t-S\}}, \ldots, x_t, y_t)$  interchangeably for notation simplicity. 1698

1699 Each pair of  $(x_t, y_t)$  is now independent conditioned on knowing the underlying  $f^{(i)}$  and  $\sigma^{(i)}$ . 1700 We omit the conditional dependencies on  $f^{(i)}$  and  $\sigma^{(i)}$  for notation simplicity. The p.d.f. of  $\tilde{H}_k^S$ 1701 conditioned on observing  $\{\tilde{H}_{\tau}^S\}_{\tau=iT+1}^{iT+t}$  and knowing  $f^{(i)}$  and  $\sigma^{(i)}$  is 1702

1703 1704

1705 1706

1707

1708

 $p(x_{\max\{1,t-S\}}, y_{\max\{1,t-S\}}, \dots, x_t, y_t | \{\tilde{H}_{\tau}^S\}_{\tau=iT+1}^{iT+t} = \{\tilde{H}_{\tau}^{S'}\}_{\tau=iT+1}^{iT+t}, f = f^{(i)}, \sigma = \sigma^{(i)})$  $= \mathbb{1}\{x_{\max\{1,t-S\}} = x'_{\max\{1,t-S\}}, \dots, y_{t-1} = y'_{t-1}\} \cdot p(x_t, y_t | f = f^{(i)}, \sigma = \sigma^{(i)})$ (by conditional independence of each pair of  $(x_{\tau}, y_{\tau})$ )  $= p(x_{\max\{1,t-S\}}, y_{\max\{1,t-S\}}, \dots, x_t, y_t | \tilde{H}_{t-1}^S = \tilde{H}_{t-1}^{S'}, f = f^{(i)}, \sigma = \sigma^{(i)})$ Thus the Markovian property holds. 

1709 1710

1711 We present the definition of *mixing time* as used in Paulin (2015). 1712

**Definition E.14** (Mixing time for inhomogeneous Markov chains). Let  $X_1, \ldots, X_N$  be a Markov 1713 chain with Polish state space  $\Omega_1 \times \cdots \times \Omega_N$  (that is,  $X_i \in \Omega_i$ ). Let  $\mathcal{L}(X_{i+t}|X_i = x)$  be the conditional 1714 distribution of  $X_{i+t}$  given  $X_i = x$ . Let us denote the minimal t such that  $\mathcal{L}(X_{i+t}|X_i = x)$  and 1715  $\mathcal{L}(X_{i+t}|X_i = y)$  are less than  $\epsilon$  away in total variational distance for every  $1 \leq i \leq N-t$  and 1716  $x, y \in \Omega_i$  by  $\tau(\epsilon)$ , that is, for  $0 < \epsilon < 1$ , let 1717

$$\bar{d}(t) \coloneqq \max_{1 \le i \le N-t} \sup_{x, y \in \Omega_i} \operatorname{TV}(\mathcal{L}(X_{i+t}|X_i=x), \mathcal{L}(X_{i+t}|X_i=y)),$$
  
$$\tau(\epsilon) \coloneqq \min\{t \in \mathbb{N} : \bar{d}(t) \le \epsilon\}.$$

1719 1720 1721

1726 1727

1722 We now upper bound the mixing time of  $(H_t^S, y_t)$ .

-

1723 Lemma E.15 (Mixing time for truncated history). Suppose we are considering the conditional 1724 distribution on knowing each  $f^{(i)}$  and  $\sigma^{(i)}$ . Then for the Markov chain  $\tilde{\mathcal{H}}_k^S$ , we have 1725

$$\tau(\epsilon) \le \min\{S, T\}$$

for any  $\epsilon \in [0, 1)$ .

1728Proof of Lemma E.15. We first consider the case when  $S \leq T$ . The mixing property inside each<br/>sequence  $\tilde{H}_k^S$  for  $k = iT + 1, \ldots, iT + T$ . Since each (x, y) is i.i.d. distributed conditioned on know-<br/>ing  $f^{(i)}$  and  $\sigma^{(i)}$ , the conditional distribution of the consecutive sequence  $(x_{t+1}, y_{t+1}), \ldots, (x_T, y_T)$ <br/>is never affected by previous t pairs  $(x_1, y_1), \ldots, (x_t, y_t)$  for any  $1 \leq t \leq T$ . We consider the<br/>conditional distribution on knowing  $f^{(i)}$  and  $\sigma^{(i)}$  from now on and omit the dependencies for notation<br/>simplicity.

For any  $1 \le t \le T - S$ , for any two points  $\tilde{H}_{iT+t}^{S'} \ne \tilde{H}_{iT+t}^{S''}$ , the distribution of  $\tilde{H}_{iT+t+S}^S$  is independent of previous t pairs of observed samples. In other words,

$$\mathcal{L}(\tilde{H}_{iT+t+t'}^S | \tilde{H}_{iT+t}^S = \tilde{H}_{iT+t}^{S\prime}) = \mathcal{L}(\tilde{H}_{iT+t+t'}^S | \tilde{H}_{iT+t}^S = \tilde{H}_{iT+t}^{S\prime\prime}),$$

1739 for any  $t' \ge S$ . Hence,

 $\bar{d}(t) = 0$ , for any  $t \ge S$ .

1741 We have 1742

1737 1738

1740

1743

1747 1748

1749

 $\tau(\epsilon) \leq S$ , for any  $\epsilon \in [0, 1)$ .

1744 When S > T, note that the flattened (truncated) history  $\tilde{H}_k^S$  restarts every time it meets the end of a 1745 sequence generated by some  $f^{(i)}$  and  $\sigma^{(i)}$ . Since the length of those sequences is T, we have 1746

$$\tau(\epsilon) \leq T$$
, for any  $\epsilon \in [0, 1)$ 

#### 1750 1751 F TECHNICAL LEMMAS

In this section, we present some technical lemmas. Note that all the notations in this section are chosen for simplicity and may have different meanings than those in other sections.

**Lemma F.1** (McDiarmid's inequality (McDiarmid et al., 1989)). Let  $X = (X_1, ..., X_N)$  be a vector of independent random variables taking values in a Polish space  $\Lambda = \Lambda_1 \times \cdots \times \Lambda_N$ . Suppose that  $f: \Lambda \to \mathbb{R}$  satisfies

1760

1762 1763

1767 1768 1769

1771 1772 1773

$$f(x) - f(y) \le \sum_{i=1}^{N} c_i \mathbb{1}\{x_i \neq y_i\}$$

1761 for any  $x, y \in \Lambda$ . Then for any  $\lambda \in \mathbb{R}$ ,

$$\mathbb{E}\Big[\exp\left(\lambda(f(X) - \mathbb{E}[f(X)])\right)\Big] \le \frac{\lambda^2 \|c\|_2^2}{2}$$

**Lemma F.2** (Corollary 2.11 in Paulin (2015)). Let  $X = (X_1, ..., X_N)$  be a Markov chain taking values in a Polish space  $\Lambda = \Lambda_1 \times \cdots \times \Lambda_N$ , with mixing time  $\tau(\epsilon)$  for  $0 \le \epsilon < 1$ . Define

$$\tau_{\min} \coloneqq \inf_{\epsilon \in [0,1)} \tau(\epsilon) \cdot \left(\frac{2-\epsilon}{1-\epsilon}\right)^2.$$

1770 Suppose that  $f: \Lambda \to \mathbb{R}$  satisfies

$$f(x) - f(y) \le \sum_{i=1}^{N} c_i \mathbb{1}\{x_i \neq y_i\},\$$

for any  $x, y \in \Lambda$ . Then for any  $\lambda \in \mathbb{R}$ ,

$$\mathbb{E}\Big[\exp\big(\lambda(f(X) - \mathbb{E}[f(X)])\big)\Big] \le \frac{\lambda^2 \tau_{\min} \|c\|_2^2}{8}$$

**Lemma F.3** (Donsker-Varadhan variational formula (Donsker & Varadhan, 1983)). Let P and Q be two probability distributions over  $(\Theta, \mathcal{F})$ . If  $Q \ll P$ , then for any real-valued function h integrable w.r.t. P, **1781 1781** 

$$\log \mathbb{E}_P[\exp h] = \sup_{Q \ll P} \{\mathbb{E}_Q[h] - D_{\mathrm{kl}}(Q \| P)\}.$$

1775 1776

**Lemma F.4** (Chernoff's bound (Chernoff, 1952)). For any random variable X, if  $\mathbb{E}[\exp(X)] \le 1$ , then for any  $\delta \in (0, 1)$ ,  $\mathbb{P}(X \le \log(1/\delta)) \ge 1 = \delta$ 

$$\mathbb{P}(X \le \log(1/\delta)) \ge 1 - \delta.$$

**Lemma F.5** (Lemma M.7 in Zhang et al. (2022)). Given any two conjugate numbers  $p, q \in [1, \infty]$ s.t. 1/p + 1/q = 1, for any  $r \in [1, \infty]$ , we have

$$||Ax||_r \le ||A||_{r,p} ||x||_q$$
, and  $||Ax||_r \le ||A^{\top}||_{p,r} ||x||_q$ 

for any matrix  $A \in \mathbb{R}^{m \times n}$  and vector  $x \in \mathbb{R}^n$ .

**Lemma F.6** (Lemma M.8 in Zhang et al. (2022)). *Given any two conjugate numbers*  $p, q \in [1, \infty]$ *s.t.* 1/p + 1/q = 1, we have

$$||AB||_{p,\infty} \le ||A||_{p,q} ||B||_{p,\infty}$$

for any matrix  $A \in \mathbb{R}^{m \times n}$  and matrix  $B \in \mathbb{R}^{n \times r}$ .

**Lemma F.7** (Lemma M.9 in Zhang et al. (2022)). Given any two vectors  $x, y \in \mathbb{R}^d$ , we have

1797 1798

1799

1801 1802 1803

1805

1812 1813

1815

1816 1817

1818

1819 1820 1821

1823

1824

1825

1826 1827

1830

1831

1788 1789

1793

 $\|softmax(x) - softmax(y)\|_1 \le 2\|x - y\|_{\infty}.$ 

**Lemma F.8** (Property of Kullback-Leibler divergence, Proposition 7.2 in Polyanskiy & Wu (2024)). Given any two probability distributions  $\mu_1$  and  $\mu_2$  over  $(\Omega, \mathcal{F})$  and any two distributions  $\nu_1$  and  $\nu_1$ over  $(\Omega', \mathcal{F}')$ , if  $\mu_1 \ll \mu_2$  and  $\nu_1 \ll \nu_2$ , then we have

$$D_{\rm kl}(\mu_1 \otimes \nu_1 \| \mu_2 \otimes \nu_2) = D_{\rm kl}(\mu_1 \| \mu_2) + D_{\rm kl}(\nu_1 \| \nu_2)$$

### G EXPERIMENT DETAILS

#### 1806 1807 G.1 TRAINING DATA GENERATION

We first describe a basic setup of all our experiments. For some experiments, we change some part(s) in below to design the corresponding "flipped" experiment or to examine the OOD ability of the trained transformer. In particular, the *i*-th thread the training data

$$\left(x_{1}^{(i)}, y_{1}^{(i)}, x_{2}^{(i)}, y_{2}^{(i)}, ..., x_{T}^{(i)}, y_{T}^{(i)}
ight)$$

1814 is generated by the following distributions:

•  $\mathcal{P}_X$ : the feature vector  $x_t^{(i)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I_d)$  where  $I_d$  is d-dimensional identity matrix.

•  $\mathcal{P}_{\epsilon}$ : the noise  $\epsilon_t^{(i)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0,1)$ .

•  $\mathcal{P}_{\sigma}$ : the noise intensity  $\sigma_i$  is sampled i.i.d. from

$$\tau_i \sim \text{Gamma}(\underline{\tau}, \overline{\tau}), \quad \sigma_i = \frac{1}{\sqrt{\tau_i}}$$

where the parameters  $\underline{\tau} = \overline{\tau} = 20$  for the basic setup of the experiment. We change these two parameters for some OOD experiments.

•  $\mathcal{P}_{\mathcal{F}}$ : The function  $f_i(x) \coloneqq w_i^\top x$  where  $w_i$  is generated from

$$w_i | \sigma_i \sim \mathcal{N}(\bar{w}, \sigma_i^2 \cdot I_d)$$

where  $I_d$  is the *d*-dimension identity matrix and  $\bar{w}$  is set to be an all-one vector of dimension *d*. The covariance matrix of  $w_i$  is related with the noise intensity  $\sigma_i$  to control the signal-to-noise ratio.

Finally, the target variable is calculated by

1834 
$$y_t^{(i)} = w_i^{\top} x_t^{(i)} + \sigma_i \epsilon_t^{(i)}.$$

Throughout the paper, we consider the dimension d = 8.

#### 1836 G.2 NUMBER OF TASKS N AND TRAINING PROCEDURE

1838 In the previous Section G.1, we define how we generate the training data. As in the previous work, we introduce the notion of *task* where each realization of  $(w_i, \sigma_i)$  is referred to as one task. The rationale is that each configuration of  $(w_i, \sigma_i)$  corresponds to one pattern of the sequence  $(x_t, y_t)$ 's. 1840 While the distribution of  $(w_i, \sigma_i)$  corresponds to infinitely many possible task configurations, we use 1841 a finite pool of tasks for training the Transformer. Specifically, we generate 1842

$$\mathcal{T} \coloneqq \{(w_i, \sigma_i)\}_{i=1}^N$$

from the distributions discussed above. Throughout the paper, we use N to refer to the total number 1845 of tasks or the pool size. 1846

1847 Training the Transformer for our setting is slightly different from the classic ML model's training. 1848 We do not use a fixed set of training data. Rather, we generate a new batch of training data freshly for 1849 each batch.

- The batch size b = 64. For each batch, we first sample with replacement b tasks from the task pool  $\mathcal{T}$ . And based on each sampled  $(w_i, \sigma_i)$ , we generate a training sequence  $\left(x_1^{(i)}, y_1^{(i)}, x_2^{(i)}, y_2^{(i)}, \dots, x_T^{(i)}, y_T^{(i)}\right)$  following the setting in the previous Section G.1.

1861

1863

1850

1851

1843

1844

- 1855
- All the numerical experiments in our paper run for 200,000 batches.

The validation and testing sets are also randomly generated instead of fixed beforehand. But unlike 1857 the training phase which draws the task configuration from the task pool  $\mathcal{T}$ , the validation and test phase samples  $(w_i, \sigma_i)$  directly from the original distribution described in the previous Section G.1. 1859 This is aimed to validate or test whether the trained model has learned the ability to solve a family of problems, or it only just memorizes a fixed pool of tasks  $\mathcal{T}$ . 1860

#### G.3 DERIVATION OF BAYES-OPTIMAL PREDICTOR 1862

In Proposition 3.1, we state the Bayes-optimal predictor in the form of a posterior expectation. Now 1864 we calculate the Bayes-optimal predictor explicitly under the generation mechanism specified in 1865 Section G.1. Conditional on history  $H_t = (x_1, y_1, \dots, x_t)$ , the posterior distribution of  $(w, \sigma)$  that 1866 governs the generation of  $H_t$  can be calculated based on the Bayesian posterior as 1867

1871

1873 1874 1875

 $\mathbb{P}(\tau|H_t) = \text{Gamma}(\tau; \underline{\tau}_t, \overline{\tau}_t), \quad \sigma = \frac{1}{\sqrt{\tau}},$ 

$$\mathbb{P}(w|\sigma, H_t) = \mathcal{N}(w_t, \sigma^2 \cdot \Sigma_t).$$

where 1872

$$\Sigma_{t} = \left(I_{d} + \sum_{s=1}^{t-1} x_{s} x_{s}^{\top}\right)^{-1}, \qquad w_{t} = \Sigma_{t} \left(\bar{w} + \sum_{s=1}^{t-1} x_{s} y_{s}\right)$$
$$\underline{\tau}_{t} = \underline{\tau} + \frac{t}{2}, \qquad \qquad \bar{\tau}_{t} = \bar{\tau} + \frac{1}{2} \sum_{s=1}^{t-1} \left(y_{s}^{2} + \bar{w}^{\top} \bar{w} - w_{s}^{\top} \Sigma_{t}^{-1} w_{s}\right).$$

1876 1877 1878

1881

1884 1885

1879 Accordingly, the Bayes-optimal predictor becomes 1880

$$y_t^*(H_t) = \mathbb{E}[y_t|H_t] = w_t^{\top} x_t,$$

and

$$\sigma_t^{*2}(H_t) = \mathbb{E}[(y_t - y_t^*(H_t))^2 | H_t] = \mathbb{E}[(f(x_t) - y_t^*(H_t))^2 | H_t] + \mathbb{E}[\sigma^2 | H_t]$$
  
=  $\frac{\bar{\tau}_t}{\tau_t - 1} \cdot (\operatorname{tr}(x_t x_t^\top \Sigma_t) + 1).$ 

These formulas are used to generate the Bayes-optimal curves in the figures.