

Appendix for “What and How does In-Context Learning Learn? Bayesian Model Averaging, Parameterization, and Generalization”

A CONCLUSION

In this paper, we investigated the theoretical foundations of ICL for the pretrained language models. We proved that the perfectly pretrained LLMs implicitly implements BMA with regret $\mathcal{O}(1/t)$ over a general response generation modeling, which subsumes the models in previous works. Based on this, we showed that the attention mechanism parameterizes the BMA algorithm. Analyzing the pretraining process, we demonstrated that the total variation between the pretrained model and the nominal distribution consists of the approximation error and the generalization error. The combination of the ICL regret and the pretraining performance gives the full description of ICL ability of pretrained LLMs. We mainly focus on the prompts that comprise several examples in this work and leave the analysis of instruction-based prompts for future works.

B MORE RELATED WORKS

Transformers. Our work is also related to the works that theoretically analyze the performance of transformers. For the analytic properties of transformers, [Vuckovic et al. \(2020\)](#) proved that attention is Lipschitz-continuous via the view of interacting particles. [Noci et al. \(2022\)](#) provided the theoretical justification of the rank collapse phenomenon in transformers. [Yun et al. \(2019\)](#) demonstrated that transformers are universal approximators. For the statistical properties of transformers, [Malladi et al. \(2022\)](#), [Hron et al. \(2020\)](#), and [Yang \(2020\)](#) analyzed the training of transformers within the neural tangent kernel framework. [Wei et al. \(2022a\)](#) presented the approximation and generalization bounds for learning boolean circuits and Turing machines with transformers. [Edelman et al. \(2021\)](#) and [Li et al. \(2023\)](#) derived the generalization error bound of transformers. In our work, we analyze transformers from both the analytic and statistical sides. We show that attention essentially implements the BMA algorithm in the ICL setting. Furthermore, we derive the approximation and generalization bounds for transformers in the pretraining phase.

Generalization. Our analysis of the pretraining is also related to the generalization analysis of the neural networks. This topic has attracted a lot of interests for a long time. [Anthony et al. \(1999\)](#) derived the uniform generalization bound for fully-connected neural networks with the help of VC dimension. [Bartlett et al. \(2017\)](#) sharpened this generalization bound for classification problem by adopting the Dudley’s integral and calculating of the covering number of neural network class. At the same time, [Neyshabur et al. \(2017\)](#) derived a similar as [Bartlett et al. \(2017\)](#) from PAC-Bayes framework. Following this line, [Liao et al. \(2020\)](#), [Ledent et al. \(2021\)](#) and [Lin and Zhang \(2019\)](#) built the generalization bound for graph neural networks and convolutional neural network. These results respected the underlying graph structure and the translation-invariance in the networks. [Edelman et al. \(2021\)](#) established the generalization bound for transformer, but this result did not reflect the permutation-invariance, still depending on the channel number. Our work focuses on the analysis of Maximum Likelihood Estimate (MLE) with transformer function class, which is not covered by previous works. Our bounds are sharper than that of [Edelman et al. \(2021\)](#) on the channel number dependency.

C EXPERIMENTAL RESULTS

We conduct five experiments to verify our theoretical findings, including the Bayesian view (Propositions 4.1 and (4.7)), the regret upper-bounded in Corollary 4.2 and Theorem 6.2, and the constant ratio between attn_\dagger and attn in Proposition 4.3. The implementation details are provided in Appendix D.

C.1 VERIFICATION OF THE BAYESIAN VIEW

To verify the Bayesian view that we adopt in the paper, we implemented two experiments. In the first experiment, we explicitly construct the hidden concept vectors that are found by LLMs.

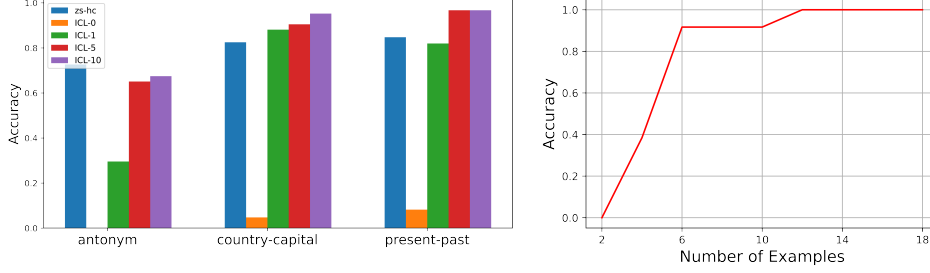


Figure 1: Accuracies of LLMs with and without explicit hidden concepts.

Figure 2: Accuracy of LLMs to find the best arm in the bandit instance with an informative arm.

Motivated by (4.7), we construct the hidden concept vector as the average sum over prompts of the values of twenty selected attention heads, i.e., we compress the hidden concept into a vector with dimension 4096. To demonstrate the effectiveness of the constructed hidden concepts, we add these hidden concept vectors at a layer of LLMs when the model resolves the prompt with zero-shot. In Figure 1, “zs-hc” refers to the results of LLMs that infers with learned hidden concept vectors and zero-shot prompt, and “ICL- i ” refers to the results of LLMs prompted with i examples. We consider the tasks of finding antonyms, finding the capitals of countries, and finding the past tense of words. The results indicate that the LLMs with learned hidden concept vectors have comparable performance with the LLMs prompted with several examples. This indicates that the learned hidden concept vectors are indeed efficient compression of the hidden concepts, which proves that LLMs deduce hidden concepts for ICL. This result strongly corroborates with (4.7).

In the second experiment, we aim to verify that LLMs implement inference with the Bayesian framework, not with gradient descent (Akyürek et al., 2022; von Oswald et al., 2022; Bai et al., 2023) on some tasks. We prompt the LLMs with the history data of a set of similar multi-armed bandit instances with 100 arms, and let LLMs indicate which arm to pull in a similar new bandit instance. In these similar bandit instances, there is an informative arm, whose reward is exactly the index of the arm with the highest rewards. We also provide the side information that “Some arm may directly tell you the arm with the highest reward, even itself does not have the highest reward”. In each example provided in the prompt, there are the rewards of six arms, including the informative arm and the best arm, in one bandit instance. As shown in Figure 2, the LLMs can efficiently implement ICL even with only 6 examples. We note that the gradient descent algorithms in the previous works cannot explain this performance, since the gradient descent algorithms need at least 100 data points, where each data point is the reward of one arm, to learn. In contrast, the Bayesian view can clearly explain Figure 2, where LLMs make use of the side information to calculate a better posterior for ICL.

C.2 VERIFICATION OF THE REGRET BOUND

To verify Corollary 4.2 and Theorem 6.2, we implement experiments to evaluate the regret in two settings. In the first setting, the LLMs is trained for the linear regression task from scratch, which is a representative setting studied in Garg et al. (2022); Akyürek et al. (2022). The examples in the prompt are $\{(x_i, y_i)\}_{i=1}^N$, where $x_i \in \mathbb{R}^d$, $d = 20$ and $y_i = w^T x_i$ for some w sampled from Gaussian distribution. Given the Gaussian model, we adopt the squared error to approximate the logarithm of the probability. Then the $t \times$ regret of the LLMs can be well approximated by the sum of the squared error till time t . The results in Figure 3 strongly corroborate our theoretical findings. First, the results verify our claim in Corollary 4.2 and Theorem 6.2 that $t \cdot$ regret can be upper bounded by a constant. Second, the line of squared error indicates that the ICL of LLMs only has a significant error when $T \leq d$, i.e., the regret only increases in this region. Thus, the regret of the ICL by LLMs is at most linear in $O(d/T)$. From the view of our theoretical result, discretizing the set $\{z \in \mathbb{R}^d \mid \|z\|_2 \leq d\}$ with approximation error $\delta > 0$ will result in a set with $(C/\delta)^d$ elements, where $C > 0$ is an absolute constant. Corollary 4.2 and Theorem 6.2 imply that the regret is the sum of the $\log |3|/T = d \log(C/\delta)/T$ and the pretraining error, which matches the simulation results.

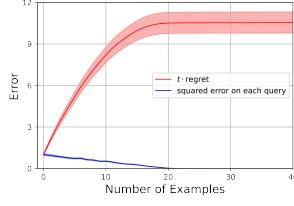


Figure 3: Squared error and regret of LLMs trained for linear log-likelihood regression.

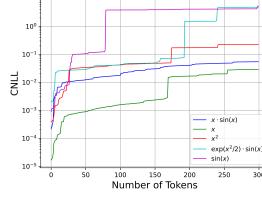


Figure 4: Cumulative negative log-likelihood of LLMs for function value prediction.

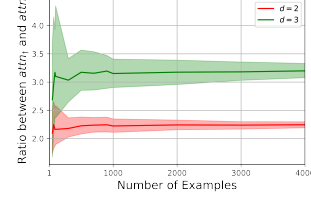


Figure 5: The ratio between pretrained attn_+ and attn .

In the second experiment, we directly evaluate the regret of pretrained LLMs on the function value prediction task. The prompt consists of the values of a function on the points with fixed intervals. Since the values are real numbers, we adopt the method in Gruver et al. (2023) to transfer a real number to a token sequence. For the pretrained model, we cannot calculate $\mathbb{P}(r_i | \text{prompt}_{i-1}, z)$ due to the unknown nominal distributions. Thus, we calculate the cumulative negative log-likelihood $\text{CNLL}_t = -\sum_{i=1}^t \hat{\mathbb{P}}(r_i | \text{prompt}_{i-1})$, and CNLL_t/t is an upper bound of the regret. In Figure 4, we indicate the cumulative negative log-likelihoods of predicting the values of five functions. The results show that the cumulative negative log-likelihoods are stepped, which means that the cumulative negative log-likelihoods are upper-bounded by constants in a long period. This corroborates with Corollary 4.2 and Theorem 6.2. In addition to the mentioned property, we also observe that there are connections between the cumulative negative log-likelihood and the prediction error. We let the LLMs to predict the value given the prompt that contains the past values. Figures 6 and 7 show that the larger cumulative negative log-likelihood implies a larger prediction error.

C.3 VERIFICATION OF THE CONSTANT RATIO BETWEEN attn_+ AND attn

To verify Proposition 4.3, we directly calculate the ratio between attn_+ and attn . We consider the case $d_v = 1$ and $d_k = d$ for some $d > 0$. The entries in K of (4.8) are i.i.d. samples of Gaussian distribution, and the i -th entry of V is calculated as the inner product between a Gaussian vector and the i -th column. Figure 5 shows the results for $d = 2$ and $d = 3$. It shows that the ratio between attn_+ and attn will converge to a constant. This constant depends on the dimension d , which originates from Proposition F.1.

D IMPLEMENTATION DETAILS OF EXPERIMENTS

In this section, we provide the implementation details of the experiments. In the hidden concepts construction experiment, we explicitly calculate the hidden concept vector for Llama2-7b with the method in Todd et al. (2023). Given the prompts generated from the same hidden concept, we calculate the average value of each attention head by prompting the LLM with different prompts. Then we select the attention head according to its average indirect effect, which is defined in Todd et al. (2023). The hidden concept vector is the sum of the average value of the selected attention heads. We test the performance of the learned hidden concept vectors on tasks: (1) Antonym: Given an input word, generate the word with the opposite meaning. (2) Country-Capital. Given a country name, generate the capital city. (3) Present-Past. Given a verb in the present tense, generate the verb’s simple past inflection. To test the effectiveness of the learned hidden concept vector, we prompt the LLM only with the query, i.e., the zero-shot case, and set the attention head values at some layer as the learned hidden concept vector.

In the bandit experiment, we ask GPT-4 for the procedures to find the arm with the highest reward. In each bandit instance, there is an informative arm, whose reward is exactly the index of the best arm. When prompting models, we provide the historical data of several bandit instances that share the same informative arm and ask models to specify how we should play in a similar bandit instance. A prompt sample with two examples is provided as follows.

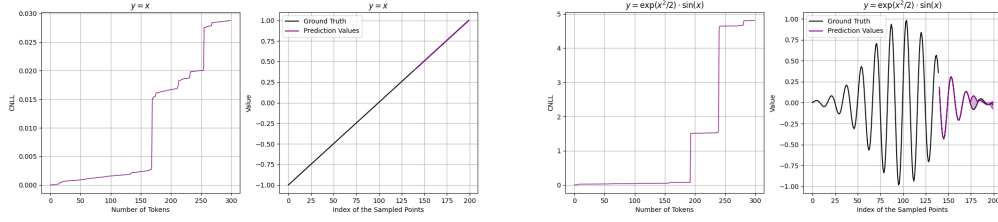


Figure 6: Cumulative negative log-likelihood and the prediction values for $y = x$. Figure 7: Cumulative negative log-likelihood and the prediction values for $y = \exp(x^2/2) \cdot \sin(x)$.

Your goal is to find the index of the arm with the highest reward, but the pulled arm may not have the highest reward. I will provide you with the past pull history on other bandits. The format of the history data on each bandit is [arm, reward]. Different pulls are separated by a comma. For example, [5,6] indicates that arm 5 will give us a reward of 6 by pulling it.

You should learn from history and tell me which arm to pull in the current bandit to find the arm with the highest reward. The history data is as follows.

Bandit:
[77, 871], [95, 613], [75, 655], [17, 449], [31, 13], [13, 1028]

Bandit:
[40, 698], [44, 88], [80, 147], [94, 265], [24, 1063], [31, 24]

Different bandits can have different rewards for each arm, but all bandits share a common pattern. Some arm may directly tells you the arm with the highest reward, even itself does not have the highest reward. Now I am playing a new bandit. This bandit will have different rewards than the bandits in history, but they share the same pattern. Tell me which arm to pull to find the arm with the highest reward. Tell me the final answer that only contains the index of the arm in a single line without any additional text.

In the above prompt, the arm 31 always returns the index of the best arm. Thus, we expect LLMs to tell us to pull arm 31 to find the best arm. The number of arms in each instance is 100, and each example only provides information about six arms in each instance. We repeat the prompt with different data ten times to plot Figure 2.

For the linear regression task, the model is trained with the loss

$$L(f) = \frac{1}{T} \sum_{t=1}^T (y_t - f(\text{prompt}_t))^2,$$

where $\text{prompt}_t = (x_1, y_1, \dots, x_{t-1}, y_{t-1}, x_t)$, $y_t = w^T x_t$, $\{x_t\}_{t=1}^T$ and w are i.i.d. samples of Gaussian distribution (Garg et al., 2022). The model is designed based on GPT-2, and we add linear layers as the first and last layers to accommodate it for the value prediction task. In the testing phase, we sample w^* and $\{x_t\}_{t=1}^T$ from the Gaussian distribution and let the model predict the response value of a query x_{t+1} given the previous examples $\{x_i, y_i\}_{i=1}^t$. We reuse the code and model in Garg et al. (2022) for the experiments. The error bar in Figure 3 is derived from 90% confidence intervals over 1000 bootstrap trials.

In the function value prediction task, we adopt the method in Gruver et al. (2023) to transfer the real number into tokens. We separate the digits with spaces and add commas ',' between the function values at different times. We calculate the negative likelihood of text-DaVinci-003 by extracting the probability value in the last layer of it. We note that the negative likelihood in Figure 4 takes every token into account, including the separating spaces between the digits.

In the experiment about the ratio between attn_i and attn , we set W_Q , W_K and W_V in attn all as the identity matrix. The entries in the K of (4.8) are i.i.d. samples of the normal distribution, and the i -th entry of V is calculated as the inner product between a Gaussian vector and the i -th column. The Gaussian vector is sampled from $\mathcal{N}(0, I)$. The error bar in Figure 3 is derived from 75th and 25th percentiles over 500 trials.

E FIGURE FOR PRETRAINING AND ICL

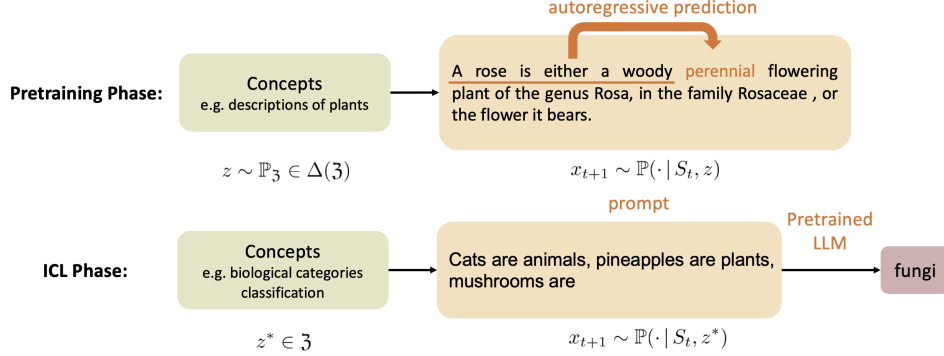


Figure 8: To form the pretraining dataset, a hidden concept z is first sampled according to \mathbb{P}_3 , and a document is generated from the concept. Taking the token sequence S_t up to position $t \in [T]$ as the input, the LLM is pretrained to maximize the next token x_{t+1} . During the ICL phase, the pretrained LLM is prompted with several examples to predict the response of the query.

F PROOFS FOR SECTION 4.1

F.1 INTRODUCTION OF CONDITIONAL MEAN EMBEDDING

Let \mathcal{H}_k and \mathcal{H}_v be the two RKHSs over the spaces Ω and \mathfrak{V} with the kernels \mathfrak{K} and \mathfrak{L} , respectively. We denote by $\phi : \Omega \rightarrow \ell_2$ and $\varphi : \mathfrak{V} \rightarrow \ell_2$ the feature mappings associated with \mathcal{H}_k and \mathcal{H}_v , respectively. Here ℓ_2 is the space of the square-integrable function class. Then it holds for any $k, k' \in \Omega$ and $v, v' \in \mathfrak{V}$ that

$$\phi(k)^\top \phi(k') = \mathfrak{K}(k, k'), \quad \varphi(v)^\top \varphi(v') = \mathfrak{L}(v, v'). \quad (\text{F.1})$$

Let $\mathbb{P}_{\mathcal{K}, \mathcal{V}}$ be the joint distribution of the two random variables \mathcal{K} and \mathcal{V} taking values in Ω and \mathfrak{V} , respectively. Then the conditional mean embedding $\text{CME}(q, \mathbb{P}_{\mathcal{K}, \mathcal{V}}) \in \mathcal{H}_v$ of the conditional distribution $\mathbb{P}_{\mathcal{V}|\mathcal{K}}$ is defined as

$$\text{CME}(q, \mathbb{P}_{\mathcal{K}, \mathcal{V}}) = \mathbb{E}[\mathfrak{L}(\mathcal{V}, \cdot) | \mathcal{K} = q].$$

The conditional mean embedding operator $C_{\mathcal{V}|\mathcal{K}} : \mathcal{H}_k \rightarrow \mathcal{H}_v$ is a linear operator such that

$$C_{\mathcal{V}|\mathcal{K}} \mathfrak{K}(q, \cdot) = \text{CME}(q, \mathbb{P}_{\mathcal{K}, \mathcal{V}}),$$

for any $q \in \Omega$. We define the (uncentered) covariance operator $C_{\mathcal{K}\mathcal{K}} : \mathcal{H}_k \rightarrow \mathcal{H}_k$ and the (uncentered) cross-covariance operator $C_{\mathcal{V}\mathcal{K}} : \mathcal{H}_k \rightarrow \mathcal{H}_v$ as follows,

$$C_{\mathcal{K}\mathcal{K}} = \mathbb{E}[\mathfrak{K}(\mathcal{K}, \cdot) \otimes \mathfrak{K}(\mathcal{K}, \cdot)], \quad C_{\mathcal{V}\mathcal{K}} = \mathbb{E}[\mathfrak{L}(\mathcal{V}, \cdot) \otimes \mathfrak{K}(\mathcal{K}, \cdot)].$$

Here \otimes is the tensor product. Song et al. (2009) shows that $C_{\mathcal{V}|\mathcal{K}} = C_{\mathcal{V}\mathcal{K}} C_{\mathcal{K}\mathcal{K}}^{-1}$. Thus, we have that

$$\text{CME}(c, \mathbb{P}_{\mathcal{K}, \mathcal{V}}) = C_{\mathcal{V}\mathcal{K}} C_{\mathcal{K}\mathcal{K}}^{-1} \mathfrak{K}(c, \cdot). \quad (\text{F.2})$$

For i.i.d. samples $\{(k^\ell, v^\ell)\}_{\ell \in [L]}$ of $\mathbb{P}_{\mathcal{K}, \mathcal{V}}$, $\|\cdot\|_{\text{HS}}$ denotes the Hilbert-Schmidt norm, we write $\phi(K) = (\phi(k^1), \dots, \phi(k^L))^\top \in \mathbb{R}^{L \times d_\phi}$ and $\varphi(V) = (\varphi(v^1), \dots, \varphi(v^L))^\top \in \mathbb{R}^{L \times d_\varphi}$. Then the empirical covariance operator $\hat{C}_{\mathcal{K}\mathcal{K}}$ and empirical cross-covariance operator $\hat{C}_{\mathcal{V}\mathcal{K}}$ are defined as

$$\begin{aligned} \hat{C}_{\mathcal{K}\mathcal{K}} &= L^{-1} \sum_{\ell=1}^L \phi(k^\ell) \phi(k^\ell)^\top = L^{-1} \phi(K)^\top \phi(K) \in \mathbb{R}^{d_\phi \times d_\phi} \\ \hat{C}_{\mathcal{V}\mathcal{K}} &= L^{-1} \sum_{\ell=1}^L \varphi(v^\ell) \phi(k^\ell)^\top = L^{-1} \varphi(V) \phi(K)^\top \in \mathbb{R}^{d_\varphi \times d_\phi}. \end{aligned} \quad (\text{F.3})$$

The empirical version of the conditional operator is

$$\hat{C}_{\mathcal{V}|\mathcal{K}}^\lambda = \varphi(V)^\top \phi(K) (\phi(K)^\top \phi(K) + \lambda \mathcal{I})^{-1} = \hat{C}_{\mathcal{V}\mathcal{K}} (\hat{C}_{\mathcal{K}\mathcal{K}} + L^{-1} \lambda \mathcal{I})^{-1} \in \mathbb{R}^{d_\varphi \times d_\phi}.$$

F.2 PROOF OF PROPOSITION 4.1

Proof. By (4.1), we have that

$$\begin{aligned}
\mathbb{P}(r_{t+1} | \text{prompt}_t) &= \int \mathbb{P}(r_{t+1} | h_{t+1}, \text{prompt}_t) \mathbb{P}(h_{t+1} | \text{prompt}_t) dh_{t+1} \\
&= \int \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, h_{t+1}) \mathbb{P}(h_{t+1} | S_t) dh_{t+1} \\
&= \int \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, h_{t+1}) \mathbb{P}(h_{t+1} | S_t, z) \mathbb{P}(z | S_t) dh_{t+1} dz \\
&= \int \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, h_{t+1}, S_t, z) \mathbb{P}(h_{t+1} | S_t, z) dh_{t+1} \mathbb{P}(z | S_t) dz \\
&= \int \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, S_t, z) \mathbb{P}(z | S_t) dz,
\end{aligned} \tag{F.4}$$

$$\tag{F.5}$$

where the first inequality results from the Bayes rule, the second equality results from the fact that r_{t+1} is conditionally independent with the previous history given h_{t+1}, \tilde{c}_{t+1} and the fact that h_{t+1} only parameterizes the transition kernel of r_{t+1} given c_{t+1} in (4.1), the fourth equality results from the fact that r_{t+1} is conditionally independent with the other variables given h_{t+1}, \tilde{c}_{t+1} , and the last equality results from the Bayes' rule.

□

F.3 PROOF OF COROLLARY 4.2

Proof. Note that

$$\mathbb{P}(z | S_t) = \frac{\mathbb{P}(S_t | z) \mathbb{P}_{\mathcal{Z}}(z)}{\int \mathbb{P}(S_t | z') \mathbb{P}_{\mathcal{Z}}(z') dz'} = \frac{\prod_{i=1}^t \mathbb{P}(r_i | z, S_t, c_i) \mathbb{P}_{\mathcal{Z}}(z)}{\int \prod_{i=1}^t \mathbb{P}(r_i | z', S_{i-1}, c_i) \mathbb{P}_{\mathcal{Z}}(z') dz'},$$

where the second equality results from the fact that the hidden variable z only parameterizes the **conditional probability** of r_t given c_t , c_t and z are independent. Then, by Bayesian model averaging, we have the following density estimation,

$$\begin{aligned}
\mathbb{P}(r_{t+1} | S_t, c_{t+1}) &= \int \mathbb{P}(r_{t+1} | z, S_t, c_{t+1}) \mathbb{P}(z | S_t) dz \\
&= \frac{\int \prod_{i=1}^{t+1} \mathbb{P}(r_i | z, S_{i-1}, c_i) \mathbb{P}_{\mathcal{Z}}(z) dz}{\int \prod_{i=1}^t \mathbb{P}(r_i | z', S_{i-1}, c_i) \mathbb{P}_{\mathcal{Z}}(z') dz'}.
\end{aligned}$$

Thus, it holds that

$$\begin{aligned}
-\sum_{t=0}^T \log \mathbb{P}(r_{t+1} | c_{t+1}, S_t) &= -\sum_{i=1}^t \left(\log \int \prod_{i=1}^{t+1} \mathbb{P}(r_i | z, S_{i-1}, c_i) \mathbb{P}_{\mathcal{Z}}(z) dz - \log \int \prod_{i=1}^t \mathbb{P}(r_i | z, S_{i-1}, c_i) \mathbb{P}_{\mathcal{Z}}(z) dz \right) \\
&= -\log \int \prod_{t=0}^T \mathbb{P}(r_t | z, S_{t-1}, c_t) \mathbb{P}_{\mathcal{Z}}(z) dz \\
&= \inf_q \mathbb{E}_{z \sim q} \left[-\sum_{i=1}^{T+1} \log \mathbb{P}(r_i | z, S_{i-1}, c_i) \right] + \mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{\mathbb{P}_{\mathcal{Z}}(z)} \right],
\end{aligned}$$

where the second equality results from the fact that $\mathbb{P}(r_{t+1} | c_{t+1}, S_t) = \frac{\int \mathbb{P}(r_1 | c_1, z) \mathbb{P}_{\mathcal{Z}}(z) dz}{1}$, and the last equality results from the standard Lagrangian arguments.

We consider q to be in the class of all Dirac measures. Then, we have that

$$-\frac{1}{T} \sum_{t=1}^T \log \mathbb{P}(r_t | c_t, S_{t-1}) \leq \frac{1}{T} \inf_z \left(-\sum_{t=1}^T \log \mathbb{P}(r_t | z, S_{t-1}, c_t) - \log \mathbb{P}_{\mathcal{Z}}(z) \right).$$

Thus, the statistical convergence rate of the Bayesian posterior averaging is $\mathcal{O}(1/T)$.

□

F.4 PROOF OF PROPOSITION 4.3

Proof. The proof of Proposition 4.3 mainly involves two steps

- Build the relationship between attn_\dagger and conditional mean embedding.
- Build the relationship between the attn and conditional mean embedding.

Step 1: Build the relationship between attn_\dagger and conditional mean embedding.

In the following, we adopt \mathcal{H}_k and \mathcal{H}_v to denote the RKHSs for the key and the value with the kernel functions \mathfrak{K} and \mathfrak{L} , respectively. Also, we use $\|\cdot\|$ to denote the norm of RKHS for an element in the corresponding RKHS and the operator norm of the operators that transform elements between RKHSs. For the value space, we adopt the Euclidean kernel $\mathfrak{L}(v, v') = v^\top v'$, and the feature mapping φ is the identity mapping. Recall the definition of the empirical covariance operator and the empirical cross-covariance operator in Appendix F.1. For keys and values, we correspondingly define them as

$$\hat{C}_{KK} = L^{-1}\phi(K)^\top\phi(K), \quad \hat{C}_{VK} = L^{-1}\varphi(V)^\top\phi(K), \quad \hat{C}_{VV} = L^{-1}\varphi(V)^\top\varphi(V),$$

where $\phi(K) = (\phi(k^1), \dots, \phi(k^L))^\top \in \mathbb{R}^{L \times d_\phi}$ and $\varphi(V) = (\varphi(v^1), \dots, \varphi(v^L))^\top \in \mathbb{R}^{L \times d_\varphi}$. By the definition of the newly defined attention in Section 4.1, we have that

$$\text{attn}_\dagger(q, K, V) = \hat{C}_{VK}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\phi(q),$$

which implies that attn_\dagger recovers the empirical conditional mean embedding. By (F.2), it holds that

$$\begin{aligned} & \|\text{attn}_\dagger(q, K, V) - \text{CME}(q, \mathbb{P}_{K,V})\| \\ & \leq \underbrace{\|\hat{C}_{VK}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\phi(q) - C_{VK}(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\phi(q)\|}_{(i)} \\ & \quad + \underbrace{\|C_{VK}(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\mathfrak{K}(q, \cdot) - C_{VK}C_{KK}^{-1}\mathfrak{K}(q, \cdot)\|}_{(ii)}. \end{aligned} \quad (\text{F.6})$$

Upper bounding term (i) of (F.6). Following the proof from Song et al. (2009), we only need to upper bound $\|\hat{C}_{VK}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1} - C_{VK}(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\|$. It holds that

$$\begin{aligned} & \|\hat{C}_{VK}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1} - C_{VK}(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\| \\ & \leq \|\hat{C}_{VK}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1}(\hat{C}_{KK} - C_{KK})(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\| + \|(\hat{C}_{VK} - C_{VK})(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\|. \end{aligned} \quad (\text{F.7})$$

Considering the first term on the right-hand side of (F.7), we have the operator decomposition $\hat{C}_{VK} = \hat{C}_{VV}^{1/2}\mathcal{W}\hat{C}_{KK}^{1/2}$ for \mathcal{W} such that $\|\mathcal{W}\| \leq 1$. This decomposition implies that

$$\begin{aligned} & \|\hat{C}_{VK}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1}(\hat{C}_{KK} - C_{KK})(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\| \\ & \leq \|\hat{C}_{VV}\|^{1/2} \cdot \|\hat{C}_{KK}^{1/2}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1/2}\| \cdot \|(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1/2}\| \cdot \|(\hat{C}_{KK} - C_{KK})(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\| \\ & \leq (L^{-1}\lambda)^{-1/2} \cdot \|(\hat{C}_{KK} - C_{KK})(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\|, \end{aligned} \quad (\text{F.8})$$

where the last inequality follows from the fact that

$$\|\hat{C}_{VV}\|^2 = L^{-1} \sum_{\ell=1}^L \|v^\ell\|_2^2 \leq 1, \quad \hat{C}_{KK}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1} \leq \mathcal{I}, \quad (\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1} \leq (L^{-1}\lambda)^{-1}\mathcal{I}.$$

Combining (F.8) and (F.7), we have

$$\begin{aligned} & \|\hat{C}_{VK}(\hat{C}_{KK} + L^{-1}\lambda\mathcal{I})^{-1} - C_{VK}(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\| \\ & \leq (L^{-1}\lambda)^{-1/2} \cdot \|(\hat{C}_{KK} - C_{KK})(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\| + \|(\hat{C}_{VK} - C_{VK})(C_{KK} + L^{-1}\lambda\mathcal{I})^{-1}\|. \end{aligned} \quad (\text{F.9})$$

In the following, we will upper bound the second term on the right-hand side of (F.9) with Lemma J.1. For this purpose, we define $\xi : \mathbb{R}^{d_v} \times \mathbb{R}^d \rightarrow \mathcal{H}_k \otimes \mathcal{H}_v$ as follows,

$$\xi(k, v) = \varphi(v)\phi(k)^\top (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}.$$

Since the operator norm of $(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}$ is upper bounded by $(L^{-1}\lambda)^{-1}$, we have that

$$\|\xi(k, v)\| = \|(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\| \cdot \|\varphi(v)\| \cdot \|\phi(k)\| \leq C \cdot (L^{-1}\lambda)^{-1},$$

where $C > 0$ is an absolute constant. Additionally, we can bound the expectation of the squared norm of $\xi(k, v)$ as

$$\begin{aligned} \mathbb{E}[\|\xi(k, v)\|^2] &= \mathbb{E}[\|\phi(k)^\top (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\|^2 \cdot \|\varphi(v)\|^2] \\ &\leq \mathbb{E}[\|(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\phi(k)\|^2] \\ &\leq (L^{-1}\lambda)^{-1} \cdot \mathbb{E}[\langle (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\phi(k), \phi(k) \rangle]. \end{aligned}$$

Using the definition of the trace operator, we have

$$\begin{aligned} \mathbb{E}[\|\xi(k, v)\|^2] &\leq \mathbb{E}[\text{tr}((C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-2}\phi(k)\phi(k)^\top)] \\ &\leq (L^{-1}\lambda)^{-1} \cdot \text{tr}((C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}}) \\ &= (L^{-1}\lambda)^{-1} \cdot \Gamma(L^{-1}\lambda). \end{aligned}$$

Here $\Gamma(L^{-1}\lambda)$ is the effective dimension of $C_{\mathcal{K}\mathcal{K}}$ in Caponnetto and De Vito (2007), which is defined as follows,

$$\Gamma(L^{-1}\lambda) = \text{tr}((C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}}).$$

We apply Lemma J.1 with $B = C(L^{-1}\lambda)^{-1}$ and $\sigma^2 = (L^{-1}\lambda)^{-1} \cdot \Gamma(L^{-1}\lambda)$, then we have that with probability at least $1 - \delta$, the following holds

$$\|\widehat{C}_{\mathcal{V}\mathcal{K}}(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1} - C_{\mathcal{V}\mathcal{K}}(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\| \leq C \cdot \left(\frac{2}{\lambda} + \sqrt{\frac{\Gamma(L^{-1}\lambda)}{\lambda}}\right) \log \frac{2}{\delta}, \quad (\text{F.10})$$

where $C > 0$ is an absolute constant. Similarly, we can prove that with probability at least $1 - \delta$, the following holds

$$\|\widehat{C}_{\mathcal{K}\mathcal{K}}(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1} - C_{\mathcal{K}\mathcal{K}}(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\| \leq C' \cdot \left(\frac{2}{\lambda} + \sqrt{\frac{\Gamma(L^{-1}\lambda)}{\lambda}}\right) \log \frac{2}{\delta}. \quad (\text{F.11})$$

Here $C' > 0$ is an absolute constant. Combining (F.9), (F.10), and (F.11), we have with probability at least $1 - \delta$ that

$$\begin{aligned} &\|\widehat{C}_{\mathcal{V}\mathcal{K}}(\widehat{C}_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1} - C_{\mathcal{V}\mathcal{K}}(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\| \\ &\leq C'' \cdot \sqrt{\frac{L}{\lambda}} \cdot \left(\frac{2}{\lambda} + \sqrt{\frac{\Gamma(L^{-1}\lambda)}{\lambda}}\right) \log \frac{2}{\delta}. \end{aligned} \quad (\text{F.12})$$

Upper bounding term (ii) of (F.6). We follow the procedures in the proof from Fukumizu (2015). For any $g \in \mathcal{H}_k$, we have that

$$\begin{aligned} \langle C_{\mathcal{V}\mathcal{K}}(g), C_{\mathcal{V}\mathcal{K}}(g) \rangle &= \mathbb{E}[\mathfrak{L}(\mathcal{V}, \bar{\mathcal{V}})g(\mathcal{K})g(\bar{\mathcal{K}})] \\ &= \left\langle (C_{\mathcal{K}\mathcal{K}} \otimes C_{\mathcal{K}\mathcal{K}}) \mathbb{E}[\mathfrak{L}(\mathcal{V}, \bar{\mathcal{V}}) \mid \mathcal{K} = \cdot, \bar{\mathcal{K}} = \ddagger], g \otimes g \right\rangle. \end{aligned}$$

Similarly, for any $q \in \mathbb{R}^{d_v}$ and any $g \in \mathcal{H}_k$, we have that

$$\begin{aligned} \left\langle C_{\mathcal{V}\mathcal{K}}, \mathbb{E}[\mathfrak{L}(\mathcal{V}, \cdot) \mid \mathcal{K} = q] \right\rangle &= \left\langle \mathbb{E}[\mathfrak{L}(\mathcal{V}, \bar{\mathcal{V}}) \mid \mathcal{K} = q, \bar{\mathcal{K}} = \ddagger], C_{\mathcal{K}\mathcal{K}}g \right\rangle \\ &= \left\langle (\mathcal{I} \otimes C_{\mathcal{K}\mathcal{K}}) \mathbb{E}[\mathfrak{L}(\mathcal{V}, \bar{\mathcal{V}}) \mid \mathcal{K} = \cdot, \bar{\mathcal{K}} = \ddagger], \mathfrak{L}(\cdot, q) \otimes g \right\rangle. \end{aligned}$$

Taking $g = (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\mathfrak{K}(q, \cdot)$, we have that

$$\begin{aligned} & \|C_{\mathcal{V}\mathcal{K}}(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\mathfrak{K}(q, \cdot) - C_{\mathcal{V}\mathcal{K}}C_{\mathcal{K}\mathcal{K}}^{-1}\mathfrak{K}(q, \cdot)\|^2 \\ &= \left\langle \left((C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}} \otimes (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}} - \mathcal{I} \otimes (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}} \right. \right. \\ & \quad \left. \left. (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}} \otimes \mathcal{I} + \mathcal{I} \otimes \mathcal{I} \right) \mathbb{E}[\mathfrak{L}(\mathcal{V}, \bar{\mathcal{V}}) \mid \mathcal{K} = \cdot, \bar{\mathcal{K}} = \dagger], \mathfrak{K}(q, \cdot) \otimes \mathfrak{K}(q, \dagger) \right\rangle. \end{aligned}$$

We note that $\mathbb{E}[\mathfrak{L}(v, \bar{v}) \mid k = \cdot, \bar{k} = \dagger] \in \mathcal{H}_k \otimes \mathcal{H}_k$ is in the range spanned by $C_{\mathcal{K}\mathcal{K}} \otimes C_{\mathcal{K}\mathcal{K}}$. Thus, we can define $\tilde{\mathcal{C}} \in \mathcal{H}_k \times \mathcal{H}_k$ such that $(C_{\mathcal{K}\mathcal{K}} \otimes C_{\mathcal{K}\mathcal{K}})\tilde{\mathcal{C}} = \mathbb{E}[\mathfrak{L}(v, \bar{v}) \mid k = \cdot, \bar{k} = \dagger]$. Let $\{\lambda_i\}_{i=1}^\infty$ and $\{\varphi_i\}_{i=1}^\infty$ be the eigenvalues and eigenvectors of $C_{\mathcal{K}\mathcal{K}}$, respectively. We then have that

$$\begin{aligned} & \|C_{\mathcal{V}\mathcal{K}}(C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}\mathfrak{K}(q, \cdot) - C_{\mathcal{V}\mathcal{K}}C_{\mathcal{K}\mathcal{K}}^{-1}\mathfrak{K}(q, \cdot)\|^4 \\ & \leq \left\| \left((C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}} \otimes (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}} - \mathcal{I} \otimes (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}} \right. \right. \\ & \quad \left. \left. (C_{\mathcal{K}\mathcal{K}} + L^{-1}\lambda\mathcal{I})^{-1}C_{\mathcal{K}\mathcal{K}} \otimes \mathcal{I} + \mathcal{I} \otimes \mathcal{I} \right) \mathbb{E}[\mathfrak{L}(\mathcal{V}, \bar{\mathcal{V}}) \mid \mathcal{K} = \cdot, \bar{\mathcal{K}} = \dagger] \right\|^2 \\ & = \sum_{i,j} \left(\frac{\lambda_i \lambda_j (L^{-1}\lambda)^2}{(\lambda_i + L^{-1}\lambda)(\lambda_j + L^{-1}\lambda)} \right)^2 \cdot \langle \varphi_i \otimes \varphi_j, \tilde{\mathcal{C}} \rangle^2 \\ & \leq (L^{-1}\lambda)^4 \cdot \|\tilde{\mathcal{C}}\|^2. \end{aligned}$$

Thus, we have

$$\|C_{\mathcal{V}\mathcal{K}}(C_{\mathcal{K}\mathcal{K}} + \lambda\mathcal{I})^{-1}\mathfrak{K}(q, \cdot) - C_{\mathcal{V}\mathcal{K}}C_{\mathcal{K}\mathcal{K}}^{-1}\mathfrak{K}(q, \cdot)\|_2 \leq C \cdot \lambda L^{-1}, \quad (\text{F.13})$$

where $C > 0$ is an absolute constant.

Combining (F.6), (F.12), and (F.13), we have with probability at least $1 - \delta$, the following holds

$$\|\text{attn}_\dagger(q, K, V) - \text{CME}(q, \mathbb{P}_{\mathcal{K}, \mathcal{V}})\| \leq \mathcal{O}\left(\sqrt{\frac{L}{\lambda}} \cdot \left(\frac{2}{\lambda} + \sqrt{\frac{\Gamma(L^{-1}\lambda)}{\lambda}}\right) \log \frac{1}{\delta} + \lambda L^{-1}\right). \quad (\text{F.14})$$

Since \mathfrak{K} is Gaussian RBF kernel, we have that $\Gamma(L^{-1}\lambda) = \mathcal{O}(L/\lambda)$.

Step 2: Build the relationship between the attn and conditional mean embedding.

We achieve our goal in two sub-steps. In the first step, we prove that there exists a constant $C > 0$ such that

$$\text{attn}_{\text{SM}}(q, K, V) = C \int_{\mathbb{S}^{d-1}} v \hat{\mathbb{P}}_{\mathcal{V}|\mathcal{K}}^{\mathfrak{K}}(v|q) dv, \quad (\text{F.15})$$

where \mathbb{S}^{d-1} is the $(d-1)$ -dimensional unit sphere. Here $\hat{\mathbb{P}}_{\mathcal{V}|\mathcal{K}}^{\mathfrak{K}}$ is the kernel conditional density estimation of $\mathbb{P}_{\mathcal{V}|\mathcal{K}}$ defined as follows,

$$\hat{\mathbb{P}}_{\mathcal{V}|\mathcal{K}}^{\mathfrak{K}}(v|q) = \frac{\sum_{\ell=1}^L \mathfrak{K}(k^\ell, q) \cdot \mathfrak{K}(v^\ell, v)}{\sum_{\ell=1}^L \mathfrak{K}(k^\ell, q)},$$

where $\iota = 1/\int_{\mathbb{S}^{d-1}} \mathfrak{K}(k, q) dq$ is a normalization constant. Note that ι does not depend on the value of k by symmetry. We transform the right-hand side of this equality as

$$\begin{aligned} \int v \hat{\mathbb{P}}_{\mathcal{V}|\mathcal{K}}^{\mathfrak{K}}(v|q) dv &= \iota \cdot \int_{\mathbb{S}^{d-1}} v \cdot \frac{\sum_{\ell=1}^L \mathfrak{K}(k^\ell, q) \cdot \mathfrak{K}(v^\ell, v)}{\sum_{\ell=1}^L \mathfrak{K}(k^\ell, q)} dv \\ &= \frac{\iota \cdot \sum_{\ell=1}^L \mathfrak{K}(k^\ell, q) \cdot \int_{\mathbb{S}^{d-1}} v \cdot \mathfrak{K}(v^\ell, v) dv}{\sum_{\ell=1}^L \mathfrak{K}(k^\ell, q)}. \end{aligned} \quad (\text{F.16})$$

Thus, it suffices to calculate the integration term $\int_{\mathbb{S}^{d-1}} v \cdot \mathfrak{K}(v^\ell, v) dv$. To this end, we have the following lemma.

Proposition F.1. Let $\mathfrak{K}(a, b) = \exp(a^\top b / \gamma)$ be the exponential kernel with a fixed $\gamma > 0$. It holds for any $b \in \mathbb{S}^{d-1}$ that

$$\int_{\mathbb{S}^{d-1}} a \cdot \mathfrak{K}(a, b) da = C_1 \cdot b,$$

where $C_1 > 0$ is an absolute constant.

Proof. See Section I.1 for a detailed proof. \square

Thus, it holds for the right-hand side of (F.16) that

$$\iota \cdot C_1 \cdot \frac{\sum_{\ell=1}^L \mathfrak{K}(k^\ell, q) \cdot v^\ell}{\sum_{\ell=1}^L \mathfrak{K}(k^\ell, q)} = \iota \cdot C_1 \cdot V^\top \text{softmax}(Kq/\gamma) = \iota \cdot C_1 \cdot \text{attn}_{\text{SM}}(q, K, V),$$

where the first equality follows from the definition of the softmax function and the second equality follows from the definition of the softmax attention.

The second step is to relate the right-hand side of (F.15) to conditional mean embedding. In fact, under the condition that $\widehat{\mathbb{P}}_{\mathcal{V}|\mathcal{K}}^{\mathfrak{K}}(v|q) \rightarrow \mathbb{P}(v|q)$ uniformly for any $q \in \mathbb{S}^{d_p-1}$ as $L \rightarrow \infty$, we have

$$\int v \widehat{\mathbb{P}}_{\mathcal{V}|\mathcal{K}}^{\mathfrak{K}}(v|q) dv \rightarrow \mathbb{E}[\mathcal{V}|\mathcal{K} = q] \quad \text{as } L \rightarrow \infty.$$

Thus, we have that

$$\text{attn}_{\text{SM}}(q, K, V) \rightarrow C \cdot \mathbb{E}[\mathcal{V}|\mathcal{K} = q] \quad \text{as } L \rightarrow \infty \quad (\text{F.17})$$

for some constant $C > 0$. Combining (F.17) and (F.14) and choosing $\lambda = L^{3/4}$, we complete the proof of Proposition 4.3. \square

G APPENDIX FOR SECTION 5

G.1 SUPPLEMENTAL DEFINITIONS FOR MARKOV CHAINS

We follow the notations in Paulin (2015). Let Ω be a Polish space. The transition kernel for a time-homogeneous Markov chain $\{X_i\}_{i=1}^\infty$ supported on Ω is a probability distribution $\mathbb{P}(x, dy)$ for every $x \in \Omega$. Given $X_1 = x_1, \dots, X_{t-1} = x_{t-1}$, the conditional distribution of X_t equals $\mathbb{P}(x_{t-1}, dy)$. A distribution π is said to be a stationary distribution of this Markov chain if $\int_{x \in \Omega} \mathbb{P}(x, dy) \pi(dx) = \pi(dy)$. We adopt $\mathbb{P}^t(x, \cdot)$ to denote the distribution of X_t conditioned on $X_1 = x$. The *mixing time* of the chain is defined by

$$d(t) = \sup_{x \in \Omega} \text{TV}(P^t(x, \cdot), \pi), \quad t_{\text{mix}}(\varepsilon) = \min\{t \mid d(t) \leq \varepsilon\}, \quad t_{\text{mix}} = t_{\text{mix}}(1/4).$$

G.2 PROOF OF THEOREM 5.3

Proof of Theorem 5.3. Our proof mainly involves three steps.

- Error decomposition with the PAC-Bayes framework.
- Control each term in the error decomposition.
- Conclude the proof.

Step 1: Error decomposition with the PAC-Bayes framework.

For ease of notation, we temporarily write T_p and N_p as T and N , respectively. Recall that the pretraining dataset is $\mathcal{D} = \{(S_t^n, x_{t+1}^n)\}_{n,t=1}^{N,T}$, which consists of N trajectories (essays), and each essay have $T + 1$ words. Given S_t^n , the next word is generated as $x_{t+1}^n \sim \mathbb{P}(\cdot | S_t^n)$, and $S_{t+1}^n =$

(S_t^n, x_{t+1}^n) . Here, we construct a ghost sample $\tilde{\mathcal{D}} = \{(\tilde{S}_t^n, \tilde{x}_{t+1}^n)\}_{n,t=1}^{N,T}$ as $\tilde{S}_t^n = S_t^n$ and $\tilde{x}_{t+1}^n \sim \mathbb{P}(\cdot | \tilde{S}_t^n)$ independently from \mathcal{D} . We define function $g(\theta) = L(\theta, \mathcal{D}) - \log \mathbb{E}_{\tilde{\mathcal{D}}}[\exp(L(\theta, \tilde{\mathcal{D}})) | \mathcal{D}]$, where

$$L(\theta, \tilde{\mathcal{D}}) = -\frac{1}{4} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}(\tilde{x}_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta}(\tilde{x}_{t+1}^n | S_t^n)}.$$

For distributions $Q, P \in \Delta(\Theta)$, where P can potentially depends on \mathcal{D} , Lemma J.3 shows that

$$\mathbb{E}_P[g(\theta)] \leq \text{KL}(P||Q) + \log \mathbb{E}_Q[\exp(g(\theta))].$$

Substituting the definition of $g(\theta)$ and taking expectation with respect to the distribution of \mathcal{D} on the both sides of the inequality, we can derive that

$$\mathbb{E}_{\mathcal{D}} \left[\exp \left\{ \mathbb{E}_P \left[L(\theta, \mathcal{D}) - \log \mathbb{E}_{\tilde{\mathcal{D}}}[\exp(L(\theta, \tilde{\mathcal{D}})) | \mathcal{D}] \right] - \text{KL}(P || Q) \right\} \right] \leq 1.$$

With Chernoff inequality, we can show that with probability at least $1 - \delta$, the following holds

$$-\mathbb{E}_{\theta \sim P} \left[\log \mathbb{E}_{\tilde{\mathcal{D}}}[\exp(L(\theta, \tilde{\mathcal{D}})) | \mathcal{D}] \right] \leq -\mathbb{E}_P[L(\theta, \mathcal{D})] + \text{KL}(P || Q) + \log \frac{1}{\delta}. \quad (\text{G.1})$$

We first cope with the left-hand side of (G.1).

$$\begin{aligned} & -\mathbb{E}_P \left[\log \mathbb{E}_{\tilde{\mathcal{D}}}[\exp(L(\theta, \tilde{\mathcal{D}})) | \mathcal{D}] \right] \\ & \geq -\frac{1}{2} \log \mathbb{E}_{\tilde{\mathcal{D}}} \left[\exp \left(-\frac{1}{2} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}(\tilde{x}_{t+1}^n | S_t^n)}{\mathbb{P}_{\hat{\theta}}(\tilde{x}_{t+1}^n | S_t^n)} \right) \middle| \mathcal{D} \right] \\ & \quad - \frac{1}{2} \mathbb{E}_{\theta \sim P} \left[\log \mathbb{E}_{\tilde{\mathcal{D}}} \left[\exp \left(-\frac{1}{2} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}_{\hat{\theta}}(\tilde{x}_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta}(\tilde{x}_{t+1}^n | S_t^n)} \right) \middle| \mathcal{D} \right] \right] \\ & = -\frac{1}{2} \sum_{n=1}^N \sum_{t=1}^T \log \mathbb{E}_{\tilde{x}_{t+1}^n \sim \mathbb{P}(\cdot | S_t^n)} \left[\exp \left(-\frac{1}{2} \log \frac{\mathbb{P}(\tilde{x}_{t+1}^n | S_t^n)}{\mathbb{P}_{\hat{\theta}}(\tilde{x}_{t+1}^n | S_t^n)} \right) \middle| \mathcal{D} \right] \\ & \quad - \frac{1}{2} \mathbb{E}_{\theta \sim P} \left[\log \mathbb{E}_{\tilde{\mathcal{D}}} \left[\exp \left(-\frac{1}{2} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}_{\hat{\theta}}(\tilde{x}_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta}(\tilde{x}_{t+1}^n | S_t^n)} \right) \middle| \mathcal{D} \right] \right] \\ & \geq \frac{1}{4} \sum_{n=1}^N \sum_{t=1}^T \text{TV}(\mathbb{P}(\cdot | S_t^n), \mathbb{P}_{\hat{\theta}}(\cdot | S_t^n))^2 - \frac{1}{2} \mathbb{E}_{\theta \sim P} \left[\log \mathbb{E}_{\tilde{\mathcal{D}}} \left[\exp \left(-\frac{1}{2} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}_{\hat{\theta}}(\tilde{x}_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta}(\tilde{x}_{t+1}^n | S_t^n)} \right) \middle| \mathcal{D} \right] \right], \end{aligned} \quad (\text{G.2})$$

where the first inequality results from the definition of $L(\theta, \mathcal{D})$ and Cauchy-Schwarz inequality, the equality results from that the transitions of \tilde{x}_{t+1}^n are independent given \mathcal{D} , and the last inequality results from Lemma J.5. The second term in the right-hand side of (G.2) can be controlled if the distribution P is chosen to concentrate around $\hat{\theta}$. This will be done in Step 2. Now we consider the right-hand side of (G.1). For any $\theta^* \in \Theta$, we can decompose it as

$$\begin{aligned} & -\mathbb{E}_P[L(\theta, \mathcal{D})] \\ & = \mathbb{E}_P \left[\frac{1}{4} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta^*}(x_{t+1}^n | S_t^n)} + \log \frac{\mathbb{P}_{\theta^*}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\hat{\theta}}(x_{t+1}^n | S_t^n)} + \log \frac{\mathbb{P}_{\hat{\theta}}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta}(x_{t+1}^n | S_t^n)} \right] \\ & \leq \frac{1}{4} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta^*}(x_{t+1}^n | S_t^n)} + \frac{1}{4} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_P \left[\log \frac{\mathbb{P}_{\hat{\theta}}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta}(x_{t+1}^n | S_t^n)} \right], \end{aligned} \quad (\text{G.3})$$

where the inequality results from the fact that $\hat{\theta}$ maximizes the likelihood. We will choose θ^* as the projection of \mathbb{P} onto $\{\mathbb{P}_{\theta} | \theta \in \Theta\}$, i.e., \mathbb{P}_{θ^*} is the best approximation of \mathbb{P} with respect to the KL divergence. Thus, the first term in the right-hand side of (G.3) is the approximation error. The

second term in the right-hand side of (G.3) can be controlled in the same way as the second term in the right-hand side of (G.2). Combining inequalities (G.1), (G.2), and (G.3), we have that

$$\begin{aligned}
& \frac{1}{4} \sum_{n=1}^N \sum_{t=1}^T \text{TV}(\mathbb{P}(\cdot | S_t^n), \mathbb{P}_{\hat{\theta}}(\cdot | S_t^n))^2 \\
& \leq \underbrace{\frac{1}{2} \mathbb{E}_{\theta \sim P} \left[\log \mathbb{E}_{\tilde{\mathcal{D}}} \left[\exp \left(-\frac{1}{2} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}_{\hat{\theta}}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta}(x_{t+1}^n | S_t^n)} \right) \middle| \mathcal{D} \right] \right]}_{\text{(I)}} + \frac{1}{4} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_P \left[\log \frac{\mathbb{P}_{\hat{\theta}}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta}(x_{t+1}^n | S_t^n)} \right] \\
& \quad + \underbrace{\frac{1}{4} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta^*}(x_{t+1}^n | S_t^n)}}_{\text{(II)}} + \underbrace{\text{KL}(P \| Q)}_{\text{(III)}} + \log \frac{1}{\delta}, \tag{G.4}
\end{aligned}$$

where term (I) is the fluctuation error induced by $\theta \sim P$, term (II) is the approximation error, and term (III) is the KL divergence between P and Q .

Step 2: Control each term in the error decomposition.

We first consider term (I). Since $\hat{\theta}$ is a deterministic function of \mathcal{D} and that $\log(\mathbb{P}_{\hat{\theta}}(x_{t+1}^n | S_t^n) / \mathbb{P}_{\theta}(x_{t+1}^n | S_t^n))$ is close to 0 if θ is close to $\hat{\theta}$, we need to design P for any $\hat{\theta} \in \Theta$ such that $\theta \sim P$ is close to $\hat{\theta}$ almost surely.

We need to quantify the fluctuation of \mathbb{P}_{θ} when θ is changing, i.e., how \mathbb{P}_{θ} is close to $\mathbb{P}_{\hat{\theta}}$ when θ is close to $\hat{\theta}$.

Proposition G.1. For any input $X \in \mathbb{R}^{L \times d}$ and $\theta, \tilde{\theta} \in \Theta$, we have that

$$\begin{aligned}
& \text{TV}(\mathbb{P}_{\theta}(\cdot | X), \mathbb{P}_{\tilde{\theta}}(\cdot | X)) \\
& \leq \frac{2}{\tau} \|A^{(D+1), \top} - \tilde{A}^{(D+1), \top}\|_{1,2} + \sum_{t=1}^D \alpha_t (\beta_t + \iota_t + \kappa_t + \rho_t),
\end{aligned}$$

where

$$\begin{aligned}
\alpha_t &= \frac{2}{\tau} B_A (1 + B_{A,1} \cdot B_{A,2}) (1 + h B_V (1 + 4 B_Q B_K))^{D-t} \\
\beta_t &= |\gamma_2^{(t)} - \tilde{\gamma}_2^{(t)}| + (1 + B_{A,1} \cdot B_{A,2}) \cdot (1 + (\|X^\top\|_{2,\infty} - 1) \mathbb{I}_{t=1}) \cdot |\gamma_1^{(t)} - \tilde{\gamma}_1^{(t)}| \\
\iota_t &= B_{A,2} \cdot \|A_1^{(t)} - \tilde{A}_1^{(t)}\|_F + B_{A,1} \cdot \|A_2^{(t)} - \tilde{A}_2^{(t)}\|_F \\
\kappa_t &= (1 + B_{A,1} \cdot B_{A,2}) \cdot (1 + (\|X^\top\|_{2,\infty} - 1) \mathbb{I}_{t=1}) \cdot \sum_{i=1}^h \|W_i^{V,(t)} - \tilde{W}_i^{V,(t)}\|_F \\
\rho_t &= 2(1 + B_{A,1} \cdot B_{A,2}) \cdot (1 + (\|X^\top\|_{2,\infty} - 1) \mathbb{I}_{t=1}) \cdot B_V \\
& \quad \cdot \sum_{i=1}^h B_K \cdot \|W_i^{Q,(t+1)} - \tilde{W}_i^{Q,(t+1)}\|_F + B_Q \cdot \|W_i^{K,(t+1)} - \tilde{W}_i^{K,(t+1)}\|_F
\end{aligned}$$

for all $t \in [D]$.

Proof of Proposition G.1. See Appendix I.3. □

Proposition G.1 implies that the difference between \mathbb{P}_{θ} and $\mathbb{P}_{\tilde{\theta}}$ can be upper-bounded by the difference between the parameters of each layer. Thus, for any $\tilde{\theta} \in \mathcal{D}$, we set the distribution P as uniform distribution on the neighborhood of parameters, and the radius of the neighborhood is set proportional to $1/NT$ shown in Figure 9.

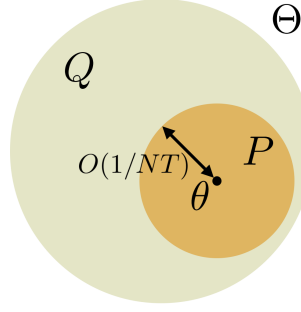


Figure 9: The distribution P in (G.5) is the uniform distribution on the neighborhood of θ with radius proportional to $1/NT$, and Q in (G.8) is the uniform distribution on Θ .

$$P = \prod_{t=1}^{D+1} \mathcal{L}_P(\theta^{(t)}) \quad (\text{G.5})$$

$$\begin{aligned} \mathcal{L}_P(\theta^{(D+1)}) &= \text{Unif}\left(\mathbb{B}(\widehat{A}^{(D+1)}, r^{(D+1)}, \|\cdot\|_{1,2})\right) \\ \mathcal{L}_P(\theta^{(t)}) &= \text{Unif}\left(\mathbb{B}(\widehat{\gamma}_1^{(t)}, r_{\gamma,1}^{(t)}, |\cdot|)\right) \cdot \text{Unif}\left(\mathbb{B}(\widehat{\gamma}_2^{(t)}, r_{\gamma,2}^{(t)}, |\cdot|)\right) \cdot \mathcal{L}_P(A^{(t)}) \cdot \mathcal{L}_P(W^{(t)}) \\ \mathcal{L}_P(A^{(t)}) &= \text{Unif}\left(\mathbb{B}(\widehat{A}_1^{(t)}, r_{A,1}^{(t)}, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(\widehat{A}_2^{(t)}, r_{A,2}^{(t)}, \|\cdot\|_F)\right) \\ \mathcal{L}_P(W^{(t)}) &= \prod_{i=1}^h \text{Unif}\left(\mathbb{B}(\widehat{W}_i^{Q,(t)}, r_Q^{(t)}, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(\widehat{W}_i^{K,(t)}, r_K^{(t)}, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(\widehat{W}_i^{V,(t)}, r_V^{(t)}, \|\cdot\|_F)\right) \end{aligned}$$

for $t \in [D]$, where Unif denotes the uniform distribution on the set, $\mathbb{B}(a, r, \|\cdot\|) = \{x \mid \|x-a\| \leq r\}$ denotes the ball centered in a with radius r , the radius is set as

$$\begin{aligned} r_{\gamma,1}^{(t)} &= R^{-1}(1 + B_{A,1} \cdot B_{A,2})^{-1} \alpha_t^{-1} / NT, & r_{\gamma,2}^{(t)} &= R^{-1} \alpha_t^{-1} / NT \\ r_{A,1}^{(t)} &= R^{-1} B_{A,2}^{-1} \alpha_t^{-1} / NT, & r_{A,2}^{(t)} &= R^{-1} B_{A,1}^{-1} \alpha_t^{-1} / NT, \\ r_V^{(t)} &= R^{-1} h^{-1} (1 + B_{A,1} \cdot B_{A,2})^{-1} \alpha_t^{-1} / NT, & r_Q^{(t)} &= R^{-1} h^{-1} (1 + B_{A,1} \cdot B_{A,2})^{-1} B_V^{-1} B_K^{-1} \alpha_t^{-1} / NT \\ r_K^{(t)} &= R^{-1} h^{-1} (1 + B_{A,1} \cdot B_{A,2})^{-1} B_V^{-1} B_Q^{-1} \alpha_t^{-1} / NT, & r^{(D+1)} &= \tau B_A^{-1} / NT. \end{aligned}$$

Under this assignment, we now bound $|\log \mathbb{P}_{\widehat{\theta}}(x|S) / \mathbb{P}_{\theta}(x|S)|$ for any $S \in \mathbb{R}^{L \times d}$ and $x \in \mathbb{R}^{d_y}$. We first note that

$$\mathbb{P}_{\widehat{\theta}}(x|S) \geq b_y = (1 + d_y \exp(B_A/\tau))^{-1} \quad (\text{G.6})$$

for any S and x , which results from the softmax layer defined below (5.1). This results from the fact that the last layer of the transformer is softmax with inverse temperature parameter τ and that

$$\left\| \frac{1}{L\tau} \mathbb{I}_L^\top X^{(D)} A^{(D+1)} \right\|_1 \leq \|A^{(D+1),\top}\|_{1,2} \leq B_A.$$

If $\text{TV}(\mathbb{P}_{\theta}(\cdot|S), \mathbb{P}_{\widehat{\theta}}(\cdot|S)) = \varepsilon \leq b_y/2$, some basic calculations show that

$$\frac{b_y}{b_y + \varepsilon} \leq \frac{\mathbb{P}_{\widehat{\theta}}(x|S)}{\mathbb{P}_{\theta}(x|S)} \leq 1 + \frac{2\varepsilon}{b_y}.$$

Thus, if we set the distribution P as the uniform distribution on the neighborhood around $\widehat{\theta}$ with radius proportional to $1/NT$, i.e., (G.5), then for $\theta \sim P$ we have that

$$\left| \log \frac{\mathbb{P}_{\widehat{\theta}}(x|S)}{\mathbb{P}_{\theta}(x|S)} \right| \leq \frac{2\varepsilon}{b_y} = \mathcal{O}\left(\frac{1}{NT}\right) \quad \text{for } P \text{ a.s.}$$

Based on this, we conclude that

$$(I) = \mathcal{O}(1). \quad (\text{G.7})$$

Next, we control term (III) in (G.4). In order to upper bound $\text{KL}(P \parallel Q)$, we need to make sure that the support of P is a subset of that of Q . Thus, we take Q as the uniform distribution on the parameter space.

$$Q = \prod_{t=1}^{D+1} \mathcal{L}_Q(\theta^{(t)}) \quad (\text{G.8})$$

$$\begin{aligned} \mathcal{L}_Q(\theta^{(D+1)}) &= \text{Unif}\left(\mathbb{B}(0, B_A, \|\cdot\|_{1,2})\right) \\ \mathcal{L}_Q(\theta^{(t)}) &= \text{Unif}\left(\mathbb{B}(1/2, 1/2, |\cdot|)\right) \cdot \text{Unif}\left(\mathbb{B}(1/2, 1/2, |\cdot|)\right) \cdot \mathcal{L}_Q(A^{(t)}) \cdot \mathcal{L}_Q(W^{(t)}) \\ \mathcal{L}_Q(A^{(t)}) &= \text{Unif}\left(\mathbb{B}(0, B_{A,1}, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(0, B_{A,2}, \|\cdot\|_F)\right) \\ \mathcal{L}_Q(W^{(t)}) &= \prod_{i=1}^h \text{Unif}\left(\mathbb{B}(0, B_Q, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(0, B_K, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(0, B_V, \|\cdot\|_F)\right). \end{aligned}$$

Then the KL divergence between P and Q is

$$\text{KL}(P \parallel Q) = \mathcal{O}\left((D^2 \cdot d \cdot (d_F + d_h + d) + d \cdot d_y) \cdot \log(1 + NT\tau^{-1}RhB_AB_{A,1}B_{A,2}B_QB_KB_V)\right). \quad (\text{G.9})$$

Finally, we control term (II) in (G.4). This term can be controlled as

$$\begin{aligned} &\frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta^*}(x_{t+1}^n | S_t^n)} \\ &= \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta^*}(x_{t+1}^n | S_t^n)} - \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_{S_t^n} \text{KL}(\mathbb{P}(\cdot | S_t^n) \parallel \mathbb{P}_{\theta^*}(\cdot | S_t^n)) \\ &\quad + \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_{S_t^n} \text{KL}(\mathbb{P}(\cdot | S_t^n) \parallel \mathbb{P}_{\theta^*}(\cdot | S_t^n)). \end{aligned}$$

The first two terms in the right-hand side of the equality is the generalization error, which can be bounded with Lemma J.4. With Assumption 5.2, we note that

$$\left| \log \frac{\mathbb{P}(x | S)}{\mathbb{P}_{\theta^*}(x | S)} \right| \leq b^* = \log \max\{c_0^{-1}, b_y^{-1}\}, \quad (\text{G.10})$$

so the function satisfies the condition in Lemma J.4 with $c_i = 2b^*$. Using the moment generating function bound in Lemma J.4 and Chernoff bound, we have that

$$\frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \log \frac{\mathbb{P}(x_{t+1}^n | S_t^n)}{\mathbb{P}_{\theta^*}(x_{t+1}^n | S_t^n)} - \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_{S_t^n} \text{KL}(\mathbb{P}(\cdot | S_t^n) \parallel \mathbb{P}_{\theta^*}(\cdot | S_t^n)) \leq \sqrt{\frac{t_{\min} b^{*,2}}{2NT}} \log \frac{1}{\delta} \quad (\text{G.11})$$

with probability at least $1 - \delta$.

Step 3: Conclude the proof.

Combining inequalities (G.4), (G.7), (G.9), and (G.11), we have that

$$\begin{aligned} &\frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \text{TV}(\mathbb{P}(\cdot | S_t^n), \mathbb{P}_{\hat{\theta}}(\cdot | S_t^n)) \\ &\leq \sqrt{\frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \text{TV}(\mathbb{P}(\cdot | S_t^n), \mathbb{P}_{\hat{\theta}}(\cdot | S_t^n))^2} \\ &= \mathcal{O}\left(\frac{t_{\min}^{1/4}}{(NT)^{1/4}} \log \frac{1}{\delta} + \frac{\sqrt{D^2 d(d_F + d_h + d) + d \cdot d_y}}{\sqrt{NT}} \cdot \log(1 + NT\bar{B})\right. \\ &\quad \left. + \inf_{\theta^* \in \Theta} \sqrt{\frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_{S_t^n} \text{KL}(\mathbb{P}(\cdot | S_t^n) \parallel \mathbb{P}_{\theta^*}(\cdot | S_t^n))}\right), \end{aligned}$$

where we take θ^* as the best approximation parameters. Finally, we will change the left-hand side of this inequality to the expectation of it. In fact, we have that

Proposition G.2. Let \mathcal{F} be the collection of functions of $f : \mathbb{R}^n \rightarrow \mathbb{R}$, and we assume that $|f| \leq b$ for any function $f \in \mathcal{F}$. For a Markov chain $X = (X_1, \dots, X_N)$, we define $f(X) = \sum_{i=1}^N f(X_i)/N$. The mixing time of this Markov chain is denoted as $t_{\text{mix}}(\varepsilon)$. Given a distribution Q on \mathcal{F} , with probability at least $1 - \delta$, we have

$$\left| \mathbb{E}_P \left[\mathbb{E}_X [f(X)] - f(X) \right] \right| \leq \sqrt{\frac{b^2 \cdot t_{\min}}{2 \log 2N}} \left[\text{KL}(P \| Q) + \log \frac{4}{\delta} \right],$$

for any distribution P on \mathcal{F} simultaneously with probability at least $1 - \delta$, where

$$t_{\min} = \inf_{0 \leq \varepsilon < 1} t_{\text{mix}}(\varepsilon) \cdot \left(\frac{2 - \varepsilon}{1 - \varepsilon} \right)^2.$$

Proof of Proposition G.2. See Appendix I.2. □

We note that Proposition G.2 is indeed an uniform convergence bound, since it holds simultaneously for all P . Thus, we can set P and Q as those in equalities (G.5) and (G.8), then we have that

$$\begin{aligned} & \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_{S_t^n} \left[\text{TV}(\mathbb{P}(\cdot | S_t^n), \mathbb{P}_{\hat{\theta}}(\cdot | S_t^n)) \right] - \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \text{TV}(\mathbb{P}(\cdot | S_t^n), \mathbb{P}_{\hat{\theta}}(\cdot | S_t^n)) \\ &= \mathcal{O} \left(\frac{\sqrt{t_{\min}}}{\sqrt{NT}} \left(\bar{D} \log(1 + NT\bar{B}) + \log \frac{1}{\delta} \right) \right). \end{aligned}$$

Thus, we have that

$$\begin{aligned} & \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_{S_t^n} \left[\text{TV}(\mathbb{P}(\cdot | S_t^n), \mathbb{P}_{\hat{\theta}}(\cdot | S_t^n)) \right] \\ &= \mathcal{O} \left(\frac{t_{\min}^{1/4}}{(NT)^{1/4}} \log \frac{1}{\delta} + \frac{\sqrt{t_{\min}}}{\sqrt{NT}} \left(\bar{D} \log(1 + NT\bar{B}) + \log \frac{1}{\delta} \right) \right. \\ & \quad \left. + \inf_{\theta^* \in \Theta} \sqrt{\frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_{S_t^n} \text{KL}(\mathbb{P}(\cdot | S_t^n) \| \mathbb{P}_{\theta^*}(\cdot | S_t^n))} \right). \end{aligned}$$

We conclude the proof of Theorem 5.3. □

G.3 FORMAL STATEMENT AND PROOF OF PROPOSITION 5.4

Denote the alphabet of the language as $\mathfrak{X} \subseteq \mathbb{R}$ ($d = 1$), then the conditional distribution \mathbb{P}^* can be viewed as a function $g^* : \mathfrak{X}^L \rightarrow \mathbb{R}^{d_y}$, where L is the maximal length of a sentence, and the output is the distribution of the next word. Since \mathcal{A} is finite, Theorem 2 in Zaheer et al. (2017) shows that there exist $\rho^* : \mathbb{R} \rightarrow \mathbb{R}^{d_y}$ and $\phi^* : \mathfrak{X} \rightarrow \mathbb{R}$ such that

$$g^*(X) = \rho^* \left(\frac{1}{L} \sum_{i=1}^L \phi^*(x_i) \right),$$

where $X = [x_1, \dots, x_L]$. The i^{th} component of ρ^* is denoted as ρ_i^* for $i \in [d_y]$. For a function f defined on Ω , the L^∞ norm of it is defined as $\|f\|_\infty = \sup_{x \in \Omega} |f(x)|$. The set of the real-valued smooth functions on it is denoted as $\mathcal{S}^\infty(\Omega, \mathbb{R})$. Then we denote the set of the smooth functions with bounded derivatives as

$$\mathcal{S}_B = \left\{ f \in \mathcal{S}^\infty([-B, B], \mathbb{R}) \mid \|f^{(n)}(x)\| \leq n! \text{ for all } n \in \mathbb{N} \right\},$$

where $f^{(n)}$ is the n^{th} -order derivative of f .

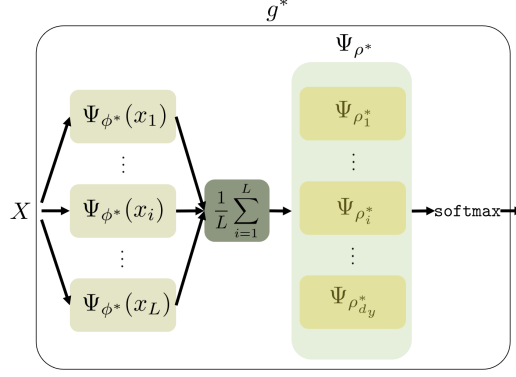


Figure 10: The construction in Proposition G.4 mainly consists of three parts: the approximation of ϕ^* , the approximation of ρ^* , and the realization of $\frac{1}{L} \sum_{i=1}^L$.

Assumption G.3. There exists $B > 0$ such that $\phi^*, \tau \log \rho_i^* \in \mathcal{S}_B$ for $i \in [d_y]$.

This assumption states that the function g^* is smooth enough for transformers to approximate.

Proposition G.4. Under Assumptions 5.2 and G.3, if $d_F \geq 16d_y$, $B_{A,1} \geq 16Rd_y$, $B_{A,2} \geq d_F$, $B_A \geq \sqrt{d_y}$, and $B_V \geq \sqrt{d}$, then

$$\max_{\|S^T\|_{2,\infty} \leq R} \text{KL}(\mathbb{P}^*(\cdot | S) \| \mathbb{P}_{\theta^*}(\cdot | S)) = \mathcal{O}\left(d_y \exp\left(-\frac{D^{1/4}}{\sqrt{C^2 B^2 \log B_{A,1}}}\right)\right),$$

for some constant $C > 0$.

Proof of Proposition G.4. Our proof mainly involves three steps.

- The high-level introduction of transformer approximator for g^* .
- Build the approximators in the transformer for ϕ^* and ρ_i^* separately.
- Conclude the proof.

Step 1: The high-level introduction of transformer approximator for g^* .

Without loss of generality, we assume that $B > 1$ in Assumption G.7. We would like to first introduce our construction in a high-level way. As shown in Figure 10, we will construct Ψ_{ϕ^*} and Ψ_{ρ^*} to respectively approximate ϕ^* and $\tau \log \rho^*$.

To approximate ϕ^* with Ψ_{ϕ^*} , we will make use of the universal approximation property of the fully-connected networks and ignore the attention module in the transformer by setting $W_i^{V,(t)} = 0$, $\gamma_1^{(t)} = 1$, $\gamma_2^{(t)} = 0$ for all $i \in [h]$. We further set $A_2^{(t)} = I_{d_F} \in \mathbb{R}^{d_F \times d_F}$, which is the identity matrix. The network structure for Ψ_{ϕ^*} is

$$X^{(t+1)} = \Pi_{\text{norm}}[\text{ReLU}(X^{(t)} A_1^{(t+1)} + b^{(t+1)} \cdot \mathbb{I}_L)],$$

where $b^{(t+1)} \in \mathbb{R}$ is the bias term. In Step 2, we will use this fully-connected network to approximate ϕ^* .

To approximate the average $\frac{1}{L} \sum_{i=1}^L \phi^*(x_i)$, we take $W_i^{Q,(t)} = 0$, $W_i^{K,(t)} = 0$, and $W_i^{V,(t)} = \mathbb{I}_d$, $\gamma_1^{(t)} = 0$, $\gamma_2^{(t)} = 1$, $A_2^{(t)} = 0$.

After this average aggregation, we still take $W_i^{V,(t)} = 0$, $\gamma_1^{(t)} = 1$, $\gamma_2^{(t)} = 0$ for all $i \in [h]$ and $A_2^{(t)} = I_{d_F} \in \mathbb{R}^{d_F \times d_F}$ to approximate ρ_i^* for $i \in [d_y]$. We stack the approximators for $\tau \log \rho_i^*$ to approximate $\tau \log \rho^*$, multiplying the width of the networks by d_F .

Step 2: Build the approximators in the transformer for ϕ^* and ρ_i^* separately.

In the first and the D^{th} layer, we take $A_1^{(1),'} = A_1^{(1)}/R$ and $A_1^{(D),'} = A_1^{(D)} \cdot R$ to normalize and retrieve the magnitudes of inputs, where R is the range of the inputs. This will keep the magnitudes of the intermediate outputs small. Next, we will use Lemma J.9 to construct the networks. In the proof of Lemma J.9, the norm of the outputs of the intermediate layers do not exceed the range of the inputs, so the layer normalization in our networks will not influence the constructed approximators. In this case, we can respectively approximate ϕ^* and $\tau \log \rho_i^*$ with fully-connected networks Ψ_{ϕ^*} and $\Psi_{\rho_i^*}$ for $i \in [d_y]$ as

$$\|\phi^* - \Psi_{\phi^*}\|_{\infty} \leq \varepsilon_{\phi}, \quad \|\tau \log \rho_i^* - \Psi_{\rho_i^*}\|_{\infty} \leq \varepsilon_{\rho} \text{ for } i \in [d_y],$$

where the depth $D(\cdot)$, the width $W(\cdot)$, and the maximal weight $B(\cdot)$ of the networks satisfy that

$$D' = D(\Psi_{\phi^*}) \leq C \cdot B \cdot (\log \varepsilon_{\phi}^{-1})^2 + \log B, \quad D'' = \max_{i \in [d_y]} D(\Psi_{\rho_i^*}) \leq C \cdot B \cdot (\log \varepsilon_{\rho}^{-1})^2 + \log B,$$

$$W(\Psi_{\phi^*}) \leq 16, \quad W(\Psi_{\rho_i^*}) \leq 16, \quad B(\Psi_{\phi^*}) \leq 1, \quad B(\Psi_{\rho_i^*}) \leq 1$$

for some constant $C > 0$. The bounds for width and maximal weight require that $d_F \geq 16d_y$ and $B_{A,1} \geq \sqrt{d_F \cdot \overline{d_F}} \geq 16d_y$. Then we have that for any $X = (x_1, \dots, x_L)$

$$\begin{aligned} & \left\| \rho^* \left(\frac{1}{L} \sum_{i=1}^L \phi^*(x_i) \right) - \text{softmax} \left(\frac{1}{\tau} \Psi_{\rho^*} \left(\frac{1}{L} \sum_{i=1}^L \Psi_{\phi^*}(x_i) \right) \right) \right\|_1 \\ & \leq \left\| \rho^* \left(\frac{1}{L} \sum_{i=1}^L \phi^*(x_i) \right) - \text{softmax} \left(\frac{1}{\tau} \Psi_{\rho^*} \left(\frac{1}{L} \sum_{i=1}^L \phi^*(x_i) \right) \right) \right\|_1 \\ & \quad + \left\| \text{softmax} \left(\frac{1}{\tau} \Psi_{\rho^*} \left(\frac{1}{L} \sum_{i=1}^L \phi^*(x_i) \right) \right) - \text{softmax} \left(\frac{1}{\tau} \Psi_{\rho^*} \left(\frac{1}{L} \sum_{i=1}^L \Psi_{\phi^*}(x_i) \right) \right) \right\|_1 \\ & \leq d_y \varepsilon_{\rho} + C' \cdot d_y \cdot (B_{A,1})^{D''} \cdot \varepsilon_{\phi}, \end{aligned} \tag{G.12}$$

where $C' > 0$ is a constant, the first inequality results from the triangle inequality, $(B_{A,1})^{D''}$ in the second inequality results from the error propagation through a depth- D'' network and the Lipschitzness of softmax in Lemma J.6. This bound reflects that the later modules will amplify the approximation error in the previous modules. In the following, we will balance the depths of different modules to handle the amplification. Lemma J.9 indicates the approximation error ε of a fully-connected network with depth D can be upper bounded as

$$\varepsilon \leq \exp\left(-\sqrt{\frac{D - \log B}{B}}\right).$$

Thus, defining the left-hand side of (G.12) as approx err, we have that

$$\text{approx err} \leq d_y \exp\left(-\sqrt{\frac{D'' - \log B}{B}}\right) + d_y B_{A,1}^{D''} \exp\left(-\sqrt{\frac{D' - \log B}{B}}\right).$$

We note the fact that: for any $l > 0, c > 0$, we have $\exp(-l\sqrt{x-c}) = O(\exp(-l\sqrt{x}))$, which follows from the direct calculation. Then we can further upper bound the approximation error as

$$\text{approx err} = O\left(d_y \exp\left(-\sqrt{\frac{D''}{B}}\right) + d_y \exp\left(\frac{1}{\sqrt{B}}[D''\sqrt{B} \log B_{A,1} - \sqrt{D'}]\right)\right).$$

To handle the second term in the right-hand side of this inequality, we require that

$$k \cdot D'' - \sqrt{D'} \leq -\sqrt{D''},$$

where $k = \sqrt{B} \log B_{A,1}$. This is equivalent to

$$D' \geq (k \cdot D'' + \sqrt{D''})^2.$$

Since $D' + D'' \leq D$, where D is the depth of the whole network, we can set

$$D'' = \sqrt{D}/(2\sqrt{B} \log B_{A,1}), \quad D' = D - 1 - D'' \geq D/2 + D^{3/4}$$

when D is large. This assignments ensure that $D' \geq (k \cdot D'' + \sqrt{D''})^2$. Thus, we have that

$$\text{approx err} = O\left(d_y \exp\left(-\sqrt{\frac{D''}{B}}\right)\right) = O\left(d_y \exp\left(-\frac{D^{1/4}}{\sqrt{C^2 B^2 \log B_{A,1}}}\right)\right)$$

for some constant $C > 0$. Here we relax the dependency on B a little for the notational clearness, and the relaxation results from the fact that $B \geq 1$ usually.

Step 3: Conclude the proof.

We denote $\Psi_{\rho^*}(\sum_{i=1}^L \Psi_{\phi^*}(x_i)/L)$ as \mathbb{P}_{θ^*} . Then if $\text{TV}(\mathbb{P}(\cdot | X), \mathbb{P}_{\theta^*}(\cdot | X)) = \varepsilon \leq c_0/2$, some basic calculations show that

$$\frac{c_0}{c_0 + \varepsilon} \leq \frac{\mathbb{P}(x | S)}{\mathbb{P}_{\theta^*}(x | S)} \leq 1 + \frac{2\varepsilon}{c_0}.$$

Thus, we have

$$\max_{\|S^\top\|_{2,\infty} \leq R} \text{KL}(\mathbb{P}(\cdot | S) \parallel \mathbb{P}_{\theta^*}(\cdot | S)) \leq \frac{2\varepsilon}{c_0} = \mathcal{O}\left(d_y \exp\left(-\frac{D^{1/4}}{\sqrt{C^2 B^2 \log B_{A,1}}}\right)\right).$$

□

G.4 PRETRAINING RESULTS FOR ℓ_2 LOSS

G.4.1 PRETRAINING ALGORITHM WITH ℓ_2 LOSS

Training with ℓ_2 loss is common in the CV community, e.g. Radford et al. (2021). The network structure is largely similar to those in Brown et al. (2020) and Devlin et al. (2018). Here, we modify the network structure of the last layer. The network derives the final output as $Y^{(D+1)} = \frac{1}{L} \mathbb{1}_L^\top X^{(D)} A^{(D+1)}$, where $\mathbb{1}_L \in \mathbb{R}^L$ is the vector with all ones, $A^{(D+1)} \in \mathbb{R}^{d \times d_y}$. The parameters in each layer are $\theta^{(t)} = (\gamma_1^{(t)}, \gamma_2^{(t)}, W^{(t)}, A^{(t)})$ for $t \in [D]$, and $\theta^{(D+1)} = A^{(D+1)}$, and the parameters of the whole network is $\theta = (\theta^{(1)}, \dots, \theta^{(D+1)})$. Similar to Section 5.1, we consider the transformer with bounded weights. The set of parameters is

$$\Theta = \left\{ \theta \mid \|A^{(D+1)}\|_F \leq B_A, \max\{|\gamma_1^{(t)}|, |\gamma_2^{(t)}|\} \leq 1, \|A_1^{(t)}\|_F \leq B_{A,1}, \|A_2^{(t)}\|_F \leq B_{A,2}, \right. \\ \left. \|W_i^{Q,(t)}\|_F \leq B_Q, \|W_i^{K,(t)}\|_F \leq B_K, \|W_i^{V,(t)}\|_F \leq B_V \text{ for all } t \in [D], i \in [h] \right\},$$

where $B_A, B_{A,1}, B_{A,2}, B_Q, B_K$, and B_V are the bounds of parameter. We only consider the non-trivial case where these bounds are larger than 1, otherwise the magnitude of the output in D^{th} layer decays exponentially with growing depth. We denote the transformer with parameter θ as f_θ .

In such case, we focus on the pretraining setting in CV tasks, i.e., the pretraining set $\mathcal{D} = \{(S^i, x^i)\}_{i=1}^N$ consists of i.i.d. pairs. The underlying distribution is denoted as $(S, x) \sim \mu \in \Delta(\mathcal{X}^* \times \mathcal{X})$. In such case, $d = d_y$, i.e., the transformer directly predicts the masked token. The training algorithm is

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \|x^i - f_\theta(S^i)\|_2^2 \quad (\text{G.13})$$

From the population version of (G.13), it is easy to see that the function $f^*(S) = \mathbb{E}[x | S]$ achieves the minimal population error, where the conditional expectation is defined from μ . In the following, we will quantify the error between $f_{\hat{\theta}}$ and f^* .

G.4.2 PERFORMANCE GUARANTEE FOR PRETRAINING WITH ℓ_2 LOSS

We first state the assumptions for the pretraining setting.

Assumption G.5. There exists a constant $R > 0$ such that for $(S, x) \sim \mu$, we have $\|S^\top\|_{2,\infty} \leq R$ and $\|x\|_2 \leq B_x$ almost surely.

Then the performance guarantee for the pretraining result $\hat{\theta}$ can be derived as following.

Theorem G.6. Let $\bar{B} = B_x R h B_A B_{A,1} B_{A,2} B_Q B_K B_V$ and $\bar{D} = D^2 d(d_F + d_h + d) + d \cdot d_y$. If Assumption G.5 holds, the pretrained model $f_{\hat{\theta}}$ by the algorithm in (G.13) satisfies

$$\mathbb{E}_{S,x} \left[\|f^*(S) - f_{\hat{\theta}}(S)\|_2^2 \right] \leq \underbrace{\frac{3}{2} \min_{\theta \in \Theta} \mathbb{E} \left[\|f^*(S) - f_{\theta}(S)\|_2^2 \right]}_{\text{approximation error}} + \underbrace{\mathcal{O} \left(\frac{B_x^2}{N} \left[\bar{D} \log(1 + N\bar{B}) + \log \frac{2}{\delta} \right] \right)}_{\text{generalization error}}$$

with probability at least $1 - \delta$.

The first term is the approximation error. It measures the proximity between the nominal function f^* and the functions induced by the parameter set Θ . The second term is the generalization error. Similar as Theorem 5.3, the generalization error is independent of the token sequence length.

Since the neural networks are universal approximators, we will explicitly approximate f^* from the transformer function class. Theorem 2 in Zaheer et al. (2017) shows that there exist $\rho^* : \mathbb{R} \rightarrow \mathbb{R}^{d_y}$ and $\phi^* : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$f^*(X) = \rho^* \left(\frac{1}{L} \sum_{i=1}^L \phi^*(x_i) \right),$$

where $X = [x_1, \dots, x_L]$. The i^{th} component of ρ^* is denoted as ρ_i^* for $i \in [d_y]$. For a function f defined on Ω , the L^∞ norm of it is defined as $\|f\|_\infty = \sup_{x \in \Omega} |f(x)|$. The set of the real-valued smooth functions on it is denoted as $\mathcal{S}^\infty(\Omega, \mathbb{R})$. Then we denote the set of the smooth functions with bounded derivatives as

$$\mathcal{S}_B = \left\{ f \in \mathcal{S}^\infty([-B, B], \mathbb{R}) \mid \|f^{(n)}(x)\| \leq n! \text{ for all } n \in \mathbb{N} \right\},$$

where $f^{(n)}$ is the n^{th} -order derivative of f .

Assumption G.7. There exists $B > 0$ such that $\phi^*, \rho_i^* \in \mathcal{S}_B$ for $i \in [d_y]$.

This assumption states that the function f^* is smooth enough. Then we have that

Proposition G.8. Under G.7, if $d_F \geq 16d_y$, $B_{A,1} \geq 16Rd_y$, $B_{A,2} \geq d_F$, $B_A \geq \sqrt{d_y}$, and $B_V \geq \sqrt{d}$, then

$$\max_{\|S^\top\|_{2,\infty} \leq R} \|f^*(S) - f_{\theta^*}(S)\|_2 = \mathcal{O} \left(d_y \exp \left(- \frac{D^{1/4}}{\sqrt{C^2 B^2 \log B_{A,1}}} \right) \right)$$

for some constant $C > 0$.

G.4.3 PROOF OF THEOREM G.6

Proof of Theorem G.6. For ease of notation, we respectively define the empirical risk and the population risk as

$$\hat{\mathcal{L}}(f, \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \|x^i - f_{\theta}(S^i)\|_2^2, \quad \mathcal{L}(f) = \mathbb{E}_{S,x} \left[\|x - f_{\theta}(S)\|_2^2 \right].$$

The our proof mainly involves three steps.

- Error decomposition for the excess population risk.
- Control each term in the error decomposition.
- Conclude the proof.

Step 1: Error decomposition for the excess population risk. The excess population risk for the estimate $\hat{\theta}$ can be decomposed to the sum of the generalization error and the approximation error as $\mathcal{L}(f_{\hat{\theta}}) - \mathcal{L}(f^*)$

$$\begin{aligned} &= \mathcal{L}(f_{\hat{\theta}}) - \mathcal{L}(f^*) - 2(\hat{\mathcal{L}}(f_{\hat{\theta}}, \mathcal{D}) - \hat{\mathcal{L}}(f^*, \mathcal{D})) + 2(\hat{\mathcal{L}}(f_{\hat{\theta}}, \mathcal{D}) - \hat{\mathcal{L}}(f_{\theta^*}, \mathcal{D})) + 2(\hat{\mathcal{L}}(f_{\theta^*}, \mathcal{D}) - \hat{\mathcal{L}}(f^*, \mathcal{D})) \\ &\leq \underbrace{\mathcal{L}(f_{\hat{\theta}}) - \mathcal{L}(f^*) - 2(\hat{\mathcal{L}}(f_{\hat{\theta}}, \mathcal{D}) - \hat{\mathcal{L}}(f^*, \mathcal{D}))}_{\text{generalization error}} + \underbrace{2(\hat{\mathcal{L}}(f_{\theta^*}, \mathcal{D}) - \hat{\mathcal{L}}(f^*, \mathcal{D}))}_{\text{approximation error}}, \end{aligned} \quad (\text{G.14})$$

where $\theta^* = \operatorname{argmin}_{\theta \in \Theta} \mathcal{L}(f_\theta)$, and the inequality results from that $\hat{\theta}$ achieves the minimal empirical risk.

Step 2: Control each term in the error decomposition.

We first consider the generalization error and will adapt Lemma J.2 to bound it. Define the function

$$g(S, x, \theta) = \|x - f_\theta(S)\|_2^2 - \|x - f^*(S)\|_2^2.$$

To verify the conditions in Lemma J.2, we notice that $|g(S, x, \theta)| \leq (B_x + B_f)^2$ and that

$$\begin{aligned} \mathbb{E}[g(S, x, \theta)] &= \mathbb{E}\left[\|x - f_\theta(S)\|_2^2 - \|x - f^*(S)\|_2^2\right] \\ &= \mathbb{E}\left[\|f^*(S) - f_\theta(S)\|_2^2\right] \\ \mathbb{E}\left[(g(S, x, \theta) - \mathbb{E}[g(S, x, \theta)])^2\right] &\leq \mathbb{E}\left[(g(S, x, \theta))^2\right] \\ &\leq \mathbb{E}\left[\|2x - f^*(S) - f_\theta(S)\|_2^2 \cdot \|f^*(S) - f_\theta(S)\|_2^2\right] \\ &\leq (3B_x + B_f)^2 \cdot \mathbb{E}\left[\|f^*(S) - f_\theta(S)\|_2^2\right], \end{aligned}$$

where the second equality results from the definition of f^* , the second inequality results from Cauchy–Schwarz inequality, and the last inequality result from the boundedness of x , f^* , and f_θ . Then Lemma J.2 shows that for a distribution $Q \in \Delta(\Theta)$ and $0 < \lambda \leq 1/(2(B_x + B_f)^2)$, the following holds with probability at least $1 - \delta$ simultaneously for all $P \in \Delta(\Theta)$

$$\begin{aligned} &\left| \mathbb{E}_{\theta \sim P} \left[\mathbb{E}[g(S, x, \theta)] - \frac{1}{N} \sum_{i=1}^N g(S^i, x^i, \theta) \right] \right| \\ &\leq \lambda(3B_x + B_f)^2 \mathbb{E}_{\theta \sim P} [\mathbb{E}[g(S, x, \theta)]] + \frac{1}{N\lambda} \left[\text{KL}(P \| Q) + \log \frac{2}{\delta} \right]. \end{aligned}$$

Taking $\lambda = 1/(2(3B_x + B_f)^2)$, we have

$$\begin{aligned} &\left| \mathbb{E}_{\theta \sim P} [\mathcal{L}(f_\theta) - \mathcal{L}(f^*) - (\hat{\mathcal{L}}(f_\theta, \mathcal{D}) - \hat{\mathcal{L}}(f^*, \mathcal{D}))] \right| \\ &\leq \frac{1}{2} \mathbb{E}_{\theta \sim P} [\mathcal{L}(f_\theta) - \mathcal{L}(f^*)] + \frac{2(3B_x + B_f)^2}{N} \left[\text{KL}(P \| Q) + \log \frac{2}{\delta} \right]. \end{aligned}$$

Next, we will take proper P and Q to relate this equation and the generalization error. For this purpose, we quantify how the perturbation of network parameters influence the output of the network.

Proposition G.9. For any input $X \in \mathbb{R}^{L \times d}$ and $\theta, \tilde{\theta} \in \Theta$, we have that

$$\|f_\theta(X) - f_{\tilde{\theta}}(X)\|_2 \leq \|A^{(D+1)} - \tilde{A}^{(D+1)}\|_F + \sum_{t=1}^D \alpha_t (\beta_t + \iota_t + \kappa_t + \rho_t),$$

where

$$\begin{aligned} \alpha_t &= B_A(1 + B_{A,1} \cdot B_{A,2})(1 + hB_V(1 + 4B_Q B_K))^{D-t} \\ \beta_t &= |\gamma_2^{(t)} - \tilde{\gamma}_2^{(t)}| + (1 + B_{A,1} \cdot B_{A,2}) \cdot (1 + (\|X^\top\|_{2,\infty} - 1)\mathbb{I}_{t=1}) \cdot |\gamma_1^{(t)} - \tilde{\gamma}_1^{(t)}| \\ \iota_t &= B_{A,2} \cdot \|A_1^{(t)} - \tilde{A}_1^{(t)}\|_F + B_{A,1} \cdot \|A_2^{(t)} - \tilde{A}_2^{(t)}\|_F \\ \kappa_t &= (1 + B_{A,1} \cdot B_{A,2}) \cdot (1 + (\|X^\top\|_{2,\infty} - 1)\mathbb{I}_{t=1}) \cdot \sum_{i=1}^h \|W_i^{V,(t)} - \tilde{W}_i^{V,(t)}\|_F \\ \rho_t &= 2(1 + B_{A,1} \cdot B_{A,2}) \cdot (1 + (\|X^\top\|_{2,\infty} - 1)\mathbb{I}_{t=1}) \cdot B_V \\ &\quad \cdot \sum_{i=1}^h B_K \cdot \|W_i^{Q,(t+1)} - \tilde{W}_i^{Q,(t+1)}\|_F + B_Q \cdot \|W_i^{K,(t+1)} - \tilde{W}_i^{K,(t+1)}\|_F \end{aligned}$$

for all $t \in [D]$.

Proof of Proposition G.9. See Appendix I.4. \square

With the help of Proposition G.9, we set the distribution P as

$$\begin{aligned}
 P &= \prod_{t=1}^{D+1} \mathcal{L}_P(\theta^{(t)}) \\
 \mathcal{L}_P(\theta^{(D+1)}) &= \text{Unif}\left(\mathbb{B}(\widehat{A}^{(D+1)}, r^{(D+1)}, \|\cdot\|_F)\right) \\
 \mathcal{L}_P(\theta^{(t)}) &= \text{Unif}\left(\mathbb{B}(\widehat{\gamma}_1^{(t)}, r_{\gamma,1}^{(t)}, |\cdot|)\right) \cdot \text{Unif}\left(\mathbb{B}(\widehat{\gamma}_2^{(t)}, r_{\gamma,2}^{(t)}, |\cdot|)\right) \cdot \mathcal{L}_P(A^{(t)}) \cdot \mathcal{L}_P(W^{(t)}) \\
 \mathcal{L}_P(A^{(t)}) &= \text{Unif}\left(\mathbb{B}(\widehat{A}_1^{(t)}, r_{A,1}^{(t)}, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(\widehat{A}_2^{(t)}, r_{A,2}^{(t)}, \|\cdot\|_F)\right) \\
 \mathcal{L}_P(W^{(t)}) &= \prod_{i=1}^h \text{Unif}\left(\mathbb{B}(\widehat{W}_i^{Q,(t)}, r_Q^{(t)}, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(\widehat{W}_i^{K,(t)}, r_K^{(t)}, \|\cdot\|_F)\right) \cdot \text{Unif}\left(\mathbb{B}(\widehat{W}_i^{V,(t)}, r_V^{(t)}, \|\cdot\|_F)\right)
 \end{aligned} \tag{G.15}$$

for $t \in [D]$, where Unif denotes the uniform distribution on the set, $\mathbb{B}(a, r, \|\cdot\|) = \{x \mid \|x - a\| \leq r\}$ denotes the ball centered in a with radius r , the radius is set as

$$\begin{aligned}
 r_{\gamma,1}^{(t)} &= (B_x + B_f)^{-1} R^{-1} (1 + B_{A,1} \cdot B_{A,2})^{-1} \alpha_t^{-1} / N, & r_{\gamma,2}^{(t)} &= (B_x + B_f)^{-1} R^{-1} \alpha_t^{-1} / N \\
 r_{A,1}^{(t)} &= (B_x + B_f)^{-1} R^{-1} B_{A,2}^{-1} \alpha_t^{-1} / N, & r_{A,2}^{(t)} &= (B_x + B_f)^{-1} R^{-1} B_{A,1}^{-1} \alpha_t^{-1} / N, \\
 r_V^{(t)} &= (B_x + B_f)^{-1} R^{-1} h^{-1} (1 + B_{A,1} \cdot B_{A,2})^{-1} \alpha_t^{-1} / N, & r^{(D+1)} &= (B_x + B_f)^{-1} B_A^{-1} / N, \\
 r_K^{(t)} &= (B_x + B_f)^{-1} R^{-1} h^{-1} (1 + B_{A,1} \cdot B_{A,2})^{-1} B_V^{-1} B_Q^{-1} \alpha_t^{-1} / N, \\
 r_Q^{(t)} &= (B_x + B_f)^{-1} R^{-1} h^{-1} (1 + B_{A,1} \cdot B_{A,2})^{-1} B_V^{-1} B_K^{-1} \alpha_t^{-1} / N.
 \end{aligned}$$

Under this assignment, we now bound $\mathbb{E}_{\theta \sim P} [\|x - f_\theta(S)\|_2^2 - \|x - f_{\hat{\theta}}(S)\|_2^2]$ as

$$\left| \mathbb{E}_{\theta \sim P} [\|x - f_\theta(S)\|_2^2 - \|x - f_{\hat{\theta}}(S)\|_2^2] \right| \leq 2(B_x + B_f) \left| \mathbb{E}_{\theta \sim P} [\|f_\theta(S) - f_{\hat{\theta}}(S)\|_2] \right| = \mathcal{O}\left(\frac{B_x + B_f}{N}\right),$$

where the inequality results from Cauchy-Schwarz inequality, and the equality results from Proposition G.9. Thus, we have that

$$\begin{aligned}
 &\mathcal{L}(f_{\hat{\theta}}) - \mathcal{L}(f^*) - (\widehat{\mathcal{L}}(f_{\hat{\theta}}, \mathcal{D}) - \widehat{\mathcal{L}}(f^*, \mathcal{D})) \\
 &\leq \frac{1}{2} (\mathcal{L}(f_{\hat{\theta}}) - \mathcal{L}(f^*)) + \mathcal{O}\left(\frac{B_x + B_f}{N}\right) + \frac{2(3B_x + B_f)^2}{N} \left[\text{KL}(P \parallel Q) + \log \frac{2}{\delta} \right].
 \end{aligned} \tag{G.16}$$

To access to the value of $\text{KL}(P \parallel Q)$, we take Q as the distribution in (G.8) except that

$$\mathcal{L}_Q(\theta^{(D+1)}) = \text{Unif}\left(\mathbb{B}(0, B_A, \|\cdot\|_F)\right). \tag{G.17}$$

Then the KL divergence between P and Q is

$$\text{KL}(P \parallel Q) = \mathcal{O}\left((D^2 \cdot d \cdot (d_F + d_h + d) + d \cdot d_y) \cdot \log(1 + NB_x R h B_A B_{A,1} B_{A,2} B_Q B_K B_V)\right).$$

Combining this equality with (G.16), we have that with probability at least $1 - \delta$, the generalization error can be bounded as

$$\mathcal{L}(f_{\hat{\theta}}) - \mathcal{L}(f^*) - 2(\widehat{\mathcal{L}}(f_{\hat{\theta}}, \mathcal{D}) - \widehat{\mathcal{L}}(f^*, \mathcal{D})) = \mathcal{O}\left(\frac{B_x^2}{N} \left[\bar{D} \log(1 + N\bar{B}) + \log \frac{2}{\delta} \right]\right). \tag{G.18}$$

Next we control the approximation error in (G.14).

$$\begin{aligned}
 &\widehat{\mathcal{L}}(f_{\theta^*}, \mathcal{D}) - \widehat{\mathcal{L}}(f^*, \mathcal{D}) \\
 &= \widehat{\mathcal{L}}(f_{\theta^*}, \mathcal{D}) - \widehat{\mathcal{L}}(f^*, \mathcal{D}) - \frac{3}{2} (\mathcal{L}(f_{\theta^*}) - \mathcal{L}(f^*)) + \frac{3}{2} (\mathcal{L}(f_{\theta^*}) - \mathcal{L}(f^*)) \\
 &= \widehat{\mathcal{L}}(f_{\theta^*}, \mathcal{D}) - \widehat{\mathcal{L}}(f^*, \mathcal{D}) - \frac{3}{2} (\mathcal{L}(f_{\theta^*}) - \mathcal{L}(f^*)) + \frac{3}{2} \mathbb{E} [\|f^*(S) - f_{\theta^*}(S)\|_2^2],
 \end{aligned} \tag{G.19}$$

where the second equality results from the definition of f^* . To bound the first two terms in the right-hand side of (G.19), we use Lemma J.2 and take P and Q as (G.15) and (G.17), replacing $\hat{\theta}$ by θ^* . Then we have that

$$\hat{\mathcal{L}}(f_{\theta^*}, \mathcal{D}) - \hat{\mathcal{L}}(f^*, \mathcal{D}) - \frac{3}{2}(\mathcal{L}(f_{\theta^*}) - \mathcal{L}(f^*)) = \mathcal{O}\left(\frac{B_x^2}{N} \left[\bar{D} \log(1 + N\bar{B}) + \log \frac{2}{\delta} \right]\right). \quad (\text{G.20})$$

Step 3: Conclude the proof.

Combining inequalities (G.14), (G.18), (G.19), and (G.20), we have that

$$\mathcal{L}(f_{\hat{\theta}}) - \mathcal{L}(f^*) = \frac{3}{2} \mathbb{E} \left[\|f^*(S) - f_{\theta^*}(S)\|_2^2 \right] + \mathcal{O}\left(\frac{B_x^2}{N} \left[\bar{D} \log(1 + N\bar{B}) + \log \frac{2}{\delta} \right]\right).$$

Thus, we conclude the proof of Theorem G.6. \square

G.4.4 PROOF OF PROPOSITION G.8

Proof of Proposition G.8. Our proof mainly involves three steps.

- Build the high-level transformer approximator for f^* .
- Build the approximators in the transformer for ϕ^* and ρ_i^* separately.
- Conclude the proof.

The first two steps follow the procedures of the proof of Proposition G.4 exactly. Now we present the final step.

Step 3: Conclude the proof.

In the final layer, we just take $A^{(D+1)} = I_{d_y}$ as the identity matrix. Denoting the derived parameters as θ^* we have that

$$\max_{\|X^\top\|_{2,\infty} \leq R} \left\| \rho^* \left(\frac{1}{L} \sum_{i=1}^L \phi^*(x_i) \right) - f_{\theta^*}(X) \right\|_2 = \mathcal{O}\left(d_y \exp\left(-\frac{D^{1/4}}{\sqrt{C^2 B^2 \log B_{A,1}}}\right)\right).$$

Thus, we conclude the proof of Proposition G.8. \square

H PROOFS AND FORMAL STATEMENTS FOR §6

H.1 PROOF OF THEOREM 6.2

Proof. By Corollary 4.2 and the fact that $\log(1/p_0(z_*)) \leq \beta$, we have that

$$T^{-1} \cdot \mathbb{E}_{\mathcal{D}_{\text{ICL}}} \left[\sum_{t=1}^T \log \mathbb{P}(r_t | z^*, \text{prompt}_{t-1}) - \sum_{t=1}^T \log \mathbb{P}(r_t | \text{prompt}_{t-1}) \right] \leq \beta/T. \quad (\text{H.1})$$

In addition, we have that

$$T^{-1} \cdot \mathbb{E}_{\mathcal{D}_{\text{ICL}}} \left[\sum_{t=1}^T \log \mathbb{P}(r_t | \text{prompt}_{t-1}) - \sum_{t=1}^T \log \mathbb{P}_{\hat{\theta}}(r_t | \text{prompt}_{t-1}) \right] = \mathbb{E}_{\mathcal{D}_{\text{ICL}}} \left[\text{KL}(\mathbb{P}(\cdot | \text{prompt}) \parallel \mathbb{P}_{\hat{\theta}}(\cdot | \text{prompt})) \right]. \quad (\text{H.2})$$

Similar to (G.10), we have that

$$\left| \log(\mathbb{P}(r | \text{prompt}) / \mathbb{P}_{\hat{\theta}}(r | \text{prompt})) \right| \leq b^* = \log \max\{c_0^{-1}, b_y^{-1}\}.$$

By Lemma J.10, we have that

$$\text{KL}(\mathbb{P}(\cdot | \text{prompt}) \parallel \mathbb{P}_{\hat{\theta}}(\cdot | \text{prompt})) \leq (3 + b^*)/2 \cdot \text{TV}(\mathbb{P}(\cdot | \text{prompt}), \mathbb{P}_{\hat{\theta}}(\cdot | \text{prompt})). \quad (\text{H.3})$$

By Assumption 6.1, we have that $\mathbb{P}_{\mathcal{D}_{\text{ICL}}}(\text{prompt}) \leq \kappa \mathbb{P}_{\mathcal{D}}(\text{prompt})$. Thus, by Theorem 5.3, we have with probability at least $1 - \delta$ that

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}_{\text{ICL}}} \left[\text{KL}(\mathbb{P}(\cdot | \text{prompt}) \| \mathbb{P}_{\hat{\theta}}(\cdot | \text{prompt})) \right] \\ & \leq C \cdot b^* \cdot \kappa \cdot \mathbb{E}_{S \sim \mathcal{D}} \left[\text{TV}(\mathbb{P}(\cdot | S), \mathbb{P}_{\hat{\theta}}(\cdot | S)) \right] \leq C \cdot b^* \cdot \kappa \cdot \Delta_{\text{pre}}(N, T, \delta). \end{aligned} \quad (\text{H.4})$$

Combining (H.4), (H.1), and (H.2), we have with probability at least $1 - \delta$ that

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}_{\text{ICL}}} \left[T^{-1} \cdot \sum_{t=1}^T \log \mathbb{P}(r_t | z^*, \text{prompt}_{t-1}) - T^{-1} \cdot \sum_{t=1}^T \log \mathbb{P}_{\hat{\theta}}(r_t | \text{prompt}_{t-1}) \right] \\ & \leq \beta/T + \mathbb{E}_{S \sim \mathcal{D}} \left[\text{KL}(\mathbb{P}(\cdot | S) \| \mathbb{P}_{\hat{\theta}}(\cdot | S)) \right] \\ & \leq \mathcal{O}(\beta/T + b^* \cdot \kappa \cdot \Delta_{\text{pre}}(N, T, \delta)), \end{aligned} \quad (\text{H.5})$$

which completes the proof of Theorem 6.2. \square

H.2 ASSUMPTIONS AND FORMAL STATEMENT FOR PROMPTING WITH WRONG INPUT-OUTPUT MAPPINGS

We first state assumptions for this setting.

Assumption H.1. Conditioned on any $z \in \mathfrak{Z}$, the input-output pairs are independent, i.e., for any two input-output pair sequences $S_t, S'_{t'} \in \mathfrak{X}^*$, we have $\mathbb{P}((S_t, S'_{t'}) | z) = \mathbb{P}(S_t | z) \cdot \mathbb{P}(S'_{t'} | z)$. This assumption states that for any task $z \in \mathfrak{Z}$, the input-output pairs are independently generated. This largely holds in realistic applications since the examples usually are independently produced. It can be relaxed when there are more structures in the token generation process, e.g. the hidden Markov model in Xie et al. (2021).

Assumption H.2. There exists a constant $c_1 > 0$ such that $\mathbb{P}_{\mathcal{Z}}(z_*) \geq c_1$. This assumption states that the prior distribution of the hidden concept z_* is strictly larger than 0, otherwise this concept can never be deduced. For two concepts $z, z' \in \mathfrak{Z}$, we define the KL divergence between the conditional distributions of input-output pair on them as $\text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z) \| \mathbb{P}(\cdot | z')) = \mathbb{E}_{X, y \sim \mathbb{P}(\cdot | z)} [\log(\mathbb{P}(X, y | z) / \mathbb{P}(X, y | z'))]$. This divergence measures the distance between distributions of input-output pairs conditioned on different tasks z and z' .

Assumption H.3. The concept z_* satisfies that $\min_{z \neq z_*} \text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z_*) \| \mathbb{P}(\cdot | z)) > 2 \log 1/c_0$, where c_0 is the constant in Assumption 5.2.

This distinguishability assumption requires that the divergence between z_* and other concepts z is large enough to infer the concept z_* from the prompt. We denote the pretraining error in Theorem 5.3 as $\Delta_{\text{pre}}(N_p, T_p, \delta)$, then we have the following result.

Proposition H.4. Under Assumptions 5.2, 6.1 H.1, H.2, and H.3, the pretrained model $\mathbb{P}_{\hat{\theta}}$ in (5.2) predicts the outputs with the prompt containing wrong mappings as

$$\begin{aligned} & \mathbb{E}_{\text{prompt}'} \left[\text{KL}(\mathbb{P}(\cdot | \tilde{c}_{t+1}, z_*) \| \mathbb{P}_{\hat{\theta}}(\cdot | S'_t, \tilde{c}_{t+1})) \right] \\ & = \mathcal{O} \left(\kappa \Delta_{\text{pre}}(N_p, T_p, \delta) + \exp \left(- \frac{\sqrt{t}}{2(1+t) \log 1/c_0} \left(\min_{z \neq z_*} \text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z_*) \| \mathbb{P}(\cdot | z)) + 2 \log c_0 \right) \right) \right) \end{aligned}$$

with probability at least $1 - \delta$.

H.3 PROOF OF PROPOSITION H.4

Proof of Proposition H.4. From Bayesian model averaging, the output distribution is

$$\begin{aligned} & \mathbb{P}(r_{t+1} | S'_t, \tilde{c}_{t+1}) \\ & = \sum_{z \in \mathfrak{Z}} \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, z) \cdot \mathbb{P}_{\mathcal{Z}}(z | S'_t) \\ & = \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, z^*) + \sum_{z \neq z^*} (\mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, z) - \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, z^*)) \cdot \mathbb{P}_{\mathcal{Z}}(z | S'_t) \\ & = \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, z^*) + \sum_{z \neq z^*} (\mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, z) - \mathbb{P}(r_{t+1} | \tilde{c}_{t+1}, z^*)) \cdot \mathbb{P}_{\mathcal{Z}}(z^* | S'_t) \cdot \frac{\mathbb{P}_{\mathcal{Z}}(z) \cdot \mathbb{P}(S'_t | z)}{\mathbb{P}_{\mathcal{Z}}(z^*) \cdot \mathbb{P}(S'_t | z^*)}, \end{aligned} \quad (\text{H.6})$$

where the first equality results from Bayesian model averaging, the last equality results from Bayes' theorem. Next, we upperbound the ratio $\mathbb{P}(S'_t | z) / \mathbb{P}(S'_t | z^*)$ in the right-hand side of Eqn. (H.6). We have that

$$\frac{1}{t} \log \frac{\mathbb{P}(S'_t | z)}{\mathbb{P}(S'_t | z^*)} = \frac{1}{t} \sum_{i=1}^t \log \frac{\mathbb{P}((\tilde{c}_i, r'_i) | z)}{\mathbb{P}((\tilde{c}_i, r'_i) | z^*)} \leq -2 \log c_0 + \frac{1}{t} \sum_{i=1}^t \log \frac{\mathbb{P}((\tilde{c}_i, r_i) | z)}{\mathbb{P}((\tilde{c}_i, r_i) | z^*)},$$

where the equality results from Assumption H.1, and the inequality results from Assumption 5.2. Assumption 5.2 also implies that $|\log \mathbb{P}((\tilde{c}_i, r_i) | z) / \mathbb{P}((\tilde{c}_i, r_i) | z^*)| \leq (1+l) \log 1/c_0$. Hoeffding inequality shows that with probability at least $1 - \delta$, we have

$$\frac{1}{t} \sum_{i=1}^t \log \frac{\mathbb{P}((\tilde{c}_i, r_i) | z)}{\mathbb{P}((\tilde{c}_i, r_i) | z^*)} + \text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z^*) \| \mathbb{P}(\cdot | z)) \leq \frac{(1+l)}{\sqrt{t}} \log \frac{1}{c_0} \cdot \log \frac{1}{\delta}.$$

Thus, we have that with probability at least $1 - \delta$, the following holds for all $z \neq z^*$

$$\frac{\mathbb{P}(S'_t | z)}{\mathbb{P}(S'_t | z^*)} \leq \exp \left(-t \left(\text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z^*) \| \mathbb{P}(\cdot | z)) + 2 \log c_0 - \frac{(1+l)}{\sqrt{t}} \log \frac{1}{c_0} \cdot \log \frac{|3|}{\delta} \right) \right).$$

Combining this inequality with Eqn. (H.6), we have that

$$\begin{aligned} & \text{TV}(\mathbb{P}(\cdot | S'_t, \tilde{c}_{t+1}), \mathbb{P}(\cdot | \tilde{c}_{t+1}, z^*)) \\ &= \mathcal{O} \left(\frac{1}{c_1} \exp \left(-t \left(\min_{z \neq z^*} \text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z^*) \| \mathbb{P}(\cdot | z)) + 2 \log c_0 - \frac{(1+l)}{\sqrt{t}} \log \frac{1}{c_0} \cdot \log \frac{|3|}{\delta} \right) \right) \right). \end{aligned} \quad (\text{H.7})$$

Taking expectations with respect to the distribution of S'_t, \tilde{c}_{t+1} on the both sides in (H.7), we have that

$$\begin{aligned} & \mathbb{E}_{\text{prompt}'} [\text{TV}(\mathbb{P}(\cdot | S'_t, \tilde{c}_{t+1}), \mathbb{P}(\cdot | \tilde{c}_{t+1}, z^*))] \\ &= \mathcal{O} \left(\frac{1}{c_1} \exp \left(-t \left(\min_{z \neq z^*} \text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z^*) \| \mathbb{P}(\cdot | z)) + 2 \log c_0 - \frac{(1+l)}{\sqrt{t}} \log \frac{1}{c_0} \cdot \log \frac{|3|}{\delta} \right) \right) \right) + \delta. \end{aligned} \quad (\text{H.8})$$

We set $\delta = |3 \exp(-a\sqrt{t}/2b)|$, where $a = \min_{z \neq z^*} \text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z^*) \| \mathbb{P}(\cdot | z)) + 2 \log c_0$, $b = -(1+l) \log c_0$. Then the right-hand side of (H.8) can be upper bounded as

$$\begin{aligned} & \mathbb{E}_{\text{prompt}'} [\text{TV}(\mathbb{P}(\cdot | S'_t, \tilde{c}_{t+1}), \mathbb{P}(\cdot | \tilde{c}_{t+1}, z^*))] \\ &= \mathcal{O} \left(\exp \left(-\frac{\sqrt{t}}{2(1+l) \log 1/c_0} \left(\min_{z \neq z^*} \text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z^*) \| \mathbb{P}(\cdot | z)) + 2 \log c_0 \right) \right) \right). \end{aligned}$$

Let $\mathbb{E}_{\text{prompt}'} [\text{TV}(\mathbb{P}(\cdot | S'_t, \tilde{c}_{t+1}), \mathbb{P}_{\hat{\theta}}(\cdot | S'_t, \tilde{c}_{t+1}))] \leq \kappa \Delta_{\text{pre}}(N_p, T_p, \delta)$, where $\Delta_{\text{pre}}(N_p, T_p, \delta)$ is the bound in Theorem 5.3. Then we have that

$$\begin{aligned} & \mathbb{E}_{\text{prompt}'} [\text{KL}(\mathbb{P}(\cdot | \tilde{c}_{t+1}, z^*) \| \mathbb{P}_{\hat{\theta}}(\cdot | S'_t, \tilde{c}_{t+1}))] \\ &\leq \mathcal{O} \left(\mathbb{E}_{\text{prompt}'} [\text{TV}(\mathbb{P}_{\hat{\theta}}(\cdot | S'_t, \tilde{c}_{t+1}), \mathbb{P}(\cdot | \tilde{c}_{t+1}, z^*))] \right) \\ &= \mathcal{O} \left(\kappa \Delta_{\text{pre}}(N_p, T_p, \delta) + \exp \left(-\frac{\sqrt{t}}{2(1+l) \log 1/c_0} \left(\min_{z \neq z^*} \text{KL}_{\text{pair}}(\mathbb{P}(\cdot | z^*) \| \mathbb{P}(\cdot | z)) + 2 \log c_0 \right) \right) \right), \end{aligned}$$

where the first equality results from Assumption 5.2. Thus, we conclude the proof of Proposition H.4. \square

I PROOF OF SUPPORTING PROPOSITIONS

I.1 PROOF OF PROPOSITION F.1

Proof. Let a, b be two vectors in the $(d-1)$ -dimensional unit sphere \mathbb{S}^{d-1} . We first define the following vector,

$$c = (a^\top b) \cdot b - (a - (a^\top b) \cdot b) \in \mathbb{S}^{d-1}. \quad (\text{I.1})$$

By direct calculation, we have the following property of c defined in (I.1),

$$c^\top b = (a^\top b) \cdot \|b\|_2^2 - a^\top b + (a^\top b) \cdot \|b\|_2^2 = a^\top b. \quad (\text{I.2})$$

By (I.1) and (I.2), we have that

$$a + c = 2(a^\top b) \cdot b = 2(c^\top b) \cdot b = (a^\top b) \cdot b + (c^\top b) \cdot b. \quad (\text{I.3})$$

We now calculate the desired integration. Note that

$$\int_{\mathbb{S}^{d-1}} a \cdot \exp(a^\top b) da = b \cdot \int_{\mathbb{S}^{d-1}} (a^\top b) \exp(a^\top b) da + \int_{\mathbb{S}^{d-1}} (a - (a^\top b) \cdot b) \cdot \exp(a^\top b) da. \quad (\text{I.4})$$

For the second term on the right-hand side of (I.4), it follows from (I.1) and (I.2) and (I.3) that

$$\int_{\mathbb{S}^{d-1}} (a - (a^\top b) \cdot b) \cdot \exp(a^\top b) da = - \int_{\mathbb{S}^{d-1}} (c - (c^\top b) \cdot b) \cdot \exp(c^\top b) dc, \quad (\text{I.5})$$

where the equality follows from the fact that $dc = 2\|b\|_2^2 da - da = da$. By replacing c by a on the right-hand side of (I.5), we have

$$\int_{\mathbb{S}^{d-1}} (a - (a^\top b) \cdot b) \cdot \exp(a^\top b) da = - \int_{\mathbb{S}^{d-1}} (a - (a^\top b) \cdot b) \cdot \exp(a^\top b) da = 0 \quad (\text{I.6})$$

Finally, by plugging (I.6) into (I.4), we obtain that

$$\int_{\mathbb{S}^{d-1}} a \cdot \exp(a^\top b) da = b \cdot \int_{\mathbb{S}^{d-1}} (a^\top b) \exp(a^\top b) da.$$

Thus, by setting

$$C_1 = \int_{\mathbb{S}^{d-1}} (a^\top b) \exp(a^\top b) da, \quad \forall b \in \mathbb{S}^{d-1},$$

we complete the proof of Proposition F.1. Note that here C_1 is an absolute constant that does not depend on b due to the symmetry on the unit sphere. \square

I.2 PROOF OF PROPOSITION G.2

Proof of Proposition G.2. We note that $f(X)$ satisfies the condition in Lemma J.4 with $c_i = 2b/N$ for $i \in [N]$. Then Lemma J.4 shows that

$$\mathbb{E}_{f \sim P_0} \left[\mathbb{E}_X \left(\exp \left[\lambda (f(X) - \mathbb{E}f(X)) \right] \right) \right] \leq \exp \left(\frac{\lambda^2 \cdot b^2 \cdot t_{\min}}{2N} \right).$$

Take $\lambda = \sqrt{2N \log 2 / (b^2 t_{\min})}$. The Markov inequality shows that

$$P \left(\mathbb{E}_{f \sim P_0} \left(\exp \left[\lambda (f(X) - \mathbb{E}f(X)) \right] \right) \geq \frac{2}{\delta} \right) \leq \delta$$

for any $0 < \delta < 1$. We note that this probability inequality does not involve P . Take the function g in Lemma J.3 as $g(f) = \lambda(f(X) - \mathbb{E}f(X))$, then it shows that

$$\log \mathbb{E}_{P_0} \left[\exp(g(X)) \right] + \text{KL}(P \| P_0) \geq \mathbb{E}_P[g(X)]$$

for any P simultaneously. Combining these inequalities, we have

$$\left| \mathbb{E}_P \left[\mathbb{E}_X[f(X)] - f(X) \right] \right| \leq \sqrt{\frac{b^2 \cdot t_{\min}}{2 \log 2N}} \left[\text{KL}(P \| P_0) + \log \frac{4}{\delta} \right],$$

for any distribution P on \mathcal{F} simultaneously with probability at least $1 - \delta$. Thus, we conclude the proof of Proposition G.2. \square

I.3 PROOF OF PROPOSITION G.1

Proof of Proposition G.1. We analyze the error layer by layer in the neural network. Denote the outputs of each layer in the networks parameterized by θ and $\tilde{\theta}$ as $X^{(t)}$ and $\tilde{X}^{(t)}$, respectively. In the final layer, we have that

$$\begin{aligned} & \text{TV}(P_\theta(\cdot | X), P_{\tilde{\theta}}(\cdot | X)) \\ & \leq 2 \left\| \frac{1}{L\tau} \mathbb{I}_L^\top X^{(D)} A^{(D+1)} - \frac{1}{L\tau} \mathbb{I}_L^\top \tilde{X}^{(D)} \tilde{A}^{(D+1)} \right\|_\infty \\ & \leq \frac{2}{\tau} \left[\|A^{(D+1),\top}\|_{1,2} \cdot \|X^{(D),\top} - \tilde{X}^{(D),\top}\|_{2,\infty} + \|A^{(D+1),\top} - \tilde{A}^{(D+1),\top}\|_{1,2} \right], \end{aligned}$$

where the first inequality results from Lemma J.6, and the second inequality results from Lemma J.7 and that $\|X^{(D),\top}\|_{2,\infty} \leq 1$ due to the layer normalization. In the following, we build the recursion relationship between $\|X^{(t),\top} - \tilde{X}^{(t),\top}\|_{2,\infty}$ for $t \in [D]$.

$$\begin{aligned} & \|X^{(t+1),\top} - \tilde{X}^{(t+1),\top}\|_{2,\infty} \\ & \leq \|\text{fhn}(Y^{(t+1)}, A^{(t+1)})^\top - \text{fhn}(\tilde{Y}^{(t+1)}, \tilde{A}^{(t+1)})^\top\|_{2,\infty} + |\gamma_2^{(t+1)} - \tilde{\gamma}_2^{(t+1)}| + \|Y^{(t+1),\top} - \tilde{Y}^{(t+1),\top}\|_{2,\infty} \\ & \leq |\gamma_2^{(t+1)} - \tilde{\gamma}_2^{(t+1)}| + \|Y^{(t+1),\top} - \tilde{Y}^{(t+1),\top}\|_{2,\infty} + B_{A,1} \cdot B_{A,2} \cdot \|Y^{(t+1),\top} - \tilde{Y}^{(t+1),\top}\|_{2,\infty} \\ & \quad + B_{A,2} \cdot \|A_1^{(t+1)} - \tilde{A}_1^{(t+1)}\|_F + B_{A,1} \cdot \|A_2^{(t+1)} - \tilde{A}_2^{(t+1)}\|_F, \end{aligned} \quad (\text{I.7})$$

where the first inequality results from the triangle inequality and that Π_{norm} is not expansive, the second inequality results from the following proposition

Proposition I.1. For any $X, \tilde{X} \in \mathbb{R}^{L \times d}$, $A_1, \tilde{A}_1 \in \mathbb{R}^{d \times d_F}$, and $A_2, \tilde{A}_2 \in \mathbb{R}^{d_F \times d}$, we have that

$$\begin{aligned} & \|\text{fhn}(X, A)^\top - \text{fhn}(\tilde{X}, \tilde{A})^\top\|_{2,\infty} \\ & \leq \|A_1\|_F \cdot \|A_2\|_F \cdot \|X^\top - \tilde{X}^\top\|_{2,\infty} + \|A_1 - \tilde{A}_1\|_F \cdot \|A_2\|_F \cdot \|\tilde{X}^\top\|_{2,\infty} \\ & \quad + \|\tilde{A}_1\|_F \cdot \|A_2 - \tilde{A}_2\|_F \cdot \|\tilde{X}^\top\|_{2,\infty}. \end{aligned}$$

Proof of Proposition I.1. See Appendix I.5. □

Next, we build the relationship between $\|Y^{(t+1),\top} - \tilde{Y}^{(t+1),\top}\|_{2,\infty}$ in the right-hand side of inequality (I.7) and $\|X^{(t),\top} - \tilde{X}^{(t),\top}\|_{2,\infty}$.

$$\begin{aligned} & \|Y^{(t+1),\top} - \tilde{Y}^{(t+1),\top}\|_{2,\infty} \\ & \leq \|\text{mha}(X^{(t)}, W^{(t+1)})^\top - \text{mha}(\tilde{X}^{(t)}, \tilde{W}^{(t+1)})^\top\|_{2,\infty} + |\gamma_1^{(t+1)} - \tilde{\gamma}_1^{(t+1)}| + \|X^{(t),\top} - \tilde{X}^{(t),\top}\|_{2,\infty} \\ & \leq |\gamma_1^{(t+1)} - \tilde{\gamma}_1^{(t+1)}| + \|X^{(t),\top} - \tilde{X}^{(t),\top}\|_{2,\infty} \\ & \quad + h \cdot B_V (1 + 4B_Q B_K) \|X^{(t),\top} - \tilde{X}^{(t),\top}\|_{2,\infty} + \sum_{i=1}^h \|W_i^{V,(t+1)} - \tilde{W}_i^{V,(t+1)}\|_F \\ & \quad + 2B_V \cdot B_K \sum_{i=1}^h \|W_i^{Q,(t+1)} - \tilde{W}_i^{Q,(t+1)}\|_F + 2B_V \cdot B_Q \sum_{i=1}^h \|W_i^{K,(t+1)} - \tilde{W}_i^{K,(t+1)}\|_F, \end{aligned} \quad (\text{I.8})$$

where the first inequality results from the triangle inequality, and the second inequality results from Lemma J.8. Combining inequalities (I.7) and (I.8), we derive that

$$\begin{aligned} & \|X^{(t+1),\top} - \tilde{X}^{(t+1),\top}\|_{2,\infty} \\ & \leq (1 + B_{A,1} \cdot B_{A,2}) (1 + hB_V (1 + 4B_Q B_K)) \|X^{(t),\top} - \tilde{X}^{(t),\top}\|_{2,\infty} + \beta_{t+1} + \iota_{t+1} + \kappa_{t+1} + \rho_{t+1}. \end{aligned}$$

This concludes the proof of Proposition G.1. □

I.4 PROOF OF PROPOSITION G.9

Proof of Proposition G.9. We analyze the error layer by layer in the neural network. Denote the outputs of each layer in the networks parameterized by θ and $\tilde{\theta}$ as $X^{(t)}$ and $\tilde{X}^{(t)}$, respectively. In the final layer, we have that

$$\begin{aligned} & \|f_\theta(X) - f_{\tilde{\theta}}(X)\|_2 \\ & \leq \|\tilde{A}^{(D+1)}\|_F \cdot \|X^{(D),\top} - \tilde{X}^{(D),\top}\|_{2,\infty} + \|A^{(D+1)} - \tilde{A}^{(D+1)}\|_F, \end{aligned}$$

where the inequality results from Lemma J.7 and that $\|X^{(D),\top}\|_{2,\infty} \leq 1$ due to the layer normalization. The remaining proof just follows the procedures in the proof of Proposition G.1, and we have that

$$\begin{aligned} & \|f_\theta(X) - f_{\tilde{\theta}}(X)\|_2 \\ & \leq \|A^{(D+1)} - \tilde{A}^{(D+1)}\|_F + \sum_{t=1}^D \alpha_t (\beta_t + \iota_t + \kappa_t + \rho_t). \end{aligned}$$

Thus, we conclude the proof of Proposition G.9. \square

I.5 PROOF OF PROPOSITION I.1

Proof of Proposition I.1. We have that

$$\begin{aligned} & \|\mathbf{f}\mathbf{f}\mathbf{n}(X, A)^\top - \mathbf{f}\mathbf{f}\mathbf{n}(\tilde{X}, \tilde{A})^\top\|_{2,\infty} \\ & \leq \max_{i \in [L]} \left[\|\text{ReLU}(X_{i,:}A_1)A_2 - \text{ReLU}(\tilde{X}_{i,:}A_1)A_2\|_2 + \|\text{ReLU}(\tilde{X}_{i,:}A_1)A_2 - \text{ReLU}(\tilde{X}_{i,:}\tilde{A}_1)\tilde{A}_2\|_2 \right] \\ & \leq \max_{i \in [L]} \left[\|A_1\|_F \cdot \|A_2\|_F \cdot \|X_{i,:} - \tilde{X}_{i,:}\|_2 + \|\text{ReLU}(\tilde{X}_{i,:}A_1)A_2 - \text{ReLU}(\tilde{X}_{i,:}\tilde{A}_1)\tilde{A}_2\|_2 \right. \\ & \quad \left. + \|\text{ReLU}(\tilde{X}_{i,:}\tilde{A}_1)\tilde{A}_2 - \text{ReLU}(\tilde{X}_{i,:}\tilde{A}_1)\tilde{A}_2\|_2 \right] \\ & \leq \max_{i \in [L]} \left[\|A_1\|_F \cdot \|A_2\|_F \cdot \|X_{i,:} - \tilde{X}_{i,:}\|_2 + \|A_1 - \tilde{A}_1\|_F \cdot \|A_2\|_F \cdot \|\tilde{X}_{i,:}\|_2 \right. \\ & \quad \left. + \|\tilde{A}_1\|_F \cdot \|A_2 - \tilde{A}_2\|_F \cdot \|\tilde{X}_{i,:}\|_2 \right], \end{aligned}$$

where the first inequality results from the triangle inequality, the second and the last inequalities result from Lemma J.7 and that ReLU is not expansive. Thus, we conclude the proof of Proposition I.1. \square

J TECHNICAL LEMMAS

Lemma J.1 (Caponnetto and De Vito (2007)). Let (Ω, ν) be a probability space and ξ be a random variable on Ω taking value in a real separable Hilbert space \mathcal{H} . We assume that there exists constants $B, \sigma > 0$ such that

$$\|\xi(w)\|_{\mathcal{H}} \leq B/2, \text{ a.s., } \mathbb{E}[\|\xi\|_{\mathcal{H}}^2] \leq \sigma^2.$$

Then, it holds with probability at least $1 - \delta$ that

$$\left\| L^{-1} \sum_{i=1}^L \xi(\omega_i) - \mathbb{E}[\xi] \right\| \leq 2 \left(\frac{B}{L} + \frac{\sigma}{\sqrt{L}} \right) \log \frac{2}{\delta}.$$

Lemma J.2 (Proposition 4.5 in Duchi (2019)). Let \mathcal{F} be the collection of functions of $f : \mathbb{R}^n \rightarrow \mathbb{R}$. For any $f \in \mathcal{F}$, we define

$$\mu(f) = \mathbb{E}_X[f(X)], \quad \sigma^2(f) = \mathbb{E}_X[(f(X) - \mathbb{E}_X[f(X)])^2],$$

where the expectation is taken with respect to a random variable $X \sim \nu$ on $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$. Assume that $|f(X) - \mu(f)| \leq b$ a.s. for some constant $b \in \mathbb{R}$ for all $f \in \mathcal{F}$. Then for any $0 < \lambda \leq 1/(2b)$, given a distribution P_0 on \mathcal{F} , with probability at least $1 - \delta$, we have

$$\left| \mathbb{E}_Q \left[\mathbb{E}_X[f(X)] - \frac{1}{n} \sum_{i=1}^n f(X_i) \right] \right| \leq \lambda \mathbb{E}_Q[\sigma^2(f)] + \frac{1}{n\lambda} \left[\text{KL}(Q \| P_0) + \log \frac{2}{\delta} \right],$$

for any distribution Q on \mathcal{F} , where X_i are i.i.d. samples of ν . If the function class \mathcal{F} further satisfies $\sigma^2(f) \leq c\mu(f)$ for some constant $c \in \mathbb{R}$ for all $f \in \mathcal{F}$, we have

$$\left| \mathbb{E}_Q \left[\mathbb{E}_X[f(X)] - \frac{1}{n} \sum_{i=1}^n f(X_i) \right] \right| \leq \lambda c \mathbb{E}_Q[\mu(f)] + \frac{1}{n\lambda} \left[\text{KL}(Q \| P_0) + \log \frac{2}{\delta} \right],$$

with probability at least $1 - \delta$.

Lemma J.3 (Donsker–Varadhan representation in [Belghazi et al. \(2018\)](#)). Let P and Q be distributions on a common space \mathcal{X} . Then

$$\text{KL}(P \| Q) = \sup_{g \in \mathcal{G}} \left\{ \mathbb{E}_P[g(X)] - \log \mathbb{E}_Q[\exp(g(X))] \right\},$$

where $\mathcal{G} = \{g : \mathcal{X} \rightarrow \mathbb{R} \mid \mathbb{E}_Q[\exp(g(X))] < \infty\}$.

Lemma J.4 (Corollary 2.11 in [Paulin \(2015\)](#)). Let $X = (X_1, \dots, X_N)$ be a Markov chain, taking values in $\Lambda = \prod_{i=1}^N \Lambda_i$ with mixing time $t_{\text{mix}}(\varepsilon)$ for $\varepsilon \in [0, 1]$. Let

$$t_{\min} = \inf_{0 \leq \varepsilon < 1} t_{\text{mix}}(\varepsilon) \cdot \left(\frac{2 - \varepsilon}{1 - \varepsilon} \right)^2.$$

If function $f : \Lambda \rightarrow \mathbb{R}$ is such that $f(x) - f(y) \leq \sum_{i=1}^N c_i \mathbb{I}_{x_i \neq y_i}$ for every $x, y \in \Lambda$, then for any $\lambda \in \mathbb{R}$,

$$\log \mathbb{E} \left(\exp [\lambda(f(X) - \mathbb{E}f(X))] \right) \leq \frac{\lambda^2 \cdot \|c\|_2^2 \cdot t_{\min}}{8}.$$

For any $t \geq 0$, we have

$$P \left(|f(X) - \mathbb{E}f(X)| \geq t \right) \leq 2 \exp \left(\frac{-2t^2}{\|c\|_2^2 \cdot t_{\min}} \right).$$

Lemma J.5 (Lemma 25 in [Agarwal et al. \(2020\)](#)). For any two conditional probability densities $P(\cdot | X), P'(\cdot | X)$ and any distribution $\nu \in \Delta(\mathcal{X})$, we have

$$\mathbb{E}_\nu \left[\text{TV}(P(\cdot | X), P'(\cdot | X))^2 \right] \leq -2 \log \left(\mathbb{E}_{X \sim \nu, Y \sim P(\cdot | X)} \left[\exp \left(-\frac{1}{2} \log \frac{P(Y | X)}{P'(Y | X)} \right) \right] \right).$$

Lemma J.6 (Corollary A.7 in [Edelman et al. \(2021\)](#)). For any $x, y \in \mathbb{R}^d$, we have

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

Lemma J.7 (Lemma 17 in [Zhang et al. \(2022a\)](#)). Given any two conjugate numbers $u, v \in [1, \infty]$, i.e., $\frac{1}{u} + \frac{1}{v} = 1$, and $1 \leq p \leq \infty$, for any $A \in \mathbb{R}^{r \times c}$ and $x \in \mathbb{R}^c$, we have

$$\|Ax\|_p \leq \|A\|_{p,u} \|x\|_v \quad \text{and} \quad \|Ax\|_p \leq \|A^\top\|_{u,p} \|x\|_v.$$

Lemma J.8 (Propositions 20 and 21 in [Zhang et al. \(2022a\)](#)). For any $X, \tilde{X} \in \mathbb{R}^{L \times d}$, and any $W_i^Q, \tilde{W}_i^Q, W_i^K, \tilde{W}_i^K \in \mathbb{R}^{d \times d_h}, W_i^V, \tilde{W}_i^V \in \mathbb{R}^{d \times d}$ for $i \in [h]$, if $\|X^\top\|_{p,\infty}, \|\tilde{X}^\top\|_{2,\infty} \leq B_X$, $\|W_i^Q\|_F, \|\tilde{W}_i^Q\|_F \leq B_Q, \|W_i^K\|_F, \|\tilde{W}_i^K\|_F \leq B_K, \|W_i^V\|_F, \|\tilde{W}_i^V\|_F \leq B_V$ for $i \in [h]$, then we have

$$\begin{aligned} & \left\| (\text{mha}(X, W) - \text{mha}(\tilde{X}, \tilde{W}))^\top \right\|_{2,\infty} \\ & \leq h \cdot B_V (1 + 4B_X^2 \cdot B_Q B_K) \|X^\top - \tilde{X}^\top\|_{2,\infty} + B_X \sum_{i=1}^h \|W_i^V - \tilde{W}_i^V\|_F \\ & \quad + 2B_X^3 \cdot B_V \cdot B_K \sum_{i=1}^h \|W_i^Q - \tilde{W}_i^Q\|_F + 2B_X^3 \cdot B_V \cdot B_Q \sum_{i=1}^h \|W_i^K - \tilde{W}_i^K\|_F. \end{aligned}$$

Lemma J.9 (Lemma A.6 in [Elbrächter et al. \(2021\)](#)). For $a, b \in \mathbb{R}$ with $a < b$, let

$$\mathcal{S}_{[a,b]} = \left\{ f \in \mathcal{S}^\infty([a,b], \mathbb{R}) \mid \|f^{(n)}(x)\| \leq n! \text{ for all } n \in \mathbb{N} \right\}.$$

There exists a constant $C > 0$ such that for all $a, b \in \mathbb{R}$ with $a < b$, $f \in \mathcal{S}_{[a,b]}$, and $\varepsilon \in (0, 1/2)$, there is a fully connect network Ψ_f such that

$$\|f - \Psi_f\|_\infty \leq \varepsilon,$$

with the depth of the network as $D(\Psi_f) \leq C \max\{2, b - a\}(\log \varepsilon^{-1})^2 + \log(\lceil \max\{|a|, |b|\} \rceil) + \log(\lceil 1/(b - a) \rceil)$, the width of the network as $W(\Psi_f) \leq 16$, and the maximal weight in the network as $B(\Psi_f) \leq 1$.

Lemma J.10. Let $b = \sup_x \log(p(x)/q(x))$. We have that

$$\text{KL}(p \parallel q) \leq 2(3 + b) \cdot \text{TV}(p, q). \quad (\text{J.1})$$

Proof. We let $f(t) = \log t$ and $g(t) = |1/t - 1|$. Then, for $0 \leq t \leq \exp(b)$, we have that

$$\sup_{0 \leq t \leq \exp(b)} \frac{f(t)}{g(t)} = \sup_{0 \leq t \leq \exp(b)} \frac{\log t}{|1/t - 1|} = \sup_{1 \leq t \leq \exp(b)} \frac{t \log t}{t - 1} \leq 2(b + 3).$$

Note that $\text{KL}(p \parallel q) = \mathbb{E}_p[f(p(x)/q(x))]$ and $\text{TV}(p, q) = \mathbb{E}_p[g(p(x)/q(x))]$, which concludes the proof. \square