

# TOWARDS INFINITE-LONG PREFIX IN TRANSFORMER

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Prompting and context-based fine-tuning methods, which we call Prefix Learning, have been proposed to enhance the performance of language models on various downstream tasks. They are empirically efficient and effective, matching the performance of full parameter fine-tuning, but the theoretical understandings are limited. In this paper, we aim to address this limitation by studying their ability from the perspective of prefix length. In particular, we provide a convergence guarantee for training an ultra-long prefix in a stylized setting using the Neural Tangent Kernel (NTK) framework. Based on this strong theoretical guarantee, we design and implement an algorithm that only needs to introduce and fine-tune a few extra trainable parameters instead of an infinite-long prefix in each layer of a transformer, and can approximate the prefix attention to a guaranteed polynomial-small error. Preliminary experimental results on vision, natural language, and math data show that our method achieves superior or competitive performance compared to existing methods like full parameters fine-tuning, P-Tuning V2, and LoRA. This demonstrates our method is promising for parameter-efficient fine-tuning.

## 1 INTRODUCTION

The advent of Large Language Models (LLMs) and Vision LLMs (vLLMs) has significantly advanced the field of Artificial Intelligence (AI), with prominent examples like ChatGPT (ChatGPT, 2022), GPT-4 (Achiam et al., 2023; Bubeck et al., 2023), Claude (Claude-3, 2024), Llama (Touvron et al., 2023a;b), Gemini (Gemini, 2024), ViT (Dosovitskiy et al., 2020), DETR (Carion et al., 2020), BLIP (Li et al., 2022; 2023a), CLIP (Radford et al., 2021). They have exhibited impressive performances across a spectrum of tasks, encompassing chat systems (Maaz et al., 2023; Xu et al., 2023a; Zheng et al., 2024), text-to-image conversion (Qiao et al., 2019; Frolov et al., 2021; Zhang et al., 2023), AI mathematical inference (Hendrycks et al., 2020; Yu et al., 2023a; Yao et al., 2023), and many more. However, despite these advancements, pre-existing LLMs often fall short in specialized domains that demand a deeper understanding of professional knowledge (Tajbakhsh et al., 2016; Devlin et al., 2018; Gururangan et al., 2020; Hu et al., 2021; Sun, 2023; Kasneci et al., 2023; Li et al., 2023b; Thirunavukarasu et al., 2023; Li et al., 2024b; Wang et al., 2024). This has led to the development of fine-tuning/adaptation (Shi et al., 2022; Xu et al., 2023b; Shi et al., 2024a) methodologies aimed at enhancing the proficiency of these models in executing more specialized tasks (Mangrulkar et al., 2022). Several notable contributions in this area, such as LoRA (Low-Rank Adaptation, Hu et al. (2021)), P-Tuning (Liu et al., 2021b; 2023), and  $(IA)^3$  (Liu et al., 2022), have displayed performances rivaling those of full-parameter fine-tuning techniques. This underscores the potential of these fine-tuning strategies to further refine the capabilities of Large Language Models.

Among the methods proposed, most context-based fine-tuning methods, e.g., Prompt-Tuning (Lester et al., 2021; Liu et al., 2021a), Prefix-Tuning (Li & Liang, 2021), P-Tuning (Liu et al., 2023; 2021b), use enhanced input sequences (or virtual prompt, a.k.a soft prompt) to optimize their model outputs. These methods are gaining significant interest due to their ease of implementation across various model architectures, and also prevention of catastrophic forgetting with static pre-trained parameters (Wang et al., 2023b; Sohn et al., 2023; Yang et al., 2024). We call the above approaches **Prefix Learning** since they improve the performance by optimizing a prefix matrix added to the input in each attention layer of the LLMs (see detailed formulation in Section 2).

Despite its wide use and strong empirical performance, we still have a limited understanding of why and how prefix learning operates (Wang et al., 2023a; Petrov et al., 2024a;b). One common phenomenon in prior empirical studies is that prefix learning results in better downstream performance

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

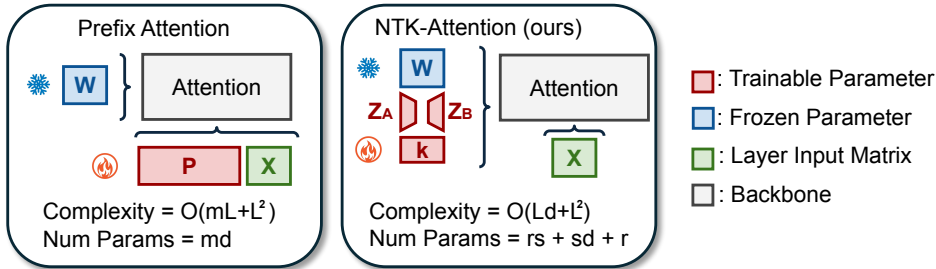


Figure 1: Illustration of existing prefix attention methods (Algorithm 1) and our NTK-Attention (Algorithm 2). Compared to the former, NTK-Attention significantly reduces the number of parameters and the time complexity. Here,  $X \in \mathbb{R}^{L \times d}$  is the input of this layer,  $W = [W_Q, W_K, W_V]$  is frozen weights of attention,  $P \in \mathbb{R}^{m \times d}$  is the trainable prefix matrix and  $Z_A \in \mathbb{R}^{r \times s}$ ,  $Z_B \in \mathbb{R}^{s \times d}$ ,  $k \in \mathbb{R}^r$  are the trainable parameters in our method.  $L$  is the input length,  $d$  the input dimension,  $m$  the prefix length, and  $r$  a hyperparameter in NTK-attention (i.e., the dimension of the constructed feature mapping; see Section 4). Note that  $m \gg L$  and  $m \gg d$ , and  $r = \text{poly}(d)$  (usually be chosen to  $d$  or  $2d$ ),  $s \leq \lfloor d/2 \rfloor$  (low-rank of  $Z_A, Z_B$ ) are used in our experiments.

when the prefix length increases (Lester et al., 2021; Liu et al., 2023). We call this phenomenon *scaling law in prefix learning*: the longer the prefix, the larger downstream dataset the model can fit, and thus the better performance the model would have. Then intuitively, we would like to ask:

*What happens when the prefix length is large or even tends to infinity?*

The answer to this cannot be directly figured out via empirical evaluations, since it is impractical to implement networks with ultra-long or even infinite prefixes in practice. Therefore, we first perform a theoretical analysis of prefix learning. We study the optimization of ultra-long prefix learning via the Neural Tangent Kernel (NTK) technique (Jacot et al., 2018), which has been used for analyzing overparameterized networks and thus is suitable for ultra-long prefix learning. Based on the insights gained from the analysis, we propose our method, NTK-attention, which reparameterizes prefix learning and can approximate infinite-long prefix learning using a finite number of parameters. We also conduct some empirical evaluations of our method on vision, natural language understanding, and math inference datasets to demonstrate its effectiveness.

Specifically, we have made the following contributions:

- We first perform a theoretical analysis of optimizing an ultra-long prefix in a stylized attention network; see Section 3. We consider a simplified attention network, and show that when prefix length  $m$  is sufficiently large (i.e., prefix learning is sufficiently over-parameterized), the training can be analyzed via NTK, which leads to our theoretical guarantee of convergence to small errors. This also provides theoretical support for scaling law in prefix learning.
- We then propose our NTK-Attention (Algorithm 2), motivated by the above strong theoretical guarantee; see Section 4. Our method approximates existing prefix attention (Algorithm 1) by utilizing three trainable parameters  $Z_A, Z_B$  and  $k$ , to replace the parameter in prefix attention (the prefix matrix  $P$ ). This allows scaling the prefix length without large memory usage and computational time that increases with the prefix length. It reduces the computation complexity from  $O(mL)$  to  $O(L^2)$ , where  $L$  is the input length and  $m$  is the prefix length. See Figure 1 for an illustration.
- We further conduct experiments on vision, language and math datasets to verify our theoretical results; see Section 5. The experiments include (1) a comparison among our NTK-Attention, full parameters fine-tuning, and LoRA on CIFAR-100, Food-101 and Tiny-Imagenet datasets with the same pretrained ViT backbone; (2) a comparison among our NTK-Attention, P-Tuning V2, and LoRA on SuperGLUE, WikiText-103, Penn TreeBank and LAMBADA datasets with the same pretrained ChatGLM3-6B and OPT- $\{125M, 350M, 1.3B, 2.7B, 6.7B\}$  family; (3) a comparison among our NTK-Attention and LoRA on GSM8K and MATH datasets with supervised fine-tune pretrained models LLAMA-3.2; (4) an ablation study to validate sensitivity of hyper-parameters in NTK-Attention; (5) a comparison of the computational costs between our method and standard prefix learning on random data. The empirical results show that on average our NTK-Attention

method achieves better performance than the competitors. For example, on SuperGLUE datasets, it achieves an average accuracy that is 1.07% higher than LoRA and 12.94% higher than P-Tuning V2. It is also observed that our method maintains low time and memory costs while those of prefix learning scales with prefix length. The experimental results demonstrate that our method is effective and efficient and supports our theoretical analysis.

## 1.1 RELATED WORK

**Prefix Learning.** Prefix Learning (Lester et al., 2021; Ding et al., 2021; Wang et al., 2022b; Zhou et al., 2022; Liu et al., 2021a; Petrov et al., 2024a; Wu et al., 2023), including Prompt-Tuning (Lester et al., 2021), Prefix-Tuning (Li & Liang, 2021), P-Tuning (Liu et al., 2023; 2021b), Reweighted In-Context Learning (RICL) (Chu et al., 2023) and so on, is proposed to enhance the performance of language models on the downstream tasks and to reduce the costs of computational resources of fine-tuning the whole model. Those methods optimize task-specific prompts for downstream task improvement. On the other hand, besides the Parameter-Efficient-Fine-Tuning (PEFT) approaches (Mangrulkar et al., 2022) we mentioned above, Retrieval Augmented Generation (RAG) (Lewis et al., 2020; Jiang et al., 2023; Gao et al., 2023b) and Chain-of-Thought (CoT) prompting (Wei et al., 2022b; Wang et al., 2022a; Fu et al., 2022) can also be considered as prefix learning. We conclude all these works to an optimization problem that improves the prefix based on task-specific measurements.

**Neural Tangent Kernel.** Neural Tangent Kernel (NTK) (Jacot et al., 2018) studies the gradient flow of neural networks in the training process. They showed neural networks are equivalent to Gaussian processes in the infinite-width limit at initialization. A bunch of works has explained the strong performance and the learning ability of neural networks at over-parameterization, such as (Li & Liang, 2018; Du et al., 2019; Song & Yang, 2019; Allen-Zhu et al., 2019; Wei et al., 2019; Bietti & Mairal, 2019; Lee et al., 2020; Chizat & Bach, 2020; Shi et al., 2021; Zhou et al., 2021; Seleznova & Kutyniok, 2022; Gao et al., 2023a; Li et al., 2024a; Shi et al., 2024c) and many more. Furthermore, Arora et al. (2019) gave the first exact algorithm on computing Convolutional NTK (CNTK), Alemohammad et al. (2020) proposed Recurrent NTK, and Hron et al. (2020) presented infinite attention via NNGP and NTK for attention networks. These works have demonstrated advanced performance by utilizing NTK in different neural network architectures. In particular, Malladi et al. (2023) have studied the training dynamic of fine-tuning LLMs via NTK and confirmed the efficiency of such methods.

**Theory of Understanding Large Language Models.** Since the complicated transformer-based architecture and stochastic optimization process of LLMs lead the study of their behaviors to be a challenge, analyzing LLMs through some theoretical guarantee helps in providing insights to improve and design the next generation of AI systems. This topic includes efficient LLMs (Alman & Song, 2023; 2024a;b; Han et al., 2024; Kacham et al., 2023; Addanki et al., 2023; Deng et al., 2024; Shi et al., 2024b), optimization of LLMs (Deng et al., 2023; Li et al., 2024a), white-box transformers (Yu et al., 2023b;c; Ferrando et al., 2024; Pai et al., 2024), analysis of emergent abilities of LLMs (Brown et al., 2020; Wei et al., 2022a; Allen-Zhu & Li, 2023a;b;c; 2024), etc. Especially, (Alman & Song, 2023) proved that the hardness of fast attention can be achieved within  $n^{1+o(1)}$  times executions, one effective way is to construct a high-order polynomial mapping based on Taylor expansion of the exponential function  $\exp(\cdot)$ , and it inspired the design of our NTK-Attention method.

## 2 PRELIMINARY: PREFIX LEARNING

In this section, we provide the detailed formulation for prefix learning, which optimizes prefix matrices in the attention layers of transformer-based LLMs. Focusing on one single-layer attention network, we formalize it as a regression problem that optimizes a prefix matrix.

**Prefix for Attention Computation.** Let  $X \in \mathbb{R}^{L \times d}$  be an input matrix to the attention network, where  $L$  and  $d$  are the input length and dimension. Prefix learning freezes the query, key, and value parameter matrices in the pretrained attention network (denoted as  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ , respectively). It introduces a trainable prefix matrix  $P \in \mathbb{R}^{m \times d}$ , which stands for  $m$  virtual token vectors (or soft prompt). Let  $S := \begin{bmatrix} P \\ X \end{bmatrix}$  be the concatenation of the prefix and the input. Then the query, key, and value matrices are given by  $Q := XW_Q, K_P := SW_K, V_P := SW_V$ , and the

attention with the prefix is:

$$\text{PrefixAttn}(X, P) := \text{Softmax}\left(\frac{QK_P^\top}{\sqrt{d}}\right) \cdot V_P \in \mathbb{R}^{L \times d}. \quad (1)$$

Here  $\text{Softmax}$  is the row-wise softmax computation, i.e., for any  $d_1, d_2 > 0$ ,  $Z \in \mathbb{R}^{d_1 \times d_2}$ ,  $\text{Softmax}(Z) := [\text{S}(Z_{1,*}), \text{S}(Z_{2,*}), \dots, \text{S}(Z_{d_1,*})]^\top \in \mathbb{R}^{d_1 \times d_2}$  where  $\text{S}(z) := \frac{\exp(z)}{\langle \exp(z), \mathbf{1}_{d_2} \rangle} \in \mathbb{R}^{d_2}$  for any  $z \in \mathbb{R}^{d_2}$ . The attention computation with prefix is summarized in Algorithm 1.

**Prefix Learning.** The prefix  $P$  is trained on a fine-tuning dataset. Denote the dataset as  $\mathcal{D}_{\text{pl}} = \{(X_i, Y_i)\}_{i=1}^n$  where  $n$  is the dataset size, and  $X_i, Y_i \in \mathbb{R}^{L \times d}$ . Let  $\ell(\cdot, \cdot)$  denote the loss function for the specific task (e.g., prompting, context-based fine-tuning, etc). The training objective of prefix learning is then:

$$\min_{P \in \mathbb{R}^{m \times d}} \mathcal{L}_{\text{pl}}(W) := \sum_{i=1}^n \ell(\text{PrefixAttn}(X_i, P), Y_i). \quad (2)$$

**Scaling Prefix Length.** A rich line of studies (Liu et al., 2021b; Lester et al., 2021; Liu et al., 2023; Reynolds & McDonell, 2021; Arora et al., 2022; Brown et al., 2020; Dong et al., 2022; Shi et al., 2023; Von Oswald et al., 2023; Xu et al., 2024; Fu et al., 2022; Agarwal et al., 2024; Kaplan et al., 2020; Hoffmann et al., 2022) have reported a common observation that as the prefix length increases, the model’s ability to master complex skills also improves. Specifically, the performance of fine-tuned models is enhanced when the prefix length grows within a certain range. A similar trend is observed in prompting methods and in-context learning (ICL), where longer and more complex prompts lead to better inference abilities in LLMs, and providing more examples in ICL results in improved LLM performance. We summarize this as the *scaling law in prefix learning*: the longer the prefix length for fine-tuning, the larger dataset the model can fit, thus, the more complicated skill it can master. This motivates investigating prefix learning with long prefixes.

In this paper, we examine the implications of using a significantly large prefix length, denoted as  $m \gg L$  and  $m \gg d$ , which is prevalent across various prompt-based methods. The primary objective of Prefix Learning is to enhance the LLMs’ outputs by identifying an advanced prefix during the generation process. For instance, the search for optimal example pairs to improve ICL (Nguyen & Wong, 2023) and the development of prompt engineering tailored for agent frameworks to address specific task requirements (dif, 2024) often necessitate the use of exceptionally long prefixes. Moreover, given the modern application demands related to long-context scenarios, optimizing previous tokens to improve next-token prediction can be framed as a prefix optimization problem. Thus, a thorough investigation into the optimization of infinitely long prefixes is essential for understanding the theoretical significance of the prefix matrix in LLMs.

### 3 THEORETICAL ANALYSIS OF PREFIX LEARNING VIA NTK

In this section, we explore the theory behind prefix learning with ultra-long prefixes. We first present the theoretical setting for a simplified model  $F(W, x, a)$  in Section 3.1, and then in Section 3.2 introduce the formal definition of the neural tangent kernel for our problem and confirm the convergence of the kernel matrices needed for performing NTK analysis. In Section 3.3 we state the main result, a convergence guarantee of prefix learning in this setting (the detailed analysis is in the appendix).

#### 3.1 PROBLEM SETUP

**Model.** The attention computation with prefix  $P$  given is by Eq. (1). Since the attention parameters are fixed, it can be rewritten as  $\text{Softmax}(\tilde{X}P^\top + b) \cdot \begin{bmatrix} PW_V \\ b' \end{bmatrix}$  where  $\tilde{X} = XW_QW_K^\top/\sqrt{d}$ ,  $b = XW_QW_K^\top X^\top/\sqrt{d}$ , and  $b' = XW_V$ . We view the input sequence as one token (i.e., assuming  $L = 1$ ) such that the input  $X$  and thus  $\tilde{X}$  become vectors, simplifying our analysis from matrix-form calculations to vector-form. Furthermore, ignoring the bias terms, and introducing notations  $x := \tilde{X}^\top$  and  $W = P^\top$ , the attention simplifies to  $\text{Softmax}(xW) \cdot W^\top W_V = \frac{\sum_{r \in [m]} \exp(w_r^\top x) w_r W_V}{\sum_{r \in [m]} \exp(w_r^\top x)}$  where

$w_r$  is the  $r$ -th column of  $W$ . We therefore consider the following two-layer attention model:

$$F(W, x, a) := m \frac{\sum_{r \in [m]} \exp(w_r^\top x) w_r a_r}{\sum_{r \in [m]} \exp(w_r^\top x)} \quad (3)$$

with the hidden-layer weights  $W = [w_1, w_2, \dots, w_m] \in \mathbb{R}^{d \times m}$  and output-layer weights  $a = [a_1, a_2, \dots, a_m]^\top \in \mathbb{R}^m$ . Such a stylized setting has been widely used for studying the learning behavior of transformer-based models (Deng et al., 2023; Chu et al., 2023; 2024; Li et al., 2024a), and they gave detailed derivations and guarantees for its connection to attention. Furthermore, our analysis can be extended to models with bias terms and matrix inputs rigorously.

**Training.** Consider a training dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  where the  $i$ -th data point  $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}^d$ . Assume  $\|x_i\|_2 \leq 1$  and  $\|y_i\|_2 \leq 1$  for any  $i \in [n]$ . The training loss is measured by the  $\ell_2$  norm of the difference between model prediction  $F(W, x_i, a)$  and ideal output vector  $y_i$ . Formally, the training objective is:

$$\mathcal{L}(W) := \frac{1}{2} \sum_{i=1}^n \|F(W, x_i, a) - y_i\|_2^2. \quad (4)$$

The weights  $W$  are initialized to  $W(0)$  as follows:  $\forall r \in [m]$ , sample  $w_r(0) \sim \mathcal{N}(0, I_d)$  independently. For output-layer  $a$ , randomly sample  $a_r \sim \text{Uniform}\{-1, +1\}$  independently for  $r \in [m]$  and fix  $a$  during the training. Then use gradient descent (GD) to update the trainable weights  $W(t)$  with a fixed learning rate  $\eta > 0$ . Then for  $t \geq 0$ :

$$W(t+1) := W(t) - \eta \cdot \nabla_W \mathcal{L}(W(t)). \quad (5)$$

### 3.2 NEURAL TANGENT KERNEL

Here, we give the formal definition of NTK in our analysis, which is a kernel function that is driven by hidden-layer weights  $W(t) \in \mathbb{R}^{d \times m}$ . To present concisely, we first introduce an operator function in the following. For all  $r \in [m]$ ,  $k \in [d]$  and  $i \in [n]$ :

$$v_{k,r}(W) := W_{k,r} \cdot a_r \cdot \mathbf{1}_m - W_{k,*} \circ a \in \mathbb{R}^m, \quad \mathcal{G}_{i,r}(W) := m S_r(W^\top x_i) \cdot \langle v_{k,r}, S(W^\top x_i) \rangle \in \mathbb{R}$$

where  $S(z) = \frac{\exp(z)}{\langle \exp(z), \mathbf{1}_m \rangle} \in \mathbb{R}^m$  for any  $z \in \mathbb{R}^m$ , and  $\circ$  denotes element-wise product.

Then, we define the kernel matrix  $H(W(t))$  as an  $nd \times nd$  Gram matrix, where its  $(k_1, k_2)$ -th block is an  $n \times n$  matrix for  $k_1, k_2 \in [d]$ , and the  $(i, j)$ -th entry of the block is:

$$[H_{k_1, k_2}]_{i,j}(W(t)) := \frac{1}{m} x_i^\top x_j \sum_{r=1}^m \mathcal{G}_{i,r}(W(t)) \cdot \mathcal{G}_{j,r}(W(t)).$$

We can show that  $S_r(W^\top x_i) = O(\frac{1}{m})$  and  $\langle v_{k,r}, S(W^\top x_i) \rangle = O(1)$ , thus  $\mathcal{G}_{i,r}(W)$  is  $O(1)$ . Then  $H(W)$  is close to  $H^* := H(W(0))$  when  $W$  is close to  $W(0)$ . This kernel convergence is the key needed for the NTK analysis and is formalized below (details in Appendix H).

**Lemma 3.1** (Kernel convergence, informal version of Lemma H.3). *For  $\delta \in (0, 0.1)$  and  $B = \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ . Let  $\tilde{W} = [\tilde{w}_1, \dots, \tilde{w}_m] \in \mathbb{R}^{d \times m}$  and satisfy  $\|\tilde{w}_r - w_r(0)\|_2 \leq R$  for any  $r \in [m]$ , where  $R$  is some constant in  $(0, 0.01)$ . Define  $\tilde{H} := H(\tilde{W}) \in \mathbb{R}^{nd \times nd}$ . Then with probability at least  $1 - \delta$ , we have  $\|H^* - \tilde{H}\| \leq 8R\sqrt{nd} \cdot \exp(22B)$ .*

### 3.3 MAIN RESULT: LOSS CONVERGENCE GUARANTEE

**Assumption on NTK  $H^*$ .** In the NTK analysis framework for the convergence of training neural networks, one widely-used and mild assumption is that  $H^*$  is a positive definite (PD) matrix, i.e., its minimum eigenvalue  $\lambda := \lambda_{\min}(H^*) > 0$  (Du et al., 2019; Oymak & Soltanolkotabi, 2020). With this, our main result is presented as follows.

**Theorem 3.2** (Main result, informal version of Theorem J.2). *Assume  $\lambda > 0$ . For any  $\epsilon, \delta \in (0, 0.1)$ ,  $B = \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ ,  $m = \lambda^{-2} \text{poly}(n, d, \exp(B))$ ,  $\eta = \lambda m^{-1} / \text{poly}(n, d, \exp(B))$  and  $\hat{T} = \Omega((m\eta\lambda)^{-1} \log(nd/\epsilon))$ . Then, after  $\hat{T}$  iterations of update (Eq. (5)), we have  $\mathcal{L}(W(\hat{T})) \leq \epsilon$  holds with probability at least  $1 - \delta$ .*

*Proof sketch of Theorem 3.2.* We use the math induction to show that the weight  $w$  perturbation is small so that the loss landscape is almost convex around the network’s initialization in Lemma J.3, Lemma J.4 and Lemma J.5, which are based on Lemma 3.1. Then, we conclude the results by standard convex optimization analysis. See the complete proof in Appendix J.1.  $\square$

**Discussion.** Theorem 3.2 mainly describes the following fact for any dataset with  $n$  data points. After initializing the prefix matrix from a normal distribution, assuming the minimum eigenvalue of NTK  $\lambda > 0$ , setting  $m$  to be a large enough value so that the network is sufficiently over-parameterized. Then with proper learning rate, the loss can be minimized in finite training time to an arbitrarily small error  $\epsilon$ . Corresponding to the real-world implementation, it explains that adequately long prefix learning can master downstream tasks when fine-tuning LLMs. Furthermore, it also helps us understand the working mechanism of prefix learning, inspiring us to explore the direction of using ultra-long prefixes.

Now we connect our theory to the *scaling law in prefix learning*. Following (Kaplan et al., 2020), we focus on the relationship between the loss and the computational cost. We prove that the loss decreases with the computational cost scaling up, providing a theoretical confirmation about the scaling law in prefix learning.

**Proposition 3.3** (Scaling Law in Prefix Learning). *We define  $N := O(md)$  as the number of parameters,  $D := O(n)$  as the size of training dataset,  $C_{\text{cpt}} := O(NDT)$  as the total compute cost, and  $\alpha := nd$ . We choose  $T$  as Theorem 3.2, then the loss of training, denotes  $L$ , satisfies:*

$$L \approx \frac{\alpha}{[\exp(\eta\lambda C_{\text{cpt}})]^{\frac{1}{\alpha}}}$$

*Proof sketch of Proposition 3.3.* This proof follows from the definitions of  $C_{\text{cpt}}$ ,  $N$ ,  $D$  and  $\alpha$  and Theorem 3.2.  $\square$

Proposition 3.3 shows that the training loss of the prefix learning converges exponentially as we increase the computational cost  $C_{\text{cpt}}$ , which primarily depends on the number of parameters and the training time in prefix learning, further indicating a possible relationship for formulating scaling law in prefix learning.

## 4 NTK-ATTENTION: APPROXIMATE INFINITE-LONG PREFIX ATTENTION

The preceding section discussed the convergence guarantee of training sufficiently long prefixes  $P$  in attention networks (recall that the trainable parameter  $W$  is just  $P^\top$ ). This strong theoretical property inspires us to scale up the prefix length  $m$ . However, such prefix learning (Algorithm 1) necessitates a time complexity of  $O(mLd + L^2d)$  in each layer of the model, this is impractical due to a large  $m$ .

This section proposes an approximate algorithm to make long prefix learning practical. Our algorithm, NTK-Attention, is designed to output an approximation of  $\text{PrefixAttn}(X, P)$  (Eq. (1)) in time within  $O(L^{1+o(1)})$  and without using the long prefix matrix  $P$ . We present the derivation and motivation of our algorithm in Section 4.1, formalize the NTK-Attention algorithm in Section 4.2, and provide an approximation guarantee in Section 4.3.

### 4.1 DERIVATION: REPLACING PREFIX $P$ WITH TRAINABLE PARAMETERS $Z, k$

There exists a wealth of attention approximation algorithms capable of executing attention computations within  $n^{1+o(1)}$  time (Han et al., 2024; Liang et al., 2024a;b). However, our focus lies predominantly with the polynomial method (Tsai et al., 2019; Katharopoulos et al., 2020; Alman & Song, 2023; 2024b). This method has exhibited exceptional performance in terms of both time and space complexity through the use of a streaming algorithm.

**Polynomial method.** In the context of attention networks, the query, key, and value state matrices, denoted as  $Q, K, V \in \mathbb{R}^{L \times d}$ , are assumed to have all entries bounded (Alman & Song, 2023). Under this condition, the polynomial method first constructs a linear mapping  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^r$ , where

$r = \text{poly}(d)$  (Alman & Song, 2023), and it satisfies the following relation ( $i, j \in [L]$ ,  $Q_i, K_j \in \mathbb{R}^d$  represent the  $i$ -th row of  $Q$  and the  $j$ -th row of  $K$  respectively):

$$\phi(Q_i)^\top \phi(K_j) \approx \exp(Q_i^\top K_j / \sqrt{d}). \quad (6)$$

Here, the mapping  $\phi(\cdot)$  is constructed based on the Taylor expansion of the exponential function, and the larger value of  $r \geq d$  would bring the approximation (Eq. (6)) with a smaller error. This is guaranteed by Lemma 3.4 in Alman & Song (2023), refer to a copy in Lemma K.7. The  $i$ -th row of the approximate attention (denoted as  $\text{PolyAttn}_i \in \mathbb{R}^{1 \times d}$ ) then can be computed as follows:

$$\text{PolyAttn}_i := \frac{\phi(Q_i)^\top \sum_{j=1}^L \phi(K_j) V_j^\top}{\phi(Q_i)^\top \sum_{j=1}^L \phi(K_j)} \in \mathbb{R}^{1 \times d}, \forall i \in [L].$$

Now recall that given an input matrix  $X \in \mathbb{R}^{L \times d}$ , thus,  $Q = XW_Q$ , and we have  $[K_P, V_P] = \begin{bmatrix} P \\ X \end{bmatrix}$ .

$[W_K, W_V] = \begin{bmatrix} PW_K & PW_V \\ XW_K & XW_V \end{bmatrix}$ . Let  $K_C := PW_K, V_C := PW_V \in \mathbb{R}^{m \times d}$  and  $K := XW_K, V := XW_V \in \mathbb{R}^{L \times d}$ . We thus expand the  $i$ -th row of the prefix attention,  $\text{PrefixAttn}_i(X, P) \in \mathbb{R}^{1 \times d}$  as:

$$\begin{aligned} \text{PrefixAttn}_i(X, P) &= \frac{\exp(Q_i^\top K^\top / \sqrt{d})V + \exp(Q_i^\top K_C^\top / \sqrt{d})V_C}{\exp(Q_i^\top K^\top / \sqrt{d})\mathbf{1}_L + \exp(Q_i^\top K_C^\top / \sqrt{d})\mathbf{1}_m} \\ &\approx \frac{\exp(Q_i^\top K^\top / \sqrt{d})V + \phi(Q_i)^\top Z}{\exp(Q_i^\top K^\top / \sqrt{d})\mathbf{1}_n + \phi(Q_i)^\top k} \end{aligned}$$

where

$$Z = \sum_{j=1}^m \phi(K_{C,j}) V_{C,j}^\top \in \mathbb{R}^{r \times d}, \quad k = \sum_{j=1}^m \phi(K_{C,j}) \in \mathbb{R}^r. \quad (7)$$

Here, the first step explicitly computes the softmax function, and the second step holds since replacing  $\exp(Q_i^\top K^\top / \sqrt{d})$  by Eq. (6), which is  $\exp(Q_i^\top K_{C,j}^\top / \sqrt{d}) \approx \phi(Q_i)^\top \phi(K_{C,j}), \forall j \in [m]$ .

Therefore, checking the training process of  $P$ , we observe that  $P$  is updating iff  $Z$  and  $k$  are updating. Hence, we can replace  $P$  by utilizing **trainable parameters**  $Z$  and  $k$  in Eq. (7) to re-parameterize the prefix attention. This is the key to how NTK-Attention approximates prefix attention without a large number of parameters.

## 4.2 ALGORITHM

To present our algorithm, based on  $\phi$ , we define:  $\Phi(A) = [\phi(A_{1,*}), \dots, \phi(A_{L,*})]^\top \in \mathbb{R}^{L \times r}, \forall A \in \mathbb{R}^{L \times d}$ . Below we present our NTK-Attention method in Algorithm 2, and for comparison also present the traditional prefix attention for prefix learning in Algorithm 1.

**Implementation Detail of  $\phi$ .** In order to find a balance between approximation and efficient computation of NTK-Attention, we use the first-order polynomial method. In particular, we choose  $r = d$ , and the function  $\phi$  is given by  $\phi(z) := d^{-\frac{1}{4}} \cdot (z \circ \mathbf{1}_{z \geq 0_d} + \exp(z) \circ \mathbf{1}_{z < 0_d}) + \mathbf{1}_d \in \mathbb{R}^d, \forall z \in \mathbb{R}^d$ , where  $\mathbf{1}_{z \geq 0_d} \in \mathbb{R}^d$  is an indicative vector and its  $i$ -th entry for  $i \in [d]$  equals 1 only when  $z_i \geq 0$ , and 0 otherwise.

**Initialization, Approximation and Training of  $Z$  and  $k$ .** In Section 3.1, we initialize the parameter  $W = P^\top$  by  $w_r(0) \sim \mathcal{N}(0, I_d)$  for  $r \in [m]$ . Since the pretrained weights  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$  are known, the initialization of  $Z$  and  $k$ , denotes  $Z(0)$  and  $k(0)$ , can then be computed by Eq. (7) using  $P(0) = W(0)^\top$ . However, consider that  $Z$  caches  $rd$  parameters for  $r = \text{poly}(d)$ , which is insufficient parameter-efficient. In response to it, we choose  $s \leq \lfloor d/2 \rfloor$  as an appropriately small integer, then  $Z(0) \approx Z_A(0) \cdot Z_B(0)$  is decomposed into two low-rank matrices  $Z_A(0) \in \mathbb{R}^{r \times s}, Z_B(0) \in \mathbb{R}^{s \times d}$ . For training, let  $g_{Z_A}(t) \in \mathbb{R}^{r \times s}, g_{Z_B}(t) \in \mathbb{R}^{s \times d}$  and  $g_k(t) \in \mathbb{R}^r$  denote the gradients of  $Z_A(t), Z_B(t)$  and  $k(t)$  at time  $t$ , and  $\eta$  denote the learning rate. Then the update rule is:

$$Z_A(t+1) := Z_A(t) - \eta \cdot g_{Z_A}(t), Z_B(t+1) := Z_B(t) - \eta \cdot g_{Z_B}(t), k(t+1) := k(t) - \eta \cdot g_k(t).$$

**Number of Trainable Parameters.** Since given  $r$  and  $s$  as two hyper-parameters in NTK-Attention, for each attention layer in transformer-based architecture, we denote  $\beta := \frac{r}{d}$ , then the number of

Table 1: Performance of different fine-tuning methods on the SuperGLUE datasets. The base model is ChatGLM3-6B. The methods include P-Tuning V2, LoRA, and our NTK-Attention method. The metric on these datasets is accuracy (measured in %). The best score on each dataset is **boldfaced**.

Method	Num Params	Task					Average
		BoolQ	CB	Copa	MultiRC	RTE	
P-Tuning V2 $m = 1$	0.12M	65.69±0.32	67.06±0.37	52.00±1.00	53.59±0.28	65.97±0.22	60.86±0.44
P-Tuning V2 $m = 10$	1.15M	66.67±0.23	74.07±0.00	54.00±0.00	54.17±0.71	66.55±0.25	63.10±0.24
P-Tuning V2 $m = 100$	11.47M	69.42±0.02	74.54±0.47	64.50±0.50	61.62±2.28	76.77±0.83	69.37±0.82
P-Tuning V2 $m = 200$	22.94M	67.51±0.15	70.11±0.28	60.00±0.50	58.37±0.91	70.83±0.44	65.36±0.46
LoRA $r' = 8$	3.67M	<b>76.52</b> ±0.10	90.23±0.39	86.50±0.50	65.09±0.41	<b>87.76</b> ±0.37	81.24±0.35
NTK-Attention (ours), $r = 128, s = 16$	3.78M	75.06±0.12	<b>96.04</b> ±0.84	<b>88.00</b> ±2.00	<b>65.85</b> ±0.33	86.59±0.52	<b>82.31</b> ±0.76

trainable parameters could be computed by  $(\beta s + \beta + s)d$  where integer  $\beta \geq 1$  and  $s \leq \lfloor d/2 \rfloor$ . This is more flexible when adjusting the practical efficiency needs. For LoRA with its hyper-parameter  $r' \leq \lfloor d/2 \rfloor$ , where  $r'$  is the rank number used for approximation, its number of trainable parameters is  $4r'd$  and for prefix attention with its hyper-parameter  $m \geq 1$ , its number of trainable parameters is  $md$  in each attention layer. By choosing  $(\beta s + \beta + s) \leq 4r'$ , the higher efficiency of NTK-Attention compared to LoRA will be satisfied.

#### Algorithm 1 Prefix Attention

**Input:** Input matrix  $X \in \mathbb{R}^{L \times d}$   
**Parameters:** Frozen query, key and value weights  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ , trainable prefix matrix  $P \in \mathbb{R}^{m \times d}$   
**Output:** Exact output Attn  $\in \mathbb{R}^{L \times d}$

- 1: **procedure** PREFIXATTEN( $X$ )
- 2:  $S \leftarrow [P^\top, X^\top]^\top$
- 3:  $Q, K_P, V_P \leftarrow XW_Q, SW_K, SW_V$
- 4:  $A \leftarrow \exp(QK_P^\top / \sqrt{d})$
- 5:  $D \leftarrow \text{diag}(A\mathbf{1}_{m+L})$
- 6: **return**  $D^{-1}AV_P$
- 7: **end procedure**

#### Algorithm 2 NTK-Attention (w/o low-rank)

**Input:** Input matrix  $X \in \mathbb{R}^{L \times d}$   
**Parameters:** Frozen query, key and value weights  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ , trainable weights  $Z \in \mathbb{R}^{r \times d}$  and  $k \in \mathbb{R}^r$   
**Output:** Approx output  $T \in \mathbb{R}^{L \times d}$

- 1: **procedure** NTK-ATTEN( $X$ )
- 2:  $Q, K, V \leftarrow XW_Q, XW_K, XW_V$ ,
- 3:  $\hat{A} \leftarrow \exp(QK^\top / \sqrt{d})$
- 4:  $\hat{D} \leftarrow \text{diag}(\hat{A}\mathbf{1}_L + \Phi(Q)k)$
- 5:  $T \leftarrow \hat{D}^{-1}(\hat{A}V + \Phi(Q)Z)$
- 6: **return**  $T$
- 7: **end procedure**

### 4.3 ERROR BOUND AND COMPLEXITY REDUCTION

Introducing an ultra-long prefix matrix  $P \in \mathbb{R}^{m \times d}$  to satisfy the conditions in Theorem J.2 requires  $md$  parameters for  $m \geq \Omega(\lambda^{-2} \text{poly}(n, d, \exp(B)))$ , while it also bring a  $O(m(m+L)d)$  time complexity to compute Algorithm 1. Our NTK-Attention relieve this by replacing  $P$  with  $Z$  and  $k$ , where we state our theoretical guarantee as follows:

**Theorem 4.1** (Error bound with reduced time complexity, informal version of Theorem K.2). *Let  $m$  denote the prefix length. Given an input matrix  $X \in \mathbb{R}^{L \times d}$  and prefix matrix  $P \in \mathbb{R}^{m \times d}$ , we denote  $Q = XW_Q$ ,  $K_C = PW_K$  and  $V_C = PW_V$ . If the condition Eq. (7),  $\|Q\|_\infty \leq o(\sqrt{\log m})$ ,  $\|K_C\|_\infty \leq o(\sqrt{\log m})$ ,  $\|V_C\|_\infty \leq o(\sqrt{\log m})$  and  $d = O(\log m)$  holds, then Algorithm 2 outputs a matrix  $T \in \mathbb{R}^{L \times d}$  within time complexity of  $O(L^2d)$  that satisfies:*

$$\|T - \text{PrefixAttn}(X, P)\|_\infty \leq 1 / \text{poly}(m). \quad (8)$$

Furthermore, if we replace the original attention operation (attention computation on input  $X$  with  $K = XW_K$  and  $V = XW_V$ ) with fast attention algorithms like HyperAttention (Han et al., 2024), then NTK-Attention can be even more efficient, achieving Eq. (8) within complexity  $O(L^{1+o(1)}d)$  (see Corollary K.3 for proofs).

## 5 EMPIRICAL EVALUATIONS

In this section, we evaluate our method NTK-Attention on natural language understanding, math inference, and fine-grained image classification tasks. All our experiments use the Huggingface (Wolf



et al., 2019) trainer with AdamW optimizer (Kingma & Ba, 2014), and all optimizer hyper-parameters are set to the defaults. We provide more details in Appendix B.

**Evaluation on Natural Language Understanding Datasets.** In this experiment, we utilize five binary classification datasets in SuperGLUE (Wang et al., 2019) for evaluation: the BoolQ, CB, Copa, MultiRC, and RTE datasets. We use a pretrained LLM ChatGLM3-6B (Zeng et al., 2022; Du et al., 2022) as the base model. For comparison, we choose P-Tuning V2 (Liu et al., 2023; 2021b) which is a standard prefix learning method, and choose LoRA (Hu et al., 2021) which is a popular parameter-efficient fine-tuning method often achieving state-of-the-art. P-Tuning V2 uses different lengths of virtual prefix  $\{1, 10, 100, 200\}$ , and LoRA uses rank  $r' = 8$ . We choose  $r = 128$  (the dimension of each head of ChatGLM3-6B) and  $s = 16$  for our NTK-Attention.

The results are provided in Table 1. Our NTK-Attention method achieves much higher performance than P-Tuning V2. Interestingly, as  $m$  increases, the performance of P-Tuning V2 also improves, which is consistent with our analysis. Our analysis also suggests that NTK-Attention approximates ultra-long prefix learning and thus can perform better than P-Tuning V2. The experimental results also show that NTK-Attention achieves better performance than LoRA on CB, Copa, and MultiRC datasets, and achieves better average performance over all the datasets. These results show that NTK-Attention can be a promising efficient fine-tuning method.

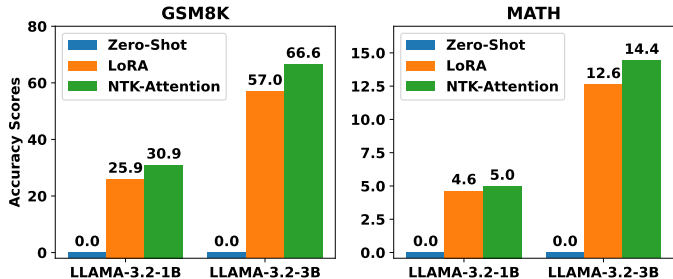


Figure 2: Compare our results with LoRA and Zero-Shot on Math inference datasets. The  $y$ -axis is the accuracy.

**Evaluation on Language Modeling Tasks.** In this experiment, we focus on the scalability of NTK-Attention on a family of language models of different sizes, the OPT family with the model sizes 125M, 350M, 1.3B, 2.7B and 6.7B (Zhang et al., 2022). We introduce three text datasets, which are WikiText-103 (Merity et al., 2016), Penn TreeBank (Marcus et al., 1993), and LAMBADA (Paperno et al., 2016), to compare the scalability of NTK-Attention with LoRA (Hu et al., 2021) and P-Tuning V2 (Liu et al., 2023; 2021b). As we choose  $r' = 8$  for LoRA,  $m = 32$  for P-Tuning V2, and  $r = 2d$  and  $s = 10$  for our NTK-Attention, the numbers of trainable parameters are aligned to the same as  $32d$  for each attention layer. The results are stated in Table 3, which shows the improvement of NTK-Attention compared to baselines when scaling the model size.

**Evaluation on Math Inference Datasets.** In order to thoroughly verify the effectiveness of NTK-Attention, we conduct experiments on the math inference task, which includes GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021) datasets. These are considered as fair benchmarks to test the complex capability of LLMs. We follow Yu et al. (2023a) to supervised fine-tune two pretrained models LLAMA-3.2-1B and LLAMA-3.2-3B (Touvron et al., 2023a;b) with dataset MetaMathQA (Yu et al., 2023a). We state our results in Figure 2, and we use accuracy scores for counting the matched answers for evaluation. As we can see, our NTK-attention ( $r = d, s = 16$ ) is better than the two baselines, LoRA and Zero-Shot, where LoRA uses  $r' = 16$  for LLAMA-3.2-1B and  $r' = 32$  for LLAMA-3.2-3B.

**Evaluation on Vision Datasets.** We evaluate the method on three image classification datasets: CIFAR-100 (Krizhevsky et al., 2009), Food-101 (Bossard et al., 2014), and Tiny-Imagenet (mn-moustafa, 2017). The base model to be fine-tuned on these datasets is ViT-Base (Dosovitskiy et al., 2020) that is pretrained on the ImageNet-21k (Deng et al., 2009). We compare our method to two baselines: (1) FFT (Full parameters Fine-Tuned) that fine-tunes all parameters; (2) LoRA that fine-tunes the base model with the popular LoRA method (Hu et al., 2021) with rank  $r' = \{16, 32\}$ .

The results are presented in Table 2. Our method performs much better than FFT: 7.40%, 5.81% and 13.26% higher accuracy on the three datasets, respectively. Note that FFT updates all parameters and has much higher computational costs than LoRA or our method. Our method has a similar

Table 2: Performance of different fine-tuning methods on the CIFAR-100, Food-101 and Tiny-Imagenet datasets. The base model is ViT-Base. The methods include FFT, LoRA, and our method NTK-Attention. The metric is accuracy (measured in %). The best score on each dataset is **boldfaced**.

Method	Num Params	Dataset			Average
		CIFAR-100	Food-101	Tiny-Imagenet	
FFT	86.39M	85.15±0.13	84.76±0.07	76.20±0.23	82.04±0.14
LoRA $r' = 16$	7.08M	92.17±0.05	89.38±0.33	88.22±0.09	89.92±0.16
LoRA $r' = 32$	14.16M	92.01±0.20	89.86±0.11	<b>90.16±0.12</b>	90.68±0.14
NTK-Attention (ours), $r = 64, s = 32$	7.09M	<b>92.55±0.03</b>	<b>90.57±0.01</b>	89.46±0.10	<b>90.86±0.05</b>

Table 3: Performance of different fine-tuning methods on OPT- $\{125M, 350M, 1.3B, 2.7B, 6.7B\}$  pretrained models with WikiText-103, Penn TreeBank and LAMBADA datasets. The metric is perplexity (PPL), with its smaller value standing for better performance. The best score on each dataset and model is **boldfaced**.

Model	Method	Num Params	Datasets			Average
			WikiText-103	Penn TreeBank	LAMBADA	
OPT-125M	LoRA, $r' = 8$	0.29M	<b>30.50</b>	35.97	46.02	37.50
	P-Tuning V2, $m = 32$		2264.22	963.09	1762.19	1663.17
	NTK-Attention, $r = 2d, s = 10$		31.41	<b>33.52</b>	<b>45.39</b>	<b>36.77</b>
OPT-350M	LoRA, $r' = 8$	0.77M	<b>24.76</b>	30.41	38.80	31.32
	P-Tuning V2, $m = 32$		7383.48	1339.43	14020.36	7581.09
	NTK-Attention, $r = 2d, s = 10$		25.67	<b>28.85</b>	<b>36.97</b>	<b>30.50</b>
OPT-1.3B	LoRA, $r' = 8$	1.57M	<b>16.71</b>	21.27	24.16	20.71
	P-Tuning V2, $m = 32$		2230.76	540.17	3480.77	2083.9
	NTK-Attention, $r = 2d, s = 10$		17.04	20.09	24.04	<b>20.39</b>
OPT-2.7B	LoRA, $r' = 8$	2.62M	15.06	19.61	22.13	18.93
	P-Tuning V2, $m = 32$		772.48	277.99	3378.18	1476.22
	NTK-Attention, $r = 2d, s = 10$		<b>14.83</b>	<b>18.52</b>	<b>21.85</b>	<b>18.40</b>
OPT-6.7B	LoRA, $r' = 8$	4.19M	12.81	17.36	19.38	16.52
	P-Tuning V2, $m = 32$		2051.10	409.37	4709.46	2389.98
	NTK-Attention, $r = 2d, s = 10$		<b>12.56</b>	<b>16.68</b>	<b>18.81</b>	<b>16.02</b>

performance to LoRA with  $r' = 32$ , achieving slightly better average accuracy. These results on vision datasets also provide positive empirical support for our method.

**Ablation Study.** We validate the sensitivity of hyper-parameters  $r$  and  $s$  and give the results in Appendix B.3. The results firstly indicate that choosing  $r = d$  and  $s = 4$  is enough for high-performance fine-tuning on LLAMA-3.1-8B. Also, we follow Table 4 to suggest choosing a larger value of  $r$  primarily instead of  $s$  to achieve supernal accuracy.

**Empirical Evaluation of Computational Cost.** We also provide experimental results of the computational costs of NTK-Attention (Algorithm 2) and the standard Prefix Attention (Algorithm 1) in Appendix B.2. The results show that Prefix Attention’s run time is quadratic and memory usage is linear in the prefix length, so its costs are typically much higher, while NTK-Attention maintains a small run time and memory usage.

## 6 CONCLUSION

In this study, we illuminated the principles of prefix learning for fine-tuning when the prefix length is large. We conducted an in-depth theoretical analysis, demonstrating that when the prefix length is sufficiently large, the attention network is over-parameterized, and the Neural Tangent Kernel technique can be leveraged to provide a convergence guarantee of prefix learning. Based on these insights, we proposed a novel efficient fine-tuning method called NTK-Attention, which approximates prefix attention using two trainable parameters to replace the large prefix matrix, thus significantly mitigating memory usage issues and reducing computational cost for long prefixes. We also provided empirical results to support our theoretical findings, demonstrating NTK-Attention’s superior performance on downstream tasks over baselines across natural language, math, and vision datasets.

## REFERENCES

- 540  
541  
542 The innovation engine for genai applications. <https://github.com/langgenius/dify>,  
543 2024.
- 544 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
545 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.  
546 *arXiv preprint arXiv:2303.08774*, 2023.
- 547 Raghav Addanki, Chenyang Li, Zhao Song, and Chiwun Yang. One pass streaming algorithm for  
548 super long token attention approximation in sublinear space. *arXiv preprint arXiv:2311.14652*,  
549 2023.
- 550 Rishabh Agarwal, Avi Singh, Lei M Zhang, Bernd Bohnet, Stephanie Chan, Ankesh Anand, Zaheer  
551 Abbas, Azade Nova, John D Co-Reyes, Eric Chu, et al. Many-shot in-context learning. *arXiv*  
552 *preprint arXiv:2404.11018*, 2024.
- 553 Sina Alemohammad, Zichao Wang, Randall Balestriero, and Richard Baraniuk. The recurrent neural  
554 tangent kernel. *arXiv preprint arXiv:2006.10246*, 2020.
- 555 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 1, context-free grammar. *arXiv*  
556 *preprint arXiv:2305.13673*, 2023a.
- 557 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.2, knowledge manipulation.  
558 *arXiv preprint arXiv:2309.14402*, 2023b.
- 559 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and  
560 extraction. *arXiv preprint arXiv:2309.14316*, 2023c.
- 561 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity scaling  
562 laws. *arXiv preprint arXiv:2404.05405*, 2024.
- 563 Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-  
564 parameterization. In *International conference on machine learning*, pp. 242–252. PMLR, 2019.
- 565 Josh Alman and Zhao Song. Fast attention requires bounded entries. *Advances in Neural Information*  
566 *Processing Systems*, 36, 2023.
- 567 Josh Alman and Zhao Song. The fine-grained complexity of gradient computation for training large  
568 language models. *arXiv preprint arXiv:2402.04497*, 2024a.
- 569 Josh Alman and Zhao Song. How to capture higher-order correlations? generalizing matrix soft-  
570 max attention to kronecker computation. In *The Twelfth International Conference on Learning*  
571 *Representations*, 2024b.
- 572 Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On  
573 exact computation with an infinitely wide neural net. *Advances in neural information processing*  
574 *systems*, 32, 2019.
- 575 Simran Arora, Avaniika Narayan, Mayee F Chen, Laurel Orr, Neel Guha, Kush Bhatia, Ines Chami,  
576 and Christopher Re. Ask me anything: A simple strategy for prompting language models. In *The*  
577 *Eleventh International Conference on Learning Representations*, 2022.
- 578 Sergei Bernstein. On a modification of chebyshev’s inequality and of the error formula of laplace.  
579 *Ann. Sci. Inst. Sav. Ukraine, Sect. Math*, 1(4):38–49, 1924.
- 580 Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. *Advances in Neural*  
581 *Information Processing Systems*, 32, 2019.
- 582 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative compo-  
583 nents with random forests. In *European Conference on Computer Vision*, 2014.
- 584 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
585 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
586 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

- 594 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,  
595 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:  
596 Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- 597  
598 Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey  
599 Zagoruyko. End-to-end object detection with transformers. In *European conference on computer  
600 vision*, pp. 213–229. Springer, 2020.
- 601 ChatGPT. Optimizing language models for dialogue. *OpenAI Blog*, November 2022. URL [https:  
602 //openai.com/blog/chatgpt/](https://openai.com/blog/chatgpt/).
- 603  
604 Herman Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of  
605 observations. *The Annals of Mathematical Statistics*, pp. 493–507, 1952.
- 606  
607 Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks  
608 trained with the logistic loss. In *Conference on learning theory*, pp. 1305–1338. PMLR, 2020.
- 609  
610 Timothy Chu, Zhao Song, and Chiwun Yang. Fine-tune language models to approximate unbiased  
611 in-context learning. *arXiv preprint arXiv:2310.03331*, 2023.
- 612  
613 Timothy Chu, Zhao Song, and Chiwun Yang. How to protect copyright data in optimization of  
614 large language models? In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp.  
615 17871–17879, 2024.
- 616  
617 Claude-3. Introducing the next generation of claude. *Anthropic News*, March 2024. URL [https:  
618 //www.anthropic.com/news/claude-3-family/](https://www.anthropic.com/news/claude-3-family/).
- 619  
620 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
621 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve  
622 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 623  
624 Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv  
625 preprint arXiv:2307.08691*, 2023.
- 626  
627 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-  
628 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,  
629 35:16344–16359, 2022.
- 630  
631 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
632 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
633 pp. 248–255. Ieee, 2009.
- 634  
635 Yichuan Deng, Zhihang Li, and Zhao Song. Attention scheme inspired softmax regression. *arXiv  
636 preprint arXiv:2304.10411*, 2023.
- 637  
638 Yichuan Deng, Zhao Song, and Chiwun Yang. Attention is naturally sparse with gaussian distributed  
639 input. *arXiv preprint arXiv:2404.02690*, 2024.
- 640  
641 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
642 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 643  
644 Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, Zhiyuan Liu, Hai-Tao Zheng, and Maosong Sun.  
645 Openprompt: An open-source framework for prompt-learning. *arXiv preprint arXiv:2111.01998*,  
646 2021.
- 647  
648 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and  
649 Zhifang Sui. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- 650  
651 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
652 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An  
653 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint  
654 arXiv:2010.11929*, 2020.
- 655  
656 Simon S Du, Xiyu Zhai, Barnabas Póczos, and Aarti Singh. Gradient descent provably optimizes  
657 over-parameterized neural networks. In *ICLR*. *arXiv preprint arXiv:1810.02054*, 2019.

- 648 Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm:  
649 General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th*  
650 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.  
651 320–335, 2022.
- 652 Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. A primer on the inner  
653 workings of transformer-based language models. *arXiv preprint arXiv:2405.00208*, 2024.
- 654 Sergey Foss, Dmitry Korshunov, Stan Zachary, et al. *An introduction to heavy-tailed and subexpo-*  
655 *656 nential distributions*, volume 6. Springer, 2011.
- 657 Stanislav Frolov, Tobias Hinz, Federico Raue, Jörn Hees, and Andreas Dengel. Adversarial text-to-  
658 image synthesis: A review. *Neural Networks*, 144:187–209, 2021.
- 659 Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting  
660 for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*,  
661 2022.
- 662 Yeqi Gao, Sridhar Mahadevan, and Zhao Song. An over-parameterized exponential regression. *arXiv*  
663 *preprint arXiv:2303.16504*, 2023a.
- 664 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and  
665 Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv*  
666 *preprint arXiv:2312.10997*, 2023b.
- 667 Gemini. Welcome to the gemini era. *Google Deepmind Technologies*, May 2024. URL <https://deepmind.google/technologies/gemini/>.
- 668 Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,  
669 and Noah A Smith. Don’t stop pretraining: Adapt language models to domains and tasks. *arXiv*  
670 *preprint arXiv:2004.10964*, 2020.
- 671 Uffe Haagerup. The best constants in the khintchine inequality. *Studia Mathematica*, 70(3):231–283,  
672 1981.
- 673 Insu Han, Rajesh Jayaram, Amin Karbasi, Vahab Mirrokni, David Woodruff, and Amir Zandieh.  
674 Hyperattention: Long-context attention in near-linear time. In *The Twelfth International Confer-*  
675 *676 ence on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Eh0Od2BJIM>.
- 677 David Lee Hanson and Farroll Tim Wright. A bound on tail probabilities for quadratic forms in  
678 independent random variables. *The Annals of Mathematical Statistics*, 42(3):1079–1083, 1971.
- 679 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
680 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
681 *arXiv:2009.03300*, 2020.
- 682 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
683 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,  
684 2021.
- 685 Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The collected*  
686 *687 works of Wassily Hoeffding*, pp. 409–426, 1994.
- 688 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
689 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.  
690 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 691 Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: Nngp and ntk  
692 for deep attention networks. In *International Conference on Machine Learning*, pp. 4376–4386.  
693 PMLR, 2020.
- 694 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,  
695 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*  
696 *arXiv:2106.09685*, 2021.

- 702 Jerry Yao-Chieh Hu, Maojiang Su, En-Jui Kuo, Zhao Song, and Han Liu. Computational limits of  
703 low-rank adaptation (lora) for transformer-based models. *arXiv preprint arXiv:2406.03136*, 2024.  
704
- 705 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and  
706 generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.  
707
- 708 Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang,  
709 Jamie Callan, and Graham Neubig. Active retrieval augmented generation. *arXiv preprint*  
710 *arXiv:2305.06983*, 2023.
- 711 Praneeth Kacham, Vahab Mirrokni, and Peilin Zhong. Polysketchformer: Fast transformers via  
712 sketches for polynomial kernels. *arXiv preprint arXiv:2310.01655*, 2023.
- 713 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott  
714 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.  
715 *arXiv preprint arXiv:2001.08361*, 2020.  
716
- 717 Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank  
718 Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, et al. Chatgpt for good?  
719 on opportunities and challenges of large language models for education. *Learning and individual*  
720 *differences*, 103:102274, 2023.
- 721 Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns:  
722 Fast autoregressive transformers with linear attention. In *International conference on machine*  
723 *learning*, pp. 5156–5165. PMLR, 2020.  
724
- 725 Aleksandr Khintchine. Über dyadische brüche. *Mathematische Zeitschrift*, 18(1):109–116, 1923.  
726
- 727 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
728 *arXiv:1412.6980*, 2014.
- 729 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.  
730
- 731 Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection.  
732 *Annals of statistics*, pp. 1302–1338, 2000.
- 733 Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and  
734 Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Advances in*  
735 *Neural Information Processing Systems*, 33:15156–15172, 2020.  
736
- 737 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt  
738 tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- 739 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
740 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-  
741 tion for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:  
742 9459–9474, 2020.
- 743 Chenyang Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Exploring the frontiers of softmax: Prov-  
744 able optimization, applications in diffusion model, and beyond. *arXiv preprint arXiv:2405.03251*,  
745 2024a.  
746
- 747 Chenyang Li, Yingyu Liang, Zhenmei Shi, Zhao Song, and Tianyi Zhou. Fourier circuits in neural  
748 networks: Unlocking the potential of large language models in mathematical reasoning and modular  
749 arithmetic. *arXiv preprint arXiv:2402.09469*, 2024b.
- 750 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-  
751 training for unified vision-language understanding and generation. In *International conference on*  
752 *machine learning*, pp. 12888–12900. PMLR, 2022.  
753
- 754 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image  
755 pre-training with frozen image encoders and large language models. In *International conference*  
*on machine learning*, pp. 19730–19742. PMLR, 2023a.

- 756 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv*  
757 *preprint arXiv:2101.00190*, 2021.  
758
- 759 Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. Large language models in finance: A survey.  
760 In *Proceedings of the Fourth ACM International Conference on AI in Finance*, pp. 374–382, 2023b.
- 761 Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient  
762 descent on structured data. *Advances in neural information processing systems*, 31, 2018.  
763
- 764 Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to  
765 solve inherently serial problems. *arXiv preprint arXiv:2402.12875*, 2024c.
- 766 Yingyu Liang, Heshan Liu, Zhenmei Shi, Zhao Song, and Junze Yin. Conv-basis: A new  
767 paradigm for efficient attention inference and gradient computation in transformers. *arXiv preprint*  
768 *arXiv:2405.05219*, 2024a.
- 769 Yingyu Liang, Zhenmei Shi, Zhao Song, and Yufa Zhou. Tensor attention training: Provably efficient  
770 learning of higher-order transformers. *arXiv preprint arXiv:2405.16411*, 2024b.  
771
- 772 Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and  
773 Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context  
774 learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.
- 775 Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig.  
776 Pre-train, prompt, and predict: a systematic survey of prompting methods in natural language  
777 processing. *arxiv. arXiv preprint arXiv:2107.13586*, 2021a.  
778
- 779 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang.  
780 P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks.  
781 *arXiv preprint arXiv:2110.07602*, 2021b.
- 782 Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt  
783 understands, too. *AI Open*, 2023.  
784
- 785 Yichao Lu, Paramveer Dhillon, Dean P Foster, and Lyle Ungar. Faster ridge regression via the  
786 subsampled randomized hadamard transform. *Advances in neural information processing systems*,  
787 26, 2013.
- 788 Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt:  
789 Towards detailed video understanding via large vision and language models. *arXiv preprint*  
790 *arXiv:2306.05424*, 2023.
- 791 Sadhika Malladi, Alexander Wettig, Dingli Yu, Danqi Chen, and Sanjeev Arora. A kernel-based  
792 view of language model fine-tuning. In *International Conference on Machine Learning*, pp.  
793 23610–23641. PMLR, 2023.  
794
- 795 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin  
796 Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>, 2022.  
797
- 798 Mitch Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus  
799 of english: The penn treebank. *Computational linguistics*, 19(2):313–330, 1993.  
800
- 801 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
802 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 803 Mohammed Ali mnmostafa. Tiny imagenet, 2017. URL [https://kaggle.com/](https://kaggle.com/competitions/tiny-imagenet)  
804 [competitions/tiny-imagenet](https://kaggle.com/competitions/tiny-imagenet).
- 805 Alexander Munteanu, Simon Omlor, Zhao Song, and David Woodruff. Bounding the width of neural  
806 networks via coupled initialization a worst case analysis. In *International Conference on Machine*  
807 *Learning*, pp. 16083–16122. PMLR, 2022.  
808
- 809 Tai Nguyen and Eric Wong. In-context example selection with influences. *arXiv preprint*  
*arXiv:2302.11042*, 2023.

- 810 Samet Oymak and Mahdi Soltanolkotabi. Toward moderate overparameterization: Global con-  
811 vergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in*  
812 *Information Theory*, 1(1):84–105, 2020.
- 813  
814 Druv Pai, Sam Buchanan, Ziyang Wu, Yaodong Yu, and Yi Ma. Masked completion via structured  
815 diffusion with white-box transformers. In *The Twelfth International Conference on Learning*  
816 *Representations*, 2024.
- 817 Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi,  
818 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset:  
819 Word prediction requiring a broad discourse context. *arXiv preprint arXiv:1606.06031*, 2016.
- 820  
821 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
822 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style,  
823 high-performance deep learning library. *Advances in neural information processing systems*, 32,  
824 2019.
- 825 Aleksandar Petrov, Philip Torr, and Adel Bibi. When do prompting and prefix-tuning work? a  
826 theory of capabilities and limitations. In *The Twelfth International Conference on Learning*  
827 *Representations*, 2024a.
- 828  
829 Aleksandar Petrov, Philip HS Torr, and Adel Bibi. Prompting a pretrained transformer can be a  
830 universal approximator. *arXiv preprint arXiv:2402.14753*, 2024b.
- 831 Tingting Qiao, Jing Zhang, Duanqing Xu, and Dacheng Tao. Mirrorgan: Learning text-to-image  
832 generation by redescription. In *Proceedings of the IEEE/CVF conference on computer vision and*  
833 *pattern recognition*, pp. 1505–1514, 2019.
- 834  
835 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
836 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
837 models from natural language supervision. In *International conference on machine learning*, pp.  
838 8748–8763. PMLR, 2021.
- 839 Laria Reynolds and Kyle McDonell. Prompt programming for large language models: Beyond the  
840 few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in*  
841 *Computing Systems*, pp. 1–7, 2021.
- 842  
843 Mark Rudelson and Roman Vershynin. Hanson-wright inequality and sub-gaussian concentration.  
844 2013.
- 845 Mariia Seleznova and Gitta Kutyniok. Neural tangent kernel beyond the infinite-width limit: Effects  
846 of depth and initialization. In *International Conference on Machine Learning*, pp. 19522–19560.  
847 PMLR, 2022.
- 848 Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao.  
849 Flashattention-3: Fast and accurate attention with asynchrony and low-precision. *arXiv preprint*  
850 *arXiv:2407.08608*, 2024.
- 851  
852 Zhenmei Shi, Junyi Wei, and Yingyu Liang. A theoretical analysis on feature learning in neural  
853 networks: Emergence from inputs and advantage over fixed features. In *International Conference*  
854 *on Learning Representations*, 2021.
- 855  
856 Zhenmei Shi, Jiefeng Chen, Kunyang Li, Jayaram Raghuram, Xi Wu, Yingyu Liang, and Somesh Jha.  
857 The trade-off between universality and label efficiency of representations from contrastive learning.  
858 In *The Eleventh International Conference on Learning Representations*, 2022.
- 859  
860 Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu Liang. Why larger language models do in-context  
861 learning differently? In *R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large*  
862 *Foundation Models*, 2023.
- 863  
864 Zhenmei Shi, Yifei Ming, Ying Fan, Frederic Sala, and Yingyu Liang. Domain generalization via  
865 nuclear norm regularization. In *Conference on Parsimony and Learning*, pp. 179–201. PMLR,  
866 2024a.



- 864 Zhenmei Shi, Yifei Ming, Xuan-Phi Nguyen, Yingyu Liang, and Shafiq Joty. Discovering the gems  
865 in early layers: Accelerating long-context llms with 1000x input token reduction. *arXiv preprint*  
866 *arXiv:2409.17422*, 2024b.
- 867 Zhenmei Shi, Junyi Wei, and Yingyu Liang. Provable guarantees for neural networks via gradient  
868 feature learning. *Advances in Neural Information Processing Systems*, 36, 2024c.
- 870 Kihyuk Sohn, Huiwen Chang, José Lezama, Luisa Polania, Han Zhang, Yuan Hao, Irfan Essa, and  
871 Lu Jiang. Visual prompt tuning for generative transfer learning. In *Proceedings of the IEEE/CVF*  
872 *Conference on Computer Vision and Pattern Recognition*, pp. 19840–19851, 2023.
- 874 Zhao Song and Xin Yang. Quadratic suffices for over-parametrization via matrix chernoff bound.  
875 *arXiv preprint arXiv:1906.03593*, 2019.
- 876 Zhongxiang Sun. A short survey of viewing large language models in legal aspect. *arXiv preprint*  
877 *arXiv:2303.09136*, 2023.
- 879 Nima Tajbakhsh, Jae Y Shin, Suryakanth R Gurudu, R Todd Hurst, Christopher B Kendall, Michael B  
880 Gotway, and Jianming Liang. Convolutional neural networks for medical image analysis: Full  
881 training or fine tuning? *IEEE transactions on medical imaging*, 35(5):1299–1312, 2016.
- 882 Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang  
883 Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*, 29(8):  
884 1930–1940, 2023.
- 886 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
887 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
888 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 889 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
890 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
891 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 893 Joel A Tropp. Improved analysis of the subsampled randomized hadamard transform. *Advances in*  
894 *Adaptive Data Analysis*, 3(01n02):115–126, 2011.
- 895 Yao-Hung Hubert Tsai, Shaojie Bai, Makoto Yamada, Louis-Philippe Morency, and Ruslan Salakhut-  
896 dinov. Transformer dissection: a unified understanding of transformer’s attention via the lens of  
897 kernel. *arXiv preprint arXiv:1908.11775*, 2019.
- 899 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,  
900 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In  
901 *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.
- 902 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer  
903 Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language  
904 understanding systems. *Advances in neural information processing systems*, 32, 2019.
- 906 Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, and Neel Joshi. Is a picture worth  
907 a thousand words? delving into spatial reasoning for vision language models. *arXiv preprint*  
908 *arXiv:2406.14852*, 2024.
- 909 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-  
910 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.  
911 *arXiv preprint arXiv:2203.11171*, 2022a.
- 913 Yihan Wang, Jatin Chauhan, Wei Wang, and Cho-Jui Hsieh. Universality and limitations of prompt  
914 tuning. *Advances in Neural Information Processing Systems*, 36, 2023a.
- 915 Zhen Wang, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, Huan Sun, and Yoon Kim. Multitask  
916 prompt tuning enables parameter-efficient transfer learning. *arXiv preprint arXiv:2303.02861*,  
917 2023b.

- 918 Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent  
919 Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings*  
920 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 139–149, 2022b.
- 921  
922 Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and  
923 optimization of neural nets vs their induced kernel. *Advances in Neural Information Processing*  
924 *Systems*, 32, 2019.
- 925 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,  
926 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.  
927 *arXiv preprint arXiv:2206.07682*, 2022a.
- 928 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
929 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
930 *neural information processing systems*, 35:24824–24837, 2022b.
- 931  
932 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
933 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers:  
934 State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- 935 Junda Wu, Tong Yu, Rui Wang, Zhao Song, Ruiyi Zhang, Handong Zhao, Chaochao Lu, Shuai Li,  
936 and Ricardo Henao. Infoprompt: Information-theoretic soft prompt tuning for natural language  
937 understanding. *Advances in Neural Information Processing Systems*, 36, 2023.
- 938  
939 Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with  
940 parameter-efficient tuning on self-chat data. *arXiv preprint arXiv:2304.01196*, 2023a.
- 941 Zhuoyan Xu, Zhenmei Shi, Junyi Wei, Fangzhou Mu, Yin Li, and Yingyu Liang. Towards few-shot  
942 adaptation of foundation models via multitask finetuning. In *The Twelfth International Conference*  
943 *on Learning Representations*, 2023b.
- 944  
945 Zhuoyan Xu, Zhenmei Shi, and Yingyu Liang. Do large language models have compositional ability?  
946 an investigation into limitations and scalability. In *ICLR 2024 Workshop on Mathematical and*  
947 *Empirical Understanding of Foundation Models*, 2024.
- 948 Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Shaochen  
949 Zhong, Bing Yin, and Xia Hu. Harnessing the power of llms in practice: A survey on chatgpt and  
950 beyond. *ACM Transactions on Knowledge Discovery from Data*, 18(6):1–32, 2024.
- 951  
952 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.  
953 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural*  
954 *Information Processing Systems*, 36, 2023.
- 955 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo  
956 Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for  
957 large language models. *arXiv preprint arXiv:2309.12284*, 2023a.
- 958 Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin  
959 Haeffele, and Yi Ma. White-box transformers via sparse rate reduction. *Advances in Neural*  
960 *Information Processing Systems*, 36:9422–9457, 2023b.
- 961  
962 Yaodong Yu, Tianzhe Chu, Shengbang Tong, Ziyang Wu, Druv Pai, Sam Buchanan, and Yi Ma. Emer-  
963 gence of segmentation with minimalistic white-box transformers. *arXiv preprint arXiv:2308.16271*,  
964 2023c.
- 965 Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,  
966 Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint*  
967 *arXiv:2210.02414*, 2022.
- 968  
969 Yuchen Zeng and Kangwook Lee. The expressive power of low-rank adaptation. *arXiv preprint*  
970 *arXiv:2310.17513*, 2023.
- 971  
972 Chenshuang Zhang, Chaoning Zhang, Mengchun Zhang, and In So Kweon. Text-to-image diffusion  
973 model in generative ai: A survey. *arXiv preprint arXiv:2303.07909*, 2023.

972 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher  
973 Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language  
974 models. *arXiv preprint arXiv:2205.01068*, 2022.

975 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
976 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
977 chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024.

978 Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for  
979 vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and  
980 pattern recognition*, pp. 16816–16825, 2022.

981 Mo Zhou, Rong Ge, and Chi Jin. A local convergence theory for mildly over-parameterized two-layer  
982 neural network. In *Conference on Learning Theory*, pp. 4577–4632. PMLR, 2021.

983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

# Appendix

1026		
1027		
1028		
1029	CONTENTS	
1030		
1031	<b>1 Introduction</b>	<b>1</b>
1032	1.1 Related Work . . . . .	3
1033		
1034	<b>2 Preliminary: Prefix Learning</b>	<b>3</b>
1035		
1036	<b>3 Theoretical Analysis of Prefix Learning via NTK</b>	<b>4</b>
1037	3.1 Problem Setup . . . . .	4
1038	3.2 Neural Tangent Kernel . . . . .	5
1039	3.3 Main Result: Loss Convergence Guarantee . . . . .	5
1040		
1041	<b>4 NTK-Attention: Approximate Infinite-Long Prefix Attention</b>	<b>6</b>
1042	4.1 Derivation: Replacing Prefix $P$ with Trainable Parameters $Z, k$ . . . . .	6
1043	4.2 Algorithm . . . . .	7
1044	4.3 Error Bound and Complexity Reduction . . . . .	8
1045		
1046	<b>5 Empirical Evaluations</b>	<b>8</b>
1047		
1048	<b>6 Conclusion</b>	<b>10</b>
1049		
1050	<b>A Algorithm Details and Computational Complexity Analysis</b>	<b>22</b>
1051		
1052	<b>B Experimental Details</b>	<b>22</b>
1053	B.1 Setup Details . . . . .	22
1054	B.2 Additional Empirical Complexity Analysis . . . . .	23
1055	B.3 Additional Ablation Study . . . . .	24
1056		
1057	<b>C Naive NTK-Attention Implementation with Flash-Attention</b>	<b>24</b>
1058		
1059	<b>D Further Discussions</b>	<b>25</b>
1060		
1061	<b>E Preliminary of Analysis</b>	<b>26</b>
1062	E.1 Facts . . . . .	26
1063	E.2 Probability . . . . .	26
1064		
1065	<b>F Definitions of NTK Analysis</b>	<b>28</b>
1066	F.1 Loss function . . . . .	29
1067		
1068	<b>G Gradient Computation</b>	<b>29</b>
1069	G.1 Computing Gradient . . . . .	29
1070	G.2 Gradient Descent . . . . .	31
1071		
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079		

1080	<b>H Neural Tangent Kernel</b>	<b>34</b>
1081		
1082	H.1 Kernel Perturbation . . . . .	34
1083	H.2 Kernel PSD during Training Process . . . . .	37
1084		
1085	<b>I Loss Decomposition</b>	<b>37</b>
1086		
1087	I.1 Bounding $C_0$ . . . . .	42
1088	I.2 Bounding $C_{1,2}$ . . . . .	46
1089	I.3 Bounding $C_2$ . . . . .	48
1090	I.4 Bounding $C_3$ . . . . .	49
1091	I.5 Bounding Loss during Training Process . . . . .	52
1092	I.6 Helpful Lemma . . . . .	52
1093		
1094		
1095		
1096	<b>J Convergence of Prefix Learning</b>	<b>55</b>
1097		
1098	J.1 Main Result . . . . .	55
1099	J.2 Induction Part 1. For Weights . . . . .	56
1100	J.3 Induction Part 2. For Loss . . . . .	56
1101	J.4 Induction Part 3. For Gradient . . . . .	57
1102	J.5 Bounding Loss at Initialization . . . . .	58
1103		
1104		
1105		
1106	<b>K NTK-Attention</b>	<b>58</b>
1107		
1108	K.1 Definitions . . . . .	58
1109	K.2 Error Bound . . . . .	58
1110	K.3 Tools from Fast Attention . . . . .	59
1111		
1112		
1113	<b>L Taylor Series</b>	<b>59</b>
1114		
1115		
1116		
1117		
1118		
1119		
1120		
1121		
1122	<b>Roadmap.</b> In Appendix A, we present the details of our method and prefix attention, and give a	
1123	complexity and memory analysis.	
1124	The experimental details for our empirical evaluation is shown in Appendix B. We give a naive	
1125	implementation of NTK-Attention within Python code in Appendix C. We provide more discussions	
1126	on our work in Appendix D, including the limitations and societal impacts of this paper.	
1127		
1128	We provide the preliminary we use in our analysis in Appendix E, including helpful probability tools.	
1129	We provide the basic definitions in Appendix F, and give helpful Lemmas about gradient computation	
1130	in Appendix G. Then we present our adaptation of NTK in our analysis in Appendix H, in Appendix I	
1131	show how to decompose the training objective to simplify proofs, and finally post our main results	
1132	and the proofs for analyzing the training in Appendix J.	
1133	In Appendix K, we compute the error bound on our NTK-Attention approximating ultra-long prefix	
	in attention. In Appendix L, we state helpful tools about the Taylor series.	

## A ALGORITHM DETAILS AND COMPUTATIONAL COMPLEXITY ANALYSIS

Here, we give the detailed version of two algorithms of this paper, which are prefix attention and NTK-Attention. Moreover, we comment on each computation step with its corresponding complexity to demonstrate our memory and complexity reduction in detail.

From Algorithm 3 and Algorithm 4, we can see the comparison analysis of memory reduction (from  $O(md)$  to  $O(rd + r)$ ) and complexity reduction (from  $O(mL + L^2)$  to  $O(Ld + L^2)$  since  $m \gg L$  and  $m \gg d$ ) between two fine-tuning methods, indicating the efficiency of our NTK-Attention.

---

### Algorithm 3 Prefix Attention (Detailed version of Algorithm 1)

---

**Input:** Input matrix  $X \in \mathbb{R}^{L \times d}$   
**Parameters:** Frozen query, key and value weights  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ , trainable prefix matrix  $P \in \mathbb{R}^{m \times d}$   $\triangleright$  Additional memory usage  $O(md)$   
**Output:** Exact output Attn  $\in \mathbb{R}^{L \times d}$

- 1: **procedure** PREFIXATTENTION( $X$ )
- 2: Concatenate input matrix with prefix matrix  $S \leftarrow \begin{bmatrix} P \\ X \end{bmatrix} \in \mathbb{R}^{(m+L) \times d}$
- 3: Compute query, key, and value matrices for attention  $Q \leftarrow XW_Q \in \mathbb{R}^{L \times d}$ ,  $K_P \leftarrow SW_K \in \mathbb{R}^{(m+L) \times d}$ ,  $V_P \leftarrow SW_V \in \mathbb{R}^{(m+L) \times d}$   $\triangleright$  Time complexity  $O(Ld^2 + 2(m+L)d^2)$
- 4: Compute exponential matrix  $A \leftarrow \exp(QK_P^\top / \sqrt{d}) \in \mathbb{R}^{L \times (m+L)}$   $\triangleright$  Time complexity  $O(L(m+L)d)$
- 5: Compute summation of exponential matrix  $D \leftarrow \text{diag}(A\mathbf{1}_{m+L}) \in \mathbb{R}^{L \times L}$   $\triangleright$  Time complexity  $O(L(m+L))$
- 6: Compute prefix attention output Attn  $\leftarrow D^{-1}AV_P \in \mathbb{R}^{L \times d}$   $\triangleright$  Here  $D^{-1}A \in \mathbb{R}^{L \times (m+L)}$  is the attention matrix (a.k.a attention scores). This step implements  $A$  multiply  $V_P$  first, then get  $D^{-1} \cdot (AV_P)$  with time complexity  $O(L(m+L)d + L^2d)$
- 7: **return** Attn
- 8: **end procedure**

---



---

### Algorithm 4 NTK-Attention (Detailed version of Algorithm 2, w low-rank)

---

**Input:** Input matrix  $X \in \mathbb{R}^{L \times d}$   
**Parameters:** Frozen query, key and value weights  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ , trainable weights  $Z_A \in \mathbb{R}^{r \times s}$ ,  $Z_B \in \mathbb{R}^{s \times d}$  and  $k \in \mathbb{R}^r$   $\triangleright$  Additional memory usage  $O(rs + sd + r)$   
**Output:** Approximating output  $T \in \mathbb{R}^{L \times d}$

- 1: **procedure** NTK-ATTENTION( $X$ )
- 2: Compute query, key, and value matrices for attention  $Q \leftarrow XW_Q \in \mathbb{R}^{L \times d}$ ,  $K \leftarrow XW_K \in \mathbb{R}^{L \times d}$ ,  $V \leftarrow XW_V \in \mathbb{R}^{L \times d}$   $\triangleright$  Time complexity  $O(3Ld^2)$
- 3: Compute approximating exponential matrix  $\hat{A} \leftarrow \exp(QK^\top / \sqrt{d}) \in \mathbb{R}^{L \times L}$   $\triangleright$  Time complexity  $O(L^2d)$
- 4: Compute approximating summation of exponential matrix  $\hat{D} \leftarrow \text{diag}(\hat{A}\mathbf{1}_L + \Phi(Q)k) \in \mathbb{R}^{L \times L}$   $\triangleright$  Time complexity  $O(L^2 + Lr)$
- 5: Compute approximation of prefix attention output  $T \leftarrow \hat{D}^{-1}(\hat{A}V + \Phi(Q)Z_A \cdot Z_B) \in \mathbb{R}^{L \times d}$   $\triangleright$  This step implements  $Z := Z_A \cdot Z_B$  first, compute  $\hat{A}V + \Phi(Q)Z$  secondly, then implements  $\hat{D}^{-1} \cdot (\hat{A}V + \Phi(Q)Z_A \cdot Z_B)$ , time complexity  $O(2L^2d + Lr^2 + rsd)$
- 6: **return**  $T$
- 7: **end procedure**

---

## B EXPERIMENTAL DETAILS

### B.1 SETUP DETAILS

Here, we give the details of the setup for the experiments in Section 5.

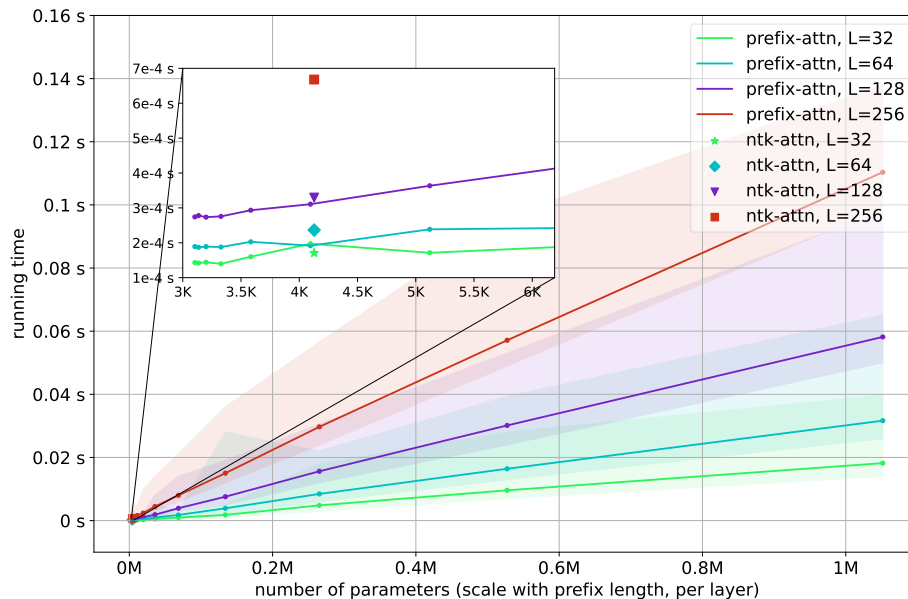
- Learning rate  $\eta = 0.001$  (default).
- Learning rate scheduler: Cosine.
- Adam hyper-parameter  $\beta_1 = 0.9$  (default).
- Adam hyper-parameter  $\beta_2 = 0.999$  (default).
- Adam hyper-parameter  $\epsilon = 1 \times 10^{-8}$  (default).
- Platform: PyTorch (Paszke et al., 2019) and Huggingface (Wolf et al., 2019).
- GPU device information: 8 V100 GPUs, 8 4090 GPUs and 4 H800 GPUs.
- Number of training epochs 30.
- Batch size for vision tasks: 256 (for best effort).
- Batch size for natural language task: 32 (for best effort).
- Max input length for natural language task: 128 for each feature, e.g. BoolQ has two dataset features: question and passage, for each data, we select the first 128 tokens in question and passage of the data respectively, and concatenate them to be the input.
- Quantization: fp16 and bf16.

## B.2 ADDITIONAL EMPIRICAL COMPLEXITY ANALYSIS

We state an additional empirical complexity analysis here to support our claim practically. We evaluate the complexity reduction on one layer to show how much efficiency our NTK-Attention will demonstrate per layer.

**Setup.** Firstly, we choose  $d = 32$  and  $r = d$ , and randomly initialize attention weights  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ . For the trainable parameters in NTK-Attention and Prefix Attention, we initialize  $P \in \mathbb{R}^{m \times d}$ ,  $Z \in \mathbb{R}^{d \times d}$  and  $k \in \mathbb{R}^d$  randomly, either. We then scale the prefix length, denotes  $m$ , within the range  $\{2^0, 2^1, \dots, 2^{16}\}$  for comparison. The input length  $L$  is chosen from  $\{32, 64, 128, 256\}$ . For computation, we initialize a new input matrix  $X \in \mathbb{R}^{L \times d}$  and compute NTK-Attention and Prefix Attention respectively. We repeat each computation with a different setup 10000 times and record the maximum, minimum, and mean values. The inference is run on an AMD CPU to compare FLOPS fairly between two algorithms (this also works on GPU devices).

Figure 3: Run time and the number of parameters of one-layer NTK-Attention and Prefix Attention (on random input data).  $x$ -axis: the number of parameters;  $y$ -axis: run time. Input length  $L$  is chosen from  $\{32, 64, 128, 256\}$ , dimension  $d = 32$  and prefix length  $m$  is chosen from  $\{2^0, 2^1, \dots, 2^{16}\}$ .



**Results.** We demonstrate our result in Figure 3. The  $x$ -axis is the number of parameters (representing memory usage), and the  $y$ -axis shows the run time in seconds. Note that the number of parameters is computed by the summation of every number in NTK-Attention or Prefix Attention. For example,  $m = 1024$ ,  $d = 32$ , the number of parameters of Prefix Attention is  $md + 3d^2 = 35840$ ; the number of parameters if NTK-Attention is  $4d^2 + d = 4128$ .

As expected, the number of parameters of Prefix Attention increases linearly with the prefix length  $m$ , and its running time increases quadratically with  $m$ . While our method, NTK-Attention, has computational costs unaffected by the prefix length. It maintains a small running time and low memory usage as shown in the figure. Roughly speaking, the cost of NTK-Attention is close to Prefix Attention with a very small prefix length  $m = 32$ .

### B.3 ADDITIONAL ABLATION STUDY

**Setup.** We provide an additional ablation study for the sensitivity of the hyper-parameters of NTK-Attention  $r$  and  $s$  here and the results are given in Table 4. In particular, this experiment is run on pretrained LLAMA-3.1-8B-Instruct model ( $d = 128$  for each head in attention) (Touvron et al., 2023a,b) with dataset WikiText-103 (Merity et al., 2016). We utilize 4 H800 GPU devices to train the model with different settings within 2 epochs on the training dataset and evaluate them on the test dataset. The metric is cross-entropy loss and its smaller value stands for better performance.

**Results.** We show the NTK-Attention with the weakest setting  $r = 128$ ,  $s = 4$  is able to achieve competitive performance with  $r = 256$ ,  $r = 64$ . This further ensures the parameter efficiency of NTK-Attention.

Moreover, Table 4 also demonstrates that choosing a big value for hyper-parameter  $r$  primarily will lead to better evaluation loss since NTK-Attention with  $(r, s) = (256, 32)$  requires 12.85M parameters but achieve superior performance compared to NTK-Attention with  $(r, s) = (128, 64)$  (requires 16.91M parameters).

However, we discover that an increased value for  $r$  might cause huge complexity - when setting  $r = 512$ , the computational complexity  $4Ld$  will lead the GPU out-of-memory (OOM) since it's usually unaffordable even for H800 (80GiB memory). Thus, we also suggest using  $r = d$  or  $r = 2d$  to make LLMs to learn downstream tasks.

Table 4: The results of ablation study to the NTK-Attention hyper-parameters  $r$  and  $s$  with pretrained LLM LLAMA-3.1-8B-Instruct and dataset WikiText-103 on H800 GPUs (80GiB).

Hyper-parameters	Num Parameters	Evaluation Loss	Training Loss
$(r, s)=(128, 4)$	1.18M	2.48	2.38
$(r, s)=(128, 8)$	2.23M	2.57	2.50
$(r, s)=(128, 16)$	4.33M	2.74	2.72
$(r, s)=(128, 32)$	8.52M	2.47	2.38
$(r, s)=(128, 64)$	16.91M	2.41	2.31
$(r, s)=(256, 4)$	1.84M	2.47	2.39
$(r, s)=(256, 8)$	3.41M	2.43	2.36
$(r, s)=(256, 16)$	6.55M	2.51	2.53
$(r, s)=(256, 32)$	12.85M	2.28	2.33
$(r, s)=(256, 64)$	25.43M	2.21	2.15
$(r, s)=(512, 4)$	3.15M (OOM since $4Ld$ complexity)	-	-

## C NAIVE NTK-ATTENTION IMPLEMENTATION WITH FLASH-ATTENTION

Below, we provide a naive Python code to implement our NTK-Attention that is written in only 10 lines, which supports the simplicity of implementation. Our code utilizes the function of Flash Attention function (Dao et al., 2022; Dao, 2023; Shah et al., 2024).



```

1296 1 def ntk_attn_forward(self, query_states, key_states, value_states,
1297   attention_mask):
1298 2   attn_outputs, lse = _flash_attention_forward(
1299 3     query_states, key_states, value_states, attention_mask,
1300 4     is_causal=self.is_causal, return_attn_probs=True
1301 5   ) # Call flash-attn function to get attn_output and logsumexp
1302 6
1302 7   Z = torch.matmul(self.Z_A, self.Z_B) # Low-rank approximate Z
1303 8   k = self.k
1304 9   phi_query_states = self.phi(query_states)
1305 10
1306 11   se = lse.exp() # Compute sumexp
1306 12   scale_factor = (se + torch.matmul(phi_query_states, k)) / se
1307 13
1308 14   attn_output = scale_factor * (attn_output * se + torch.matmul(
1309   phi_query_states, Z))
1310 15
1311 16   return attn_output

```

## D FURTHER DISCUSSIONS

Prior works (Arora et al., 2019; Alemohammad et al., 2020; Hron et al., 2020) had already given exact algorithms for computing the extension of NTK to neural nets and conducted experiments showing enhanced performance from adding NTK into models, while in this paper, our contributions are not limited to this. Our theory about NTK of attention with the infinite-long prefix provides more insights. We clarify this further in the following.

**Can LLMs master any advanced reasoning skill through self-planning and prompting?** We will answer that it may be possible. Since an attention network can converge on any dataset with the infinite-long prefix, we can tell that for any advanced reasoning skill that is equivalent to training on a well-constructed dataset, there exists an ultra-long prefix matrix satisfying the training objective smaller than any positive value  $\epsilon > 0$ . It’s noteworthy that this conclusion is not only suitable for LLMs with outstanding performance but also can be worked on those small language models with common performance.

**What is NTK-Attention used for? What is the meaning of proposing this method?** The attention with an infinite-long prefix is superior due to its over-parameterization phenomenon, whereas it is nearly impossible to implement practically, our NTK-Attention method gives us a chance to approximate the infinite-long prefix and makes it possible for us to study its empirical properties in experiments. Besides, any form of prefix learning can be formulated into the training of  $Z \in \mathbb{R}^{d \times d}$  and  $k \in \mathbb{R}^d$  in NTK-Attention, we can compress prompts into  $Z$  and  $k$  if  $\phi(\cdot)$  by utilizing Lemma K.7, hence, the approaches in Prefix Learning would be much more efficient.

**Comparison between NTK-Attention and LoRA.** LoRA in (Hu et al., 2021; Zeng & Lee, 2023; Hu et al., 2024) is a popular efficient fine-tuning method for large base models. Usually, LoRA makes adaptation on Query and Value projections  $W_Q, W_V \in \mathbb{R}^{d \times d}$ ; denote the adaptation as  $W_{\Delta Q}, W_{\Delta V} \in \mathbb{R}^{d \times d}$ . Given an input  $X \in \mathbb{R}^{L \times d}$ , LoRA computes  $\tilde{D}^{-1} \tilde{A} X (W_V + W_{\Delta V})$ , where  $\tilde{A} := \exp(X(W_Q + W_{\Delta Q})W_K^\top X^\top)$ ,  $\tilde{D} := \text{diag}(\tilde{A} \mathbf{1}_L)$ , and  $W_K \in \mathbb{R}^{d \times d}$  is the Key projection weights. So LoRA updates query and value weights during training, while our NTK-Attention compresses the additional prefix  $P$  into  $Z$  and  $k$  (Algorithm 2), which is a completely different mechanism. Our method also achieves comparable performance to LoRA in our experiments in Section 5. Also, note that the two methods are orthogonal to each other and can be used together.

**Connection to the newest SOTA LLM on math inference tasks, GPT-o1<sup>1</sup>.** On September 12-th, 2024, OpenAI released the newest SOTA LLM on math inference tasks, GPT-o1, which is trained by Reinforcement Learning (RL) methods to enhance the Chain-of-Thought (CoT) ability. Li et al. (2024c) explained the necessity of CoT for LLM on complicated inference tasks, meanwhile, they also emphasized how the embedding size and the CoT length affect the capability to solve high-order problems. Connecting to our work, we believe that these empirical and theoretical results support the

<sup>1</sup><https://openai.com/o1/>

conclusion of our work since we consider CoT as a specific application of Prefix Learning. Moreover, we think our *scaling law in prefix learning* is more universal for explaining the LLMs’ context-based advanced skills. However, even when we present our theory, we still have a limited understanding of prefix learning, for example, what is the relationship between prefix length and complexity of problems that aim to solve; if we want to solve an NP problem by LLM, how long is the prefix needed for inference? We don’t know the answers. Thus, explaining prefix learning, or particularly, CoT, is still a fascinating and challenging problem for future work.

**Limitations.** The work has limited experimental analysis and results. While empirical evaluations have been provided for some datasets and LLM models, the proposed method is widely applicable to different data and models, so comprehensive evaluations on more datasets and more practical methods can provide stronger empirical support.

Besides, the computational efficiency of NTK-Attention is insufficiently better than prefix attention when  $m < d$ , since the design of NTK-Attention is toward the ultra-big value of  $m$ , such we only compare to the prefix attention with prefix length  $m \gg d$  to meet the over-parameterization setting in our analysis.

**Societal impact.** This paper presents work whose goal is to advance the understanding of context-based fine-tuning methods (prefix learning) theoretically. There are many positive potential societal consequences of our work, such as inspiring new algorithm design. Since our work is theoretical in nature, we do not foresee any potential negative societal impacts which worth pointing out.

## E PRELIMINARY OF ANALYSIS

We provide our notations for this paper as follows:

**Notations** In this paper, we use integer  $d$  to denote the dimension of networks. We use integer  $m$  to denote the prefix length in prefix learning, we think  $m$  is an ultra-big number. We use  $L$  to denote the input length in language models.  $\nabla_x f(x)$  and  $\frac{df(x)}{dx}$  are both means to take the derivative of  $f(x)$  with  $x$ . Let a vector  $z \in \mathbb{R}^n$ . We denote the  $\ell_2$  norm as  $\|z\|_2 := (\sum_{i=1}^n z_i^2)^{1/2}$ , the  $\ell_1$  norm as  $\|z\|_1 := \sum_{i=1}^n |z_i|$ ,  $\|z\|_0$  as the number of non-zero entries in  $z$ ,  $\|z\|_\infty$  as  $\max_{i \in [n]} |z_i|$ . We use  $z^\top$  to denote the transpose of a  $z$ . We use  $\langle \cdot, \cdot \rangle$  to denote the inner product. Let  $A \in \mathbb{R}^{n \times d}$ , we use  $\text{vec}(A)$  to denote a length  $nd$  vector. We denote the Frobenius norm as  $\|A\|_F := (\sum_{i \in [n], j \in [d]} A_{i,j}^2)^{1/2}$ . For any positive integer  $n$ , we use  $[n]$  to denote set  $\{1, 2, \dots, n\}$ . We use  $\mathbb{E}[\cdot]$  to denote the expectation. We use  $\Pr[\cdot]$  to denote the probability. We use  $\epsilon$  to denote the error. We define  $\lambda_{\min}(\cdot)$  as a function that outputs the minimum eigenvalues of the input matrix, e.g. matrix  $A \in \mathbb{R}^{n \times n}$  has eigenvalues  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ ,  $\lambda_{\min}(A) = \min\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ .

### E.1 FACTS

**Fact E.1.** For any  $x \in (-0.01, 0.01)$ , we have

$$\exp(x) = 1 + x + \Theta(1)x^2.$$

**Fact E.2.** For any  $x \in (0, 0.1)$ , we have

$$\sum_{i=1}^n x^i \leq \frac{1}{1-x}.$$

### E.2 PROBABILITY

Here, we state a probability toolkit in the following, including several helpful lemmas we’d like to use. Firstly, we provide the lemma about Chernoff bound in (Chernoff, 1952) below.

**Lemma E.3** (Chernoff bound, (Chernoff, 1952)). Let  $X = \sum_{i=1}^n X_i$ , where  $X_i = 1$  with probability  $p_i$  and  $X_i = 0$  with probability  $1 - p_i$ , and all  $X_i$  are independent. Let  $\mu = \mathbb{E}[X] = \sum_{i=1}^n p_i$ . Then

- $\Pr[X \geq (1 + \delta)\mu] \leq \exp(-\delta^2\mu/3), \forall \delta > 0;$

$$\bullet \Pr[X \leq (1 - \delta)\mu] \leq \exp(-\delta^2\mu/1), \forall 0 < \delta < 1.$$

Next, we offer the lemma about Hoeffding bound as in (Hoeffding, 1994).

**Lemma E.4** (Hoeffding bound, (Hoeffding, 1994)). *Let  $X_1, \dots, X_n$  denote  $n$  independent bounded variables in  $[a_i, b_i]$  for  $a_i, b_i \in \mathbb{R}$ . Let  $X := \sum_{i=1}^n X_i$ , then we have*

$$\Pr[|X - \mathbb{E}[X]| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

We show the lemma of Bernstein inequality as (Bernstein, 1924).

**Lemma E.5** (Bernstein inequality, (Bernstein, 1924)). *Let  $X_1, \dots, X_n$  denote  $n$  independent zero-mean random variables. Suppose  $|X_i| \leq M$  almost surely for all  $i$ . Then, for all positive  $t$ ,*

$$\Pr\left[\sum_{i=1}^n X_i \geq t\right] \leq \exp\left(-\frac{t^2/2}{\sum_{j=1}^n \mathbb{E}[X_j^2] + Mt/3}\right)$$

Then, we give the Khintchine's inequality in (Khintchine, 1923; Haagerup, 1981) as follows:

**Lemma E.6** (Khintchine's inequality, (Khintchine, 1923; Haagerup, 1981)). *Let  $\sigma_1, \dots, \sigma_n$  be i.i.d sign random variables, and let  $z_1, \dots, z_n$  be real numbers. Then there are constants  $C > 0$  so that for all  $t > 0$*

$$\Pr\left[\left|\sum_{i=1}^n z_i \sigma_i\right| \geq t \|z\|_2\right] \leq \exp(-Ct^2).$$

We give Hason-wright inequality from (Hanson & Wright, 1971; Rudelson & Vershynin, 2013) below.

**Lemma E.7** (Hason-wright inequality, (Hanson & Wright, 1971; Rudelson & Vershynin, 2013)). *Let  $x \in \mathbb{R}^n$  denote a random vector with independent entries  $x_i$  with  $\mathbb{E}[x_i] = 0$  and  $|x_i| \leq K$ . Let  $A$  be an  $n \times n$  matrix. Then, for every  $t \geq 0$*

$$\Pr[|x^\top Ax - \mathbb{E}[x^\top Ax]| > t] \leq 2 \exp(-c \min\{t^2/(K^4 \|A\|_F^2), t/(K^2 \|A\|)\}).$$

We state Lemma 1 on page 1325 of Laurent and Massart (Laurent & Massart, 2000).

**Lemma E.8** (Lemma 1 on page 1325 of Laurent and Massart, (Laurent & Massart, 2000)). *Let  $X \sim \mathcal{X}_k^2$  be a chi-squared distributed random variable with  $k$  degrees of freedom. Each one has zero mean and  $\sigma^2$  variance. Then*

$$\begin{aligned} \Pr[X - k\sigma^2 \geq (2\sqrt{kt} + 2t)\sigma^2] &\leq \exp(-t) \\ \Pr[X - k\sigma^2 \leq 2\sqrt{kt}\sigma^2] &\leq \exp(-t). \end{aligned}$$

Here, we provide a tail bound for sub-exponential distribution (Foss et al., 2011).

**Lemma E.9** (Tail bound for sub-exponential distribution, (Foss et al., 2011)). *We say  $X \in \text{SE}(\sigma^2, \alpha)$  with parameters  $\sigma > 0, \alpha > 0$ , if*

$$\mathbb{E}[e^{\lambda X}] \leq \exp(\lambda^2 \sigma^2 / 2), \forall |\lambda| < 1/\alpha.$$

*Let  $X \in \text{SE}(\sigma^2, \alpha)$  and  $\mathbb{E}[X] = \mu$ , then:*

$$\Pr[|X - \mu| \geq t] \leq \exp(-0.5 \min\{t^2/\sigma^2, t/\alpha\}).$$

In the following, we show the helpful lemma of matrix Chernoff bound as in (Tropp, 2011; Lu et al., 2013).

**Lemma E.10** (Matrix Chernoff bound, (Tropp, 2011; Lu et al., 2013)). *Let  $\mathcal{X}$  be a finite set of positive-semidefinite matrices with dimension  $d \times d$ , and suppose that*

$$\max_{X \in \mathcal{X}} \lambda_{\max}(X) \leq B.$$

1458 *Sample  $\{X_1, \dots, X_n\}$  uniformly at random from  $\mathcal{X}$  without replacement. We define  $\mu_{\min}$  and  $\mu_{\max}$*   
 1459 *as follows:*

$$\begin{aligned} 1460 \mu_{\min} &:= n \cdot \lambda_{\min}(\mathbb{E}_{X \in \mathcal{X}}(X)) \\ 1461 & \\ 1462 \mu_{\max} &:= n \cdot \lambda_{\max}(\mathbb{E}_{X \in \mathcal{X}}(X)). \\ 1463 & \end{aligned}$$

1464 *Then*

$$\begin{aligned} 1465 \Pr[\lambda_{\min}(\sum_{i=1}^n X_i) \leq (1 - \delta)\mu_{\min}] &\leq d \cdot \exp(-\delta^2 \mu_{\min}/B) \text{ for } \delta \in (0, 1], \\ 1466 & \\ 1467 & \\ 1468 \Pr[\lambda_{\max}(\sum_{i=1}^n X_i) \geq (1 + \delta)\mu_{\max}] &\leq d \cdot \exp(-\delta^2 \mu_{\max}/(4B)) \text{ for } \delta \geq 0. \\ 1469 & \\ 1470 & \end{aligned}$$

## 1471 F DEFINITIONS OF NTK ANALYSIS

1472 This section provides the fundamental definitions of our NTK analysis in this paper.

1473 To begin with, we re-denote our weight of prefix in attention as  $W \in \mathbb{R}^{d \times m}$  and  $a \in \{-1, +1\}^m$  as  
 1474 follows<sup>2</sup>:

1475 **Definition F.1.** *We choose  $a \in \{-1, +1\}^m$  to be weights that each entry  $a_r$  is randomly sampled*  
 1476 *from  $-1$  with probability  $1/2$  and  $+1$  with probability  $1/2$ .*

1477 *Let  $W \in \mathbb{R}^{d \times m}$  denote random Gaussian weights, i.e., each entry independently draws from  $\mathcal{N}(0, \sigma^2)$ .*  
 1478 *For each  $r \in [m]$ , we use  $w_r \in \mathbb{R}^d$  to denote the  $r$ -th column of  $W$ .*

1479 Since we have established the equivalence between the ultra-long prefix matrix in attention and our  
 1480 theory in Section 3.1, it's reasonable we utilize the following definition of F to decompose the model  
 1481 function and facilitate our analysis.

1482 **Definition F.2.** *We define function  $F : \mathbb{R}^{d \times m} \times \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}^d$*

$$1483 F(W, x, a) = m \frac{\sum_{r \in [m]} a_r \exp(w_r^\top x) w_r}{\sum_{r \in [m]} \exp(w_r^\top x)}$$

1484 *Here we use  $w_r \in \mathbb{R}^d$  to denote the  $r$ -th column of  $W \in \mathbb{R}^{d \times m}$ .*

1485 To further break down the complicated F for more convenience analysis. We give an operator function  
 1486  $\alpha$  as follows:

1487 **Definition F.3.** *We define  $\alpha(x)$  as follows*

$$1488 \alpha(x) := \langle \underbrace{\exp(W^\top x)}_{m \times d \times d \times 1}, \mathbf{1}_m \rangle$$

1489 Thus, we can rewrite F in the following claim.

1490 **Claim F.4.** *We can rewrite  $F(W, x, a) \in \mathbb{R}^d$  as follows*

$$1491 F(W, x, a) = m \underbrace{\alpha(x)^{-1}}_{\text{scalar}} \underbrace{W}_{d \times m} \underbrace{(a)}_{m \times 1} \underbrace{\circ \exp(W^\top x)}_{m \times 1}$$

1492 *Proof.* We can show

$$\begin{aligned} 1493 F(W, x, a) &= m \frac{\sum_{r \in [m]} a_r \exp(w_r^\top x) w_r}{\sum_{r \in [m]} \exp(w_r^\top x)} \\ 1494 &= m \alpha(x)^{-1} \sum_{r \in [m]} a_r \exp(w_r^\top x) w_r \end{aligned}$$

1495 <sup>2</sup>Note that the proof of the case with  $a$  and without  $a$  are similar. We mainly focus on the proofs under the  
 1496 setting that includes  $a$ .

$$= m\alpha(x)^{-1}W(a \circ \exp(W^\top x))$$

where the first step follows from Definition F.2, the second step follows from Definition F.3 and simple algebras, the third step follows from  $w_r \in \mathbb{R}^d$  is denoting the  $r$ -th column of  $W \in \mathbb{R}^{d \times m}$  and simple algebras.  $\square$

In the following Definition F.6 and Definition F.5, we further derive and define two operator functions to convenient our analysis.

**Definition F.5.** We define  $\beta$  as follows

$$\beta_k := W_{k,*} \circ a, \forall k \in [d]$$

Let  $\beta \in \mathbb{R}^{d \times m}$  be defined as  $\beta = \underbrace{W}_{d \times m} \underbrace{\text{diag}(a)}_{m \times m}$

Here, we define softmax.

**Definition F.6.** We define  $S \in \mathbb{R}^m$  as follows

$$S := \underbrace{\alpha(x)^{-1}}_{\text{scalar}} \cdot \underbrace{\exp(W^\top x)}_{m \times 1}.$$

Here, we use  $\beta$  and  $S$  to re-denote the model function  $F$ .

**Definition F.7.** For each  $k \in [d]$ , let  $W_{k,*}^\top$  denote the  $k$ -th row of  $W$ , we define

$$F_k(W, x, a) := m \underbrace{\alpha(x)^{-1}}_{\text{scalar}} \langle \underbrace{W_{k,*}^\top}_{m \times 1} \circ \underbrace{a}_{m \times 1}, \underbrace{\exp(W^\top x)}_{m \times 1} \rangle$$

Then, we can rewrite it as

$$F_k(W, x, a) := m \langle \beta_k, S \rangle.$$

## F.1 LOSS FUNCTION

Here, we state the training objective that we aim to solve in the analysis.

**Definition F.8.** Given a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}^d$ . Let function  $F : \mathbb{R}^{d \times m} \times \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}^d$  be defined as Definition F.2, we define the training objective  $\mathcal{L} : \mathbb{R}^{m \times d} \rightarrow \mathbb{R}$  as follows:

$$\mathcal{L}(W) := 0.5 \sum_{i=1}^n \|F(W, x_i, a) - y_i\|_2^2$$

## G GRADIENT COMPUTATION

In this section, we first compute the gradients that we need for the analysis of NTK. Then we define the training dynamic of our model in the process of gradient descent.

### G.1 COMPUTING GRADIENT

We give our computation of the gradients as the following lemma.

**Lemma G.1.** If the following conditions hold

- Let  $W \in \mathbb{R}^{d \times m}$  and  $a \in \mathbb{R}^m$  be defined as Definition F.1.
- Let  $\alpha(x) \in \mathbb{R}$  be defined as Definition F.3
- Let  $S \in \mathbb{R}^m$  be defined as Definition F.6
- Let  $F \in \mathbb{R}^d$  be defined as Definition F.7

Then, we can show that for each  $r \in [m]$

1566 • **Part 1.** For  $k_1 \in [d]$ , we have

$$1567 \frac{dW^\top x}{dw_{r,k_1}} = x_{k_1} e_r$$

1571 • **Part 2.** For  $k_1 \in [d]$ , we have

$$1572 \frac{d \exp(W^\top x)}{dw_{r,k_1}} = (x_{k_1} e_r) \circ \exp(W^\top x)$$

1575 • **Part 3.** For  $k_1 \in [d]$ , we have

$$1577 \frac{d\alpha(x)}{dw_{r,k_1}} = \langle x_{k_1} e_r, \exp(W^\top x) \rangle$$

1580 • **Part 4.** For  $k_1 \in [d]$ , we have

$$1581 \frac{d\alpha(x)^{-1}}{dw_{r,k_1}} = -\alpha(x)^{-1} \langle x_{k_1} e_r, S \rangle$$

1584 • **Part 5.** For  $k_1 \in [d]$ , we have

$$1585 \frac{dS}{dw_{r,k_1}} = -\langle x_{k_1} e_r, S \rangle \cdot S + (x_{k_1} e_r) \circ S$$

1589 • **Part 6.** For  $k_1, k \in [d]$  and  $k_1 \neq k$ , we have

$$1590 \frac{dF(W, x, a)_k}{dw_{r,k_1}} = +0 - mx_{k_1} \cdot S_r \cdot \langle \beta_k, S \rangle + mx_{k_1} S_r \beta_{k,r}$$

1593 • **Part 7.** For  $k_1, k \in [d]$  and  $k_1 = k$ , we have

$$1594 \frac{dF(W, x, a)_k}{dw_{r,k}} = +m \langle a \circ e_r, S \rangle - mx_k \cdot S_r \cdot \langle \beta_k, S \rangle + mx_k S_r \beta_{k,r}$$

1598 • **Part 8.** For  $k \in [d]$ , we have

$$1599 \frac{dF(W, x, a)_k}{dw_r} = ma_r S_r \cdot e_k - m \langle \beta_k, S \rangle S_r \cdot x + m \beta_{k,r} S_r \cdot x$$

1602 **Proof. Proof of Part 1.**

$$1603 \frac{dW^\top x}{dw_{r,k_1}} = x_{k_1} e_r$$

1604 where this step follows from simple differential rules.

1608 **Proof of Part 2.**

$$1609 \frac{d \exp(W^\top x)}{dw_{r,k_1}} = \exp(W^\top x) \circ \frac{dW^\top x}{dw_{r,k_1}}$$

$$1610 = (x_{k_1} e_r) \circ \exp(W^\top x)$$

1611 where the first step follows from chain rules, the second step follows from Part 1 of this Lemma.

1616 **Proof of Part 3.**

$$1617 \frac{d\alpha(x)}{dw_{r,k_1}} = \left\langle \frac{d \exp(W^\top x)}{dw_{r,k_1}}, \mathbf{1}_m \right\rangle$$

1620

$$= \langle x_{k_1} e_r, \exp(W^\top x) \rangle$$

1621

where the first step follows from Definition F.3 and simple algebras, the second step follows from Part 2 of this Lemma.

1622

1623

1624

**Proof of Part 4.**

1625

1626

$$\begin{aligned} \frac{d\alpha(x)^{-1}}{dw_{r,k_1}} &= -\alpha(x)^{-2} \frac{d\alpha(x)}{dw_{r,k_1}} \\ &= -\alpha(x)^{-1} \langle x_{k_1} e_r, S \rangle \end{aligned}$$

1627

1628

1629

where this step follows from chain rules, the second step follows from Part 3 of this Lemma.

1630

1631

**Proof of Part 5.**

1632

1633

$$\begin{aligned} \frac{dS}{dw_{r,k_1}} &= \frac{d\alpha(x)^{-1}}{dw_{r,k_1}} \cdot \exp(W^\top x) + \alpha(x)^{-1} \cdot \frac{d\exp(W^\top x)}{dw_{r,k_1}} \\ &= -\alpha(x)^{-1} \langle x_{k_1} e_r, S \rangle \cdot \exp(W^\top x) + \alpha(x)^{-1} \cdot (x_{k_1} e_r) \circ \exp(W^\top x) \\ &= -\langle x_{k_1} e_r, S \rangle \cdot S + (x_{k_1} e_r) \circ S \end{aligned}$$

1634

1635

1636

1637

1638

where the first step follows from Definition F.6 and differential rules, the second step follows from Part 2 and Part 4 of this Lemma, the last step follows from simple algebras.

1639

1640

1641

**Proof of Part 6.** For  $k_1 \neq k$

1642

$$\begin{aligned} \frac{dF(W, x, a)_k}{dw_{r,k_1}} &= +m \langle \frac{d\beta_k}{dw_{r,k_1}}, S \rangle + m \langle \beta_k, \frac{dS}{dw_{r,k_1}} \rangle \\ &= -m \langle x_{k_1} e_r, S \rangle \cdot \langle \beta_k, S \rangle + m \langle \beta_k, (x_{k_1} e_r) \circ S \rangle \\ &= +0 - mx_{k_1} \cdot S_r \cdot \langle \beta_k, S \rangle + mx_{k_1} S_r \beta_{k,r} \end{aligned}$$

1643

1644

1645

1646

where the first step follows from Definition F.7 and simple algebras, the second step follows from Definition F.5, simple algebras and Part 5 of this Lemma, the last step follows from simple algebras.

1647

1648

1649

**Proof of Part 7.** For  $k_1 = k$

1650

$$\begin{aligned} \frac{dF(W, x, a)_k}{dw_{r,k}} &= +m \langle \frac{d\beta_k}{dw_{r,k}}, S \rangle + m \langle \beta_k, \frac{dS}{dw_{r,k}} \rangle \\ &= +m \langle a \circ e_r, S \rangle - m \langle x_k e_r, S \rangle \cdot \langle \beta_k, S \rangle + m \langle \beta_k, (x_k e_r) \circ S \rangle \\ &= +m \langle a \circ e_r, S \rangle - mx_k \cdot S_r \cdot \langle \beta_k, S \rangle + mx_k S_r \beta_{k,r} \end{aligned}$$

1651

1652

1653

1654

where the first step follows from Definition F.7 and simple algebras, the second step follows from Definition F.5, simple algebras and Part 5 of this Lemma, the last step follows from simple algebras.

1655

1656

1657

**Proof of Part 8.**

1658

This part of proof follows from the combination of Part 6 and Part 7 of this Lemma.  $\square$

1659

1660

## G.2 GRADIENT DESCENT

1661

1662

After we computed the gradient of the model function above, we are now able to define the training dynamic of F by updating weight using gradient descent.

1663

1664

1665

We use  $e_r$  to denote a vector where the  $r$ -th coordinate is 1 and everywhere else is 0.  $\forall r \in [m], \forall k \in [d]$ , we have  $\frac{dF(W, x, a)_k}{dw_r} \in \mathbb{R}^d$  can be written as

1666

1667

1668

1669

$$\underbrace{\frac{dF_k(W, x, a)}{dw_r}}_{d \times 1} = ma_r S_r \cdot e_k - m \langle \beta_k, S \rangle S_r \cdot x + m \beta_{k,r} S_r \cdot x. \quad (9)$$

1670

1671

Hence, by defining several following dynamical operator functions, we can further convenient our proofs.

1672

1673

We first define  $u_i(\tau) \in \mathbb{R}^m$  for simplification as follows:

1674 **Definition G.2.** For each  $i \in [n]$ , we define  $\mathbf{u}_i(\tau) \in \mathbb{R}^m$  as

$$1675 \underbrace{\mathbf{u}_i(\tau)}_{m \times 1} := \exp\left(\underbrace{W(\tau)}_{m \times d}^\top \underbrace{x_i}_{d \times 1}\right)$$

1679 Secondly, we re-denote  $\alpha_i(\tau) \in \mathbb{R}$  below, which holds due to the definition of  $\alpha(x)$  and the updating  
1680 of  $W \in \mathbb{R}^{d \times m}$ .

1681 **Definition G.3.** For each  $i \in [n]$ , we define  $\alpha_i(\tau) \in \mathbb{R}$  as

$$1682 \underbrace{\alpha_i(\tau)}_{\text{scalar}} := \langle \underbrace{\mathbf{u}_i(\tau)}_{m \times 1}, \underbrace{\mathbf{1}_m}_{m \times 1} \rangle.$$

1686 We define  $\beta_k(\tau) \in \mathbb{R}^m$  for convenience.

1687 **Definition G.4.** For each  $k \in [d]$ , we define  $\beta_k(\tau) \in \mathbb{R}^m$  as

$$1688 \underbrace{\beta_k(\tau)}_{m \times 1} = \underbrace{(W_{k,*}(\tau))}_{m \times 1} \circ \underbrace{a}_{m \times 1}$$

1691 **Remark G.5.** The purpose of defining notation  $\beta$  is to make our proofs more aligned with softmax  
1692 NTK proofs in previous work (Li et al., 2024a).

1694 We define  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  for convenience as follows :

1695 **Definition G.6.** For each  $i \in [n]$ , for each  $k \in [d]$ , we define  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  as follows

$$1697 \underbrace{\theta_{k,i}(\tau)}_{m \times 1} := \underbrace{\beta_k(\tau)}_{m \times 1} \cdot \underbrace{\alpha_i(\tau)^{-1}}_{\text{scalar}}$$

1700 We denote  $S_r(\tau)$ .

1701 **Definition G.7.** For each  $i \in [n]$ . Let  $S_i(\tau) \in \mathbb{R}^m$  be defined as

$$1702 \underbrace{S_i(\tau)}_{m \times 1} := \underbrace{\alpha_i(\tau)^{-1}}_{\text{scalar}} \cdot \underbrace{\mathbf{u}_i(\tau)}_{m \times 1}$$

1706 for integer  $\tau \geq 0$ . For  $r \in [m]$ , we denote  $S_{i,r}(\tau) \in \mathbb{R}$  as the  $r$ -th entry of vector  $S_i(\tau)$ .

1707 Now, we can define F at different timestamps.

1708 **Definition G.8** (F( $\tau$ ), dynamic prediction). For each  $k \in [d]$ , for each  $i \in [n]$ , we define  $F_i(\tau) \in \mathbb{R}^d$ ,  
1709 for any timestamp  $\tau$ , as

$$1711 F_{k,i}(\tau) := m \langle \mathbf{u}(\tau), \mathbf{1}_m \rangle^{-1} \langle W(\tau)_{k,*} \circ a, \mathbf{u}(\tau) \rangle.$$

1712 Here  $x_i \in \mathbb{R}^d$ . It can be rewritten as

$$1713 F_{k,i}(\tau) = m \cdot \langle \underbrace{\beta_k(\tau)}_{m \times 1}, \underbrace{S_i(\tau)}_{m \times 1} \rangle.$$

1716 and also

$$1717 F_{k,i}(\tau) = m \cdot \langle \underbrace{\theta_{k,i}(\tau)}_{m \times 1}, \underbrace{\mathbf{u}_i(\tau)}_{m \times 1} \rangle$$

1721 We consider  $d$ -dimensional MSE loss.

1722 **Definition G.9** (Loss function over time). We define the objective function  $\mathcal{L}$  as below:

$$1723 \mathcal{L}(W(\tau)) := \frac{1}{2} \sum_{i \in [n]} \sum_{k \in [d]} (F_{k,i}(\tau) - y_{k,i})^2.$$

1727 Thus, we define the gradient of  $w$ .



**Definition G.10** ( $\Delta w_r(\tau)$ ). For any  $r \in [m]$ , we define  $\Delta w_r(\tau) \in \mathbb{R}^d$  as below:

$$\begin{aligned} & \Delta w_r(\tau) \\ & := m \sum_{i=1}^n \sum_{k=1}^d (\mathbb{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( a_r \mathbf{S}_{i,r}(\tau) \cdot e_k - \langle \beta_k(\tau), \mathbf{S}_i(\tau) \rangle \mathbf{S}_{i,r}(\tau) \cdot x + \beta_{k,r} \mathbf{S}_{i,r}(\tau) \cdot x \right) \end{aligned}$$

Here, we utilize  $v$  to simplify  $\Delta w_r(\tau)$ , we have the following:

**Definition G.11.** For each  $k \in [d]$ , for each  $r \in [m]$ , we define  $v_{k,r}(\tau) \in \mathbb{R}^m$  as follows

$$v_{k,r}(\tau) := \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau).$$

Note that we can simplify the gradient calculation by the fact  $1 = \langle \mathbf{1}_m, \mathbf{S}_i(\tau) \rangle$  for  $i \in [n]$ . Thus, we have the following claim.

**Claim G.12.** We can rewrite  $\Delta w_r(\tau)$  as follows

$$\Delta w_r(\tau) = m \sum_{i=1}^n \sum_{k=1}^d (\mathbb{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( \langle v_{k,r}(\tau), \mathbf{S}_i(\tau) \rangle \cdot \mathbf{S}_{i,r}(\tau) \cdot x_i + a_r \mathbf{S}_{i,r}(\tau) e_k \right)$$

*Proof.* We have

$$\begin{aligned} & \Delta w_r(\tau) \\ & = m \sum_{i=1}^n \sum_{k=1}^d (\mathbb{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( a_r \mathbf{S}_{i,r}(\tau) \cdot e_k - \langle \beta_k(\tau), \mathbf{S}_i(\tau) \rangle \mathbf{S}_{i,r}(\tau) \cdot x + \beta_{k,r} \mathbf{S}_{i,r}(\tau) \cdot x \right) \\ & = m \sum_{i=1}^n \sum_{k=1}^d (\mathbb{F}_{k,i}(\tau) - y_{k,i}) \\ & \quad \cdot \left( a_r \mathbf{S}_{i,r}(\tau) \cdot e_k - \langle \beta_k(\tau), \mathbf{S}_i(\tau) \rangle \mathbf{S}_{i,r}(\tau) \cdot x + \beta_{k,r} \langle \mathbf{1}_m, \mathbf{S}_i(\tau) \rangle \mathbf{S}_{i,r}(\tau) \cdot x \right) \\ & = m \sum_{i=1}^n \sum_{k=1}^d (\mathbb{F}_{k,i}(\tau) - y_{k,i}) \\ & \quad \cdot \left( a_r \mathbf{S}_{i,r}(\tau) \cdot e_k - \langle \beta_k(\tau), \mathbf{S}_i(\tau) \rangle \mathbf{S}_{i,r}(\tau) \cdot x + \langle \beta_{k,r} \cdot \mathbf{1}_m, \mathbf{S}_i(\tau) \rangle \mathbf{S}_{i,r}(\tau) \cdot x \right) \\ & = m \sum_{i=1}^n \sum_{k=1}^d (\mathbb{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( a_r \mathbf{S}_{i,r}(\tau) \cdot e_k + \langle \beta_{k,r} \cdot \mathbf{1}_m - \beta_k(\tau), \mathbf{S}_i(\tau) \rangle \mathbf{S}_{i,r}(\tau) \cdot x \right) \\ & = m \sum_{i=1}^n \sum_{k=1}^d (\mathbb{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( a_r \mathbf{S}_{i,r}(\tau) \cdot e_k + \langle v_{k,r}(\tau), \mathbf{S}_i(\tau) \rangle \mathbf{S}_{i,r}(\tau) \cdot x \right) \end{aligned}$$

where the first step follows from Definition G.10, the second step follows from the fact  $1 = \langle \mathbf{1}_m, \mathbf{S}_i(\tau) \rangle$  for  $i \in [n]$ , the third and fourth steps follow from simple algebras, the last step follows from Definition G.11.

□

We use the gradient descent (GD) algorithm with the learning rate  $\eta$  to train the network. As we only train the hidden layer  $W$  and fix  $a$ , we have the following gradient update rule.

**Definition G.13** (Gradient descent). The gradient descent algorithm for optimizing the weight matrix  $W$  is defined as:

$$W(\tau + 1) = W(\tau) - \eta \Delta W(\tau).$$

where  $\Delta W(\tau) \in \mathbb{R}^{d \times m}$  and  $\Delta w_r(\tau) \in \mathbb{R}^d$  is the  $r$ -th column of  $\Delta W(\tau)$  defined in Definition G.10.

## 1782 H NEURAL TANGENT KERNEL

1783  
1784 Now in this section, we give the exact computation of NTK in our analysis below.

1785 **Definition H.1** (Kernel function, Definition 3.6 in (Li et al., 2024a)). *For simplicity, we denote*  
1786  $S(W^\top x_i)$  *as*  $S_i \in \mathbb{R}_{\geq 0}^m$  *and*  $v_{k,r} = \beta_{k,r} \cdot \mathbf{1}_m - \beta_k \in \mathbb{R}^m$ . *We define the function (Gram matrix)*  
1787  $H : \mathbb{R}^{d \times m} \rightarrow \mathbb{R}^{nd \times nd}$  *as following*

$$1788 \quad H(W) := \begin{bmatrix} H_{1,1} & H_{1,2} & \cdots & H_{1,d} \\ H_{2,1} & H_{2,2} & \cdots & H_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ H_{d,1} & H_{d,2} & \cdots & H_{d,d} \end{bmatrix},$$

1789  
1790 and for each  $k_1, k_2 \in [d]$ , we have  $H_{k_1, k_2} \in \mathbb{R}^{n \times n}$  is defined as

$$1791 \quad [H_{k_1, k_2}]_{i,j}(W) := \frac{1}{m} x_i^\top x_j \sum_{r=1}^m \langle v_{k_1, r}, S_i \rangle \cdot m S_{i,r} \cdot \langle v_{k_2, r}, S_j \rangle \cdot m S_{j,r}.$$

1792  
1793 For any timestamp  $\tau$ , for simplicity, we denote  $H(\tau) := H(W(\tau))$  and denote  $H(0)$  as  $H^*$ .

### 1800 H.1 KERNEL PERTURBATION

1801 The purpose of this section is to prove Lemma H.3. In the proof, we do not use concentration  
1802 inequality. Please see Remark H.2 for more details.

1803 **Remark H.2.** *In the proof of Lemma H.3, we do not use concentration bound as previous work (Song*  
1804 *& Yang, 2019; Munteanu et al., 2022; Gao et al., 2023a). The reason is that we consider the worst*  
1805 *case. In general,  $\mathbb{E}[H(W) - H(\widetilde{W})] \neq \mathbf{0}_{nd \times nd}$ . Thus, using the concentration bound may not gain*  
1806 *any benefits.*

1807 **Lemma H.3.** *If the following conditions hold*

- 1808 • Let  $C > 10$  denote a sufficiently large constant
- 1809 • Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- 1810 • Let  $R \in (0, 0.01)$ .
- 1811 • Let  $x_i \in \mathbb{R}^d$  and  $\|x_i\|_2 \leq 1$  for all  $i \in [n]$ .
- 1812 • Let  $\widetilde{W} = [\widetilde{w}_1, \dots, \widetilde{w}_m] \in \mathbb{R}^{d \times m}$ , where  $\widetilde{w}_1, \dots, \widetilde{w}_m$  are i.i.d. draw from  $\mathcal{N}(0, \sigma^2 I_d)$ .
- 1813 • Let  $W = [w_1, \dots, w_m] \in \mathbb{R}^{d \times m}$  and satisfy  $\|\widetilde{w}_r - w_r\|_2 \leq R$  for any  $r \in [m]$ .
- 1814 • Let  $v_{k,r} = \beta_{k,r} \cdot \mathbf{1}_m - \beta_k \in \mathbb{R}^m$ , for any  $k \in [d]$  and for any  $r \in [m]$ . Note that  $\beta_{k,r}$  is the  
1815  $r$ -th in  $\beta_k$ .
- 1816 • Let  $\alpha_i = \langle \mathbf{1}_m, \exp(W^\top x_i) \rangle$  and  $\tilde{\alpha}_i = \langle \mathbf{1}_m, \exp(\widetilde{W}^\top x_i) \rangle, \forall i \in [n]$ .
- 1817 • Let  $H$  be defined as Definition H.1.

1818 Then, we have

- 1819 • Part 1. Then with probability at least  $1 - \delta / \text{poly}(nd)$ ,

$$1820 \quad |[H_{k_1, k_2}]_{i,j}(W) - [H_{k_1, k_2}]_{i,j}(\widetilde{W})| \leq 8R \cdot \exp(22B).$$

- 1821 • Part 2. Then with probability at least  $1 - \delta$ , we have

$$1822 \quad \|H(W) - H(\widetilde{W})\|_F \leq 8R\sqrt{nd} \cdot \exp(22B).$$

1823 *Proof.* For simplicity, we give the following notations:

- Note that  $\tilde{S}_i := \exp(\tilde{W}(\tau)^\top x_i) \cdot \tilde{\alpha}_i^{-1}$ .
- Note that  $\tilde{\beta}_k := \tilde{W}_{k,*} \circ a$ .
- Note that  $\tilde{v}_{k,r} := \tilde{\beta}_{k,r} \cdot \mathbf{1}_m - \tilde{\beta}_k$ .

**Proof of Part 1.** We have

$$|[H_{k_1, k_2}]_{i,j}(W) - [H_{k_1, k_2}]_{i,j}(\tilde{W})| = mx_i^\top x_j \sum_{r=1}^m (B_{1,r} + B_{2,r} + B_{3,r} + B_{4,r} + B_{5,r} + B_{6,r})$$

here, we define:

$$\begin{aligned} B_{1,r} &:= \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, S_j \rangle \cdot S_{j,r} - \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, S_j \rangle \cdot \tilde{S}_{j,r} \\ B_{2,r} &:= \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, S_j \rangle \cdot \tilde{S}_{j,r} - \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} \\ B_{3,r} &:= \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} - \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} \\ B_{4,r} &:= \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} - \langle v_{k_1,r}, S_i \rangle \cdot \tilde{S}_{i,r} \cdot \langle \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} \\ B_{5,r} &:= \langle v_{k_1,r}, S_i \rangle \cdot \tilde{S}_{i,r} \cdot \langle \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} - \langle v_{k_1,r}, \tilde{S}_i \rangle \cdot \tilde{S}_{i,r} \cdot \langle \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} \\ B_{6,r} &:= \langle v_{k_1,r}, \tilde{S}_i \rangle \cdot \tilde{S}_{i,r} \cdot \langle \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} - \langle \tilde{v}_{k_1,r}, \tilde{S}_i \rangle \cdot \tilde{S}_{i,r} \cdot \langle \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} \end{aligned}$$

Before we bound all terms, we provide a tool as follows:

$$\begin{aligned} \|v_{k,r} - \tilde{v}_{k,r}\|_2^2 &= \sum_{r_1=1}^m (v_{k,r,r_1} - \tilde{v}_{k,r,r_1})^2 \\ &= \sum_{r_1=1}^m (\beta_{k,r} - \beta_{k,r_1} - \tilde{\beta}_{k,r} + \tilde{\beta}_{k,r_1})^2 \\ &= \sum_{r_1=1}^m (a_r W_{k,r} - a_{r_1} W_{k,r} - a_r \tilde{W}_{k,r} + a_{r_1} \tilde{W}_{k,r})^2 \\ &= \sum_{r_1=1}^m (a_r (W_{k,r} - \tilde{W}_{k,r}) + a_{r_1} (\tilde{W}_{k,r_1} - W_{k,r_1}))^2 \\ &\leq \sum_{r_1=1}^m (|W_{k,r} - \tilde{W}_{k,r}| + |\tilde{W}_{k,r_1} - W_{k,r_1}|)^2 \\ &\leq \sum_{r_1=1}^m 4R^2 \\ &\leq m4R^2 \end{aligned} \tag{10}$$

where the first step follows from the definition of  $\ell_2$  norm, the second step follows from the definition of  $v_{k,r}$ , the third step follows from Definition F.5, the fourth and fifth steps follow from simple algebras, the sixth step follows from  $\|w_r - v_r\|_\infty \leq \|w_r - v_r\|_2 \leq R$ , the last step follows from simple algebras.

To bound  $B_{1,r}$ , we have

$$\begin{aligned} |B_{1,r}| &:= |\langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, S_j \rangle \cdot S_{j,r} - \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, S_j \rangle \cdot \tilde{S}_{j,r}| \\ &= |\langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, S_j \rangle \cdot (S_{j,r} - \tilde{S}_{j,r})| \\ &\leq \frac{\exp(15B)}{m} \cdot |S_{j,r} - \tilde{S}_{j,r}| \\ &\leq \frac{R \exp(19B + 3R)}{m^2} \end{aligned}$$

where the first step follows from the definition of  $B_{1,r}$ , the second step follows from simple algebras, the third step follows from Part 6 of Lemma L.2 and  $0 \leq S_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1, the last step follows from Part 12 of Lemma L.1.

To bound  $B_{2,r}$ , we have

$$\begin{aligned} |B_{2,r}| &:= |\langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, S_j \rangle \cdot \tilde{S}_{j,r} - \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r}| \\ &= |\langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, S_j - \tilde{S}_j \rangle \cdot \tilde{S}_{j,r}| \\ &\leq \frac{2B \exp(12B)}{m^2} \cdot |\langle \frac{1}{2B} v_{k_2,r}, S_j - \tilde{S}_j \rangle| \\ &\leq \frac{2BR \exp(16B + 3R)}{m^2} \end{aligned}$$

where the first step follows from the definition of  $B_{2,r}$ , the second step follows from simple algebras, the third step follows from Part 6 of Lemma L.2 and  $0 \leq S_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1, the last step follows from Part 13 of Lemma L.1 and  $\|v_{k,r}\|_\infty \leq 2B$  by simple algebras.

To bound  $B_{3,r}$ , we have

$$\begin{aligned} |B_{3,r}| &:= |\langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r} - \langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r}| \\ &= |\langle v_{k_1,r}, S_i \rangle \cdot S_{i,r} \cdot \langle v_{k_2,r} - \tilde{v}_{k_2,r}, \tilde{S}_j \rangle \cdot \tilde{S}_{j,r}| \\ &\leq \frac{\exp(12B)}{m^2} \cdot |\langle v_{k_2,r} - \tilde{v}_{k_2,r}, \tilde{S}_j \rangle| \\ &\leq \frac{2R \exp(15B)}{m^2} \end{aligned}$$

where the first step follows from the definition of  $B_{3,r}$ , the second step follows from simple algebras, the third step follows from Part 6 of Lemma L.2 and  $0 \leq S_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1, the last step follows from Cauchy-Schwarz inequality, Eq. (10) and  $\|S_i\|_2 \leq \frac{\exp(3B)}{\sqrt{m}}$ .

The proof of bounding  $B_{4,r}$  is similar to the proof of bounding  $B_{1,r}$ , we have  $|B_{4,r}| \leq \frac{R \exp(19B+3R)}{m^2}$ .

The proof of bounding  $B_{5,r}$  is similar to the proof of bounding  $B_{2,r}$ , we have  $|B_{5,r}| \leq \frac{2BR \exp(16B+3R)}{m^2}$ .

The proof of bounding  $B_{6,r}$  is similar to the proof of bounding  $B_{3,r}$ , we have  $|B_{6,r}| \leq \frac{2R \exp(15B)}{m^2}$ .

Now we combine all terms, we have

$$\begin{aligned} |[H_{k_1,k_2}]_{i,j}(W) - [H_{k_1,k_2}]_{i,j}(\tilde{W})| &= m x_i^\top x_j \sum_{r=1}^m (B_{1,r} + B_{2,r} + B_{3,r} + B_{4,r} + B_{5,r} + B_{6,r}) \\ &\leq m \sum_{r=1}^m (B_{1,r} + B_{2,r} + B_{3,r} + B_{4,r} + B_{5,r} + B_{6,r}) \\ &\leq m \sum_{r=1}^m (|B_{1,r}| + |B_{2,r}| + |B_{3,r}| + |B_{4,r}| + |B_{5,r}| + |B_{6,r}|) \\ &\leq m \sum_{r=1}^m \frac{8R \exp(22B)}{m^2} \\ &\leq 8R \cdot \exp(22B) \end{aligned}$$

where the second step follows from  $\|x_i\|_2 \leq 1$ , the third step follows from simple algebras, the fourth step follows from  $R \leq B$ ,  $B \leq \exp(B)$  and the combination of all terms, the last step follows from simple algebras.

**Proof of Part 2.** This proof follows from Part 1 of this Lemma and the definition of Frobenius norm.  $\square$

## 1944 H.2 KERNEL PSD DURING TRAINING PROCESS

1945  
1946 **Claim H.4.** *If the following conditions hold:*

- 1947 • *Let  $\lambda = \lambda_{\min}(H^*)$*
- 1948 • *Let  $C > 10$  denote a sufficiently large constant*
- 1949 • *Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .*
- 1950 • *Let  $\delta \in (0, 0.1)$ .*
- 1951 • *Let timestamp  $\tau \geq 0$  denotes as a integer.*
- 1952 • *Denote  $H^*$  as  $H(W)$  in Definition H.1.*
- 1953 • *Denote  $H(\tau)$  as  $H(\widetilde{W})$  in Definition H.1.*
- 1954 • *Let  $D := 2\lambda^{-1} \cdot \exp(20B) \frac{\sqrt{nd}}{m} \|Y - F(0)\|_F$*
- 1955 • *Let  $\|w_r(t) - w_r(0)\|_2 \leq D < R = \lambda / \text{poly}(n, d, \exp(B))$ ,  $\forall r \in [m], \forall t \geq 0$*

1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963 *Then, with a probability at least  $1 - \delta$ , we have*

$$1964 \lambda_{\min}(H(\tau)) \geq \lambda/2$$

1965  
1966 *Proof.* By Lemma H.3, with a probability at least  $1 - \delta$ , we have

$$1967 \begin{aligned} 1968 \|H^* - H(\tau)\|_F &\leq 8R\sqrt{nd}\exp(22B) \\ 1969 &\leq \lambda/2 \end{aligned} \tag{11}$$

1970  
1971 where the first step follows from Part 2 of Lemma H.3, the second step follows by choice of  $R$ .

1972 By eigenvalue perturbation theory, we have

$$1973 \begin{aligned} 1974 \lambda_{\min}(H(\tau)) &\geq \lambda_{\min}(H^*) - \|H(\tau) - H^*\| \\ 1975 &\geq \lambda_{\min}(H^*) - \|H(\tau) - H^*\|_F \\ 1976 &\geq \lambda_{\min}(H^*) - \lambda/2 \\ 1977 &\geq \lambda/2 \end{aligned}$$

1978  
1979 where the first step comes from triangle inequality, the second step is due to Frobenius norm, the  
1980 third step is due to Eq. (11), the last step follows from  $\lambda_{\min}(H^*) = \lambda$ .  $\square$

## 1981 I LOSS DECOMPOSITION

1982  
1983  
1984 In this section, we provide the lemma below to decompose it into five terms, and then we will give  
1985 bounds to four terms.

1986 **Lemma I.1.** *Assuming the following condition is met:*

- 1987 • *Let  $W \in \mathbb{R}^{d \times m}$  and  $a \in \mathbb{R}^m$  as Definition F.1.*
- 1988 • *Let  $\lambda = \lambda_{\min}(H^*)$*
- 1989 • *For  $i, j \in [n]$  and  $k_1, k_2 \in [d]$ .*
- 1990 • *Let  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  be defined as Definition G.6.*
- 1991 • *Let  $u_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.2.*
- 1992 • *Denote  $F(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.*
- 1993 • *Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.*

- 1998
- 1999
- 2000
- 2001
- 2002
- 2003
- 2004
- 2005
- 2006
- 2007
- 2008
- 2009
- 2010
- 2011
- 2012
- 2013
- 2014
- 2015
- 2016
- 2017
- 2018
- 2019
- 2020
- 2021
- 2022
- 2023
- 2024
- 2025
- 2026
- 2027
- 2028
- 2029
- 2030
- 2031
- 2032
- 2033
- 2034
- 2035
- 2036
- 2037
- 2038
- 2039
- 2040
- 2041
- 2042
- 2043
- 2044
- 2045
- 2046
- 2047
- 2048
- 2049
- 2050
- 2051
- Let  $\eta > 0$  denote the learning rate.
  - Let scalar  $v_{0,k,i} \in \mathbb{R}$  be defined as follows

$$v_{0,k,i} := m \sum_{r \in [m]} (\theta_{k,i,r}(\tau+1) - \theta_{k,i,r}(\tau)) \cdot \mathbf{u}_{i,r}(\tau+1)$$

- Let scalar  $v_{1,k,i} \in \mathbb{R}$  be defined as follows

$$v_{1,k,i} := m \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot (-\eta \langle \Delta w_r(\tau), x_i \rangle)$$

- Let scalar  $v_{2,k,i} \in \mathbb{R}$  be defined as follows

$$v_{2,k,i} := m \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \eta^2 \cdot \Theta(1) \cdot \langle \Delta w_r(\tau), x_i \rangle^2$$

- **Gradient Property.**  $\eta \|\Delta w_r(i)\|_2 \leq 0.01, \forall r \in [m], \forall i \in [\tau]$
- $C_0 = 2 \langle \text{vec}(\mathbf{F}(\tau) - Y), \text{vec}(v_0) \rangle$
- $C_1 = 2 \langle \text{vec}(\mathbf{F}(\tau) - Y), \text{vec}(v_1) \rangle$
- $C_2 = 2 \langle \text{vec}(\mathbf{F}(\tau) - Y), \text{vec}(v_2) \rangle$
- $C_3 = \|\mathbf{F}(\tau+1) - \mathbf{F}(\tau)\|_F^2$

Then, we can show

$$\|\mathbf{F}(\tau+1) - Y\|_F^2 = \|\mathbf{F}(\tau) - Y\|_F^2 + C_0 + C_1 + C_2 + C_3.$$

*Proof.* The expression  $\|Y - \mathbf{F}(\tau+1)\|_F^2 = \|\text{vec}(Y - \mathbf{F}(\tau+1))\|_2^2$  can be rewritten in the following:

$$\begin{aligned} & \|\text{vec}(Y - \mathbf{F}(\tau+1))\|_2^2 \\ &= \|\text{vec}(Y - \mathbf{F}(\tau) - (\mathbf{F}(\tau+1) - \mathbf{F}(\tau)))\|_2^2 \\ &= \|\text{vec}(Y - \mathbf{F}(\tau))\|_2^2 - 2 \text{vec}(Y - \mathbf{F}(\tau))^\top \text{vec}(\mathbf{F}(\tau+1) - \mathbf{F}(\tau)) \\ & \quad + \|\text{vec}(\mathbf{F}(\tau+1) - \mathbf{F}(\tau))\|_2^2. \end{aligned} \tag{12}$$

where the first step follows from simple algebra, the last step follows from simple algebra.

Recall the update rule (Definition G.13),

$$w_r(\tau+1) = w_r(\tau) - \eta \cdot \Delta w_r(\tau)$$

In the following manner,  $\forall k \in [d]$ , we can express  $\mathbf{F}_k(\tau+1) - \mathbf{F}_k(\tau) \in \mathbb{R}^n$ :

$$\begin{aligned} & \mathbf{F}_{k,i}(\tau+1) - \mathbf{F}_{k,i}(\tau) \\ &= m \sum_{r \in [m]} (\theta_{k,i,r}(\tau+1) \mathbf{u}_{i,r}(\tau+1) - \theta_{k,i,r}(\tau) \mathbf{u}_{i,r}(\tau)) \\ &= + \sum_{r \in [m]} (\theta_{k,i,r}(\tau+1) - \theta_{k,i,r}(\tau)) \cdot \mathbf{u}_{i,r}(\tau+1) \\ & \quad + m \sum_{r \in [m]} \theta_{k,i,r} \cdot (\mathbf{u}_{i,r}(\tau+1) - \mathbf{u}_{i,r}(\tau)) \\ &= + \sum_{r \in [m]} (\theta_{k,i,r}(\tau+1) - \theta_{k,i,r}(\tau)) \cdot \mathbf{u}_{i,r}(\tau+1) \end{aligned}$$

$$\begin{aligned}
& + m \sum_{r \in [m]} \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot (\exp(-\eta \langle \Delta w_r(\tau), x_i \rangle) - 1) \\
& = + \sum_{r \in [m]} (\theta_{k,i,r}(\tau + 1) - \theta_{k,i,r}(\tau)) \cdot \mathbf{u}_{i,r}(\tau + 1) \\
& \quad + m \sum_{r \in [m]} \theta_{k,i,r}(\tau) \mathbf{u}_{i,r}(\tau) \cdot (-\eta \langle \Delta w_r(\tau), x_i \rangle + \Theta(1) \eta^2 \langle \Delta w_r(\tau), x_i \rangle^2) \\
& = v_{0,k,i} + v_{1,k,i} + v_{2,k,i} \tag{13}
\end{aligned}$$

where the first step is due to the definition of  $F_{k,i}(\tau)$ , the second step is from the simple algebra, the third step is due to  $|\eta \Delta w_r(\tau)^\top x_i| \leq 0.01$  (due to **Gradient Property** and  $\|x_i\|_2 \leq 1$ ), the fourth step follows from the Taylor series approximation, the last step follows from

$$\begin{aligned}
v_{0,k,i} & := m \sum_{r \in [m]} (\theta_{k,i,r}(\tau + 1) - \theta_{k,i,r}(\tau)) \cdot \mathbf{u}_{i,r}(\tau + 1) \\
v_{1,k,i} & := m \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot (-\eta \langle \Delta w_r(\tau), x_i \rangle) \\
v_{2,k,i} & := m \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \eta^2 \cdot \Theta(1) \cdot \langle \Delta w_r(\tau), x_i \rangle^2
\end{aligned}$$

Here  $v_{0,k,i}$  and  $v_{1,k,i}$  are linear in  $\eta$  and  $v_{2,k,i}$  is quadratic in  $\eta$ . Thus,  $v_{0,k,i}$  and  $v_{1,k,i}$  are the first order term, and  $v_{2,k,i}$  is the second order term.

We can rewrite the second term in the Eq. (12) above as below:

$$\begin{aligned}
& \langle \text{vec}(Y - F(\tau)), \text{vec}(F(\tau + 1) - F(\tau)) \rangle \\
& = \langle \text{vec}(Y - F(\tau)), \text{vec}(v_0 + v_1 + v_2) \rangle \\
& = \langle \text{vec}(Y - F(\tau)), \text{vec}(v_0) \rangle + \langle \text{vec}(Y - F(\tau)), \text{vec}(v_1) \rangle + \langle \text{vec}(Y - F(\tau)), \text{vec}(v_2) \rangle
\end{aligned}$$

where the first step follows from Eq.(13), the second step follows from simple algebras.

Therefore, we can conclude that

$$\|F(\tau + 1) - Y\|_F^2 = \|F(\tau) - Y\|_F^2 + C_0 + C_1 + C_2 + C_3.$$

□

The below lemma analyzes the first-order term that is making progress.

**Lemma I.2** (Progress terms). *If the following conditions hold*

- Let  $\lambda = \lambda_{\min}(H^*)$
- Let  $C > 10$  denote a sufficiently large constant
- Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- Let  $\delta \in (0, 0.1)$ .
- Let  $m \geq \Omega(\lambda^{-2}n^2d^2 \exp(30B)\sqrt{\log(nd/\delta)})$
- Let  $r \in [m]$ , let  $i, j \in [n]$ , let  $k, k_2 \in [d]$ .
- Let  $\beta_k(\tau) \in \mathbb{R}^m$  be defined as Definition F.5.
- Let  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  be defined as Definition G.6.
- Let  $\mathbf{u}_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.2.
- Let  $S_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.7.
- Let  $v_{k,r} := \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau) \in \mathbb{R}^m$

- Denote  $F(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.
- Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.
- Let  $\eta > 0$  denote the learning rate.
- Let scalar  $v_{1,1,k,i} \in \mathbb{R}$  be defined as follows (we omit  $(\tau)$  in the following terms)

$$v_{1,1,k,i} = m^2 \sum_{r \in [m]} \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (F_{k_2,j}(\tau) - y_{k_2,j}) \cdot \left( \langle v_{k_2,r}, \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \right) \cdot x_j^\top \right) x_i$$

- Let  $C_{1,1} := 2 \langle \text{vec}(F(\tau) - Y), \text{vec}(v_{1,1}) \rangle$

Then, we have

- $C_{1,1} \leq -1.6m\eta \text{vec}(F(\tau) - Y)^\top H(\tau) \text{vec}(F(\tau) - Y)$

*Proof.* We have

$$\begin{aligned} v_{1,1,k,i} &= m^2 \sum_{r \in [m]} \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \\ &\quad \cdot \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (F_{k_2,j}(\tau) - y_{k_2,j}) \cdot \left( \langle v_{k_2,r}, \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \right) \cdot x_j^\top \right) x_i \\ &= m^2 \sum_{r \in [m]} \beta_{k,r}(\tau) \cdot \alpha_i(\tau)^{-1} \cdot \mathbf{u}_{i,r}(\tau) \\ &\quad \cdot \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (F_{k_2,j}(\tau) - y_{k_2,j}) \cdot \left( \langle v_{k_2,r}, \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \right) \cdot x_j^\top \right) x_i \\ &= m^2 \sum_{r \in [m]} \beta_{k,r}(\tau) \cdot \mathbf{S}_{i,r}(\tau) \\ &\quad \cdot \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (F_{k_2,j}(\tau) - y_{k_2,j}) \cdot \left( \langle v_{k_2,r}, \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \right) \cdot x_j^\top \right) x_i \\ &= m^2 \sum_{r \in [m]} \langle \beta_{k,r}(\tau) \cdot \mathbf{1}_m, \mathbf{S}_i(\tau) \rangle \cdot \mathbf{S}_{i,r}(\tau) \\ &\quad \cdot \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (F_{k_2,j}(\tau) - y_{k_2,j}) \cdot \left( \langle v_{k_2,r}, \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \right) \cdot x_j^\top \right) x_i \\ &= m^2 \sum_{r \in [m]} \left( \langle v_{k,r}, \mathbf{S}_i(\tau) \rangle + \langle \beta_k(\tau), \mathbf{S}_i(\tau) \rangle \right) \cdot \mathbf{S}_{i,r}(\tau) \\ &\quad \cdot \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (F_{k_2,j}(\tau) - y_{k_2,j}) \cdot \left( \langle v_{k_2,r}, \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \right) \cdot x_j^\top \right) x_i \\ &= m^2 (Q_{1,1,k,i} + Q_{1,2,k,i}) \end{aligned}$$

where the first step follows from the definition of  $v_{1,1,k,i}$ , the second step follows from Definition G.6, the third step follows from Definition G.7, the fourth step follows from  $\langle \beta_{k,r}(\tau) \cdot \mathbf{1}_m, \mathbf{S}_i \rangle = \beta_{k,r}(\tau)$ , the fifth step follows from the definition of  $v_k$  for  $k \in [d]$  and simple algebras, the last step holds since we define

$$Q_{1,1,k,i} := \sum_{r \in [m]} \langle v_{k,r}, \mathbf{S}_i(\tau) \rangle \cdot \mathbf{S}_{i,r}(\tau)$$



$$\begin{aligned}
& \cdot (-\eta \sum_{j=1}^n \sum_{k_2=1}^d (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \cdot \left( (\langle v_{k_2,r}, \mathbf{S}_j(\tau) \rangle) \cdot \mathbf{S}_{j,r}(\tau) \right) \cdot x_j^\top) x_i, \\
Q_{1,2,k,i} & := \sum_{r \in [m]} \langle \beta_k(\tau), \mathbf{S}_i(\tau) \rangle \cdot \mathbf{S}_{i,r}(\tau) \\
& \cdot (-\eta \sum_{j=1}^n \sum_{k_2=1}^d (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \cdot \left( (\langle v_{k_2,r}, \mathbf{S}_j(\tau) \rangle) \cdot \mathbf{S}_{j,r}(\tau) \right) \cdot x_j^\top) x_i.
\end{aligned}$$

**Bounding first term.** Then for the first term  $Q_{1,1,k,i}$ , we have its quantity

$$\sum_{i=1}^n \sum_{k=1}^d Q_{1,1,k,i} (\mathbf{F}_{k,i}(\tau) - y_{k,i}) = -\frac{1}{m} \eta \text{vec}(\mathbf{F}(\tau) - Y)^\top H(\tau) \text{vec}(\mathbf{F}(\tau) - Y)$$

where this step follows from Definition H.1.

**Bounding second term.** On the other hand, for the second term  $Q_{1,2,k,i}$ , we have its quantity,

$$\begin{aligned}
& \left| \sum_{i=1}^n \sum_{k=1}^d Q_{1,2,k,i} (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \eta \left| \frac{\exp(9B)}{m^3} \sum_{i=1}^n \sum_{j=1}^n \sum_{r=1}^m \sum_{k=1}^d \sum_{k_2=1}^d \sigma_r C_{k,k_2,r} (\mathbf{F}_{k,i}(\tau) - y_{k,i}) (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \right| \\
& \leq \eta \frac{\exp(9B)}{m^3} \cdot \left| \sum_{r=1}^m \sigma_r \max_{k,k_2 \in [d]} C_{k,k_2,r} \cdot \|(\mathbf{F}(\tau) - Y) \otimes (\mathbf{F}(\tau) - Y)\|_1 \right| \\
& \leq \eta \frac{\exp(9B)}{m^3} \cdot \left| \sum_{r=1}^m \sigma_r \max_{k,k_2 \in [d]} C_{k,k_2,r} \cdot \|\mathbf{F}(\tau) - Y\|_1^2 \right| \\
& \leq \eta \frac{nd \exp(9B)}{m^3} \cdot \left| \sum_{r=1}^m \sigma_r \max_{k,k_2 \in [d]} C_{k,k_2,r} \cdot \|\mathbf{F}(\tau) - Y\|_F^2 \right| \\
& \leq \eta \frac{\exp(9B)}{m^3 \lambda} \left| \sum_{r=1}^m \sigma_r \max_{k,k_2 \in [d]} C_{k,k_2,r} \cdot \text{vec}(\mathbf{F}(\tau) - Y)^\top H(\tau) \text{vec}(\mathbf{F}(\tau) - Y) \right|
\end{aligned}$$

where the first step follows from  $0 \leq \mathbf{S}_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1,  $\|\mathbf{S}_i\|_2 \leq \frac{\exp(3B)}{\sqrt{m}}$ ,  $\|x_i\| \leq 1$  and

$$C_{k,k_2,r} := \|\beta_k(\tau)\|_2 \cdot \|v_{k_2,r}\|_2, \sigma_r \in \{+1, -1\}$$

the second and third steps follow from the definition of Kronecker product, the fourth step follows from  $\|U\|_1 \leq \sqrt{nd} \|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ , the last step follows from  $\text{vec}(\mathbf{F}(\tau) - Y)^\top H(\tau) \text{vec}(\mathbf{F}(\tau) - Y) \geq \lambda \|\mathbf{F} - Y\|_F^2$ .

Thus, by following Part 2 and Part 3 of Lemma L.2, we have

$$C_{k,k_2,r} = \|\beta_k(\tau)\|_2 \cdot \|v_{k_2,r}\|_2 \leq 2mB^2.$$

Besides, we apply Hoeffding inequality (Lemma E.4) to all random variables  $\sigma_r \max_{k,k_2 \in [d]} C_{k,k_2,r}$  for  $r \in [m]$ , especially  $\mathbb{E}[\sum_{r=1}^m \sigma_r \max_{k,k_2 \in [d]} C_{k,k_2,r}] = 0$  due to the symmetry of  $a_r$ , we have

$$\begin{aligned}
& \left| \sum_{i=1}^n \sum_{k=1}^d Q_{1,2,k,i} (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq C \eta \frac{nd \exp(9B)}{m^3 \lambda} \cdot \text{vec}(\mathbf{F}(\tau) - Y)^\top H(\tau) \text{vec}(\mathbf{F}(\tau) - Y) \cdot mB^2 \sqrt{m \log(nd/\delta)}
\end{aligned}$$

with probability at least  $1 - \delta / \text{poly}(nd)$ .

Note that by Lemma condition, we have

$$C \frac{nd \exp(9B)}{m^3 \lambda} \cdot mB^2 \sqrt{m \log(nd/\delta)} \leq 0.2 \frac{1}{m}.$$

Finally, we complete the proof with the result

$$C_{1,1} \leq -1.6m\eta \text{vec}(\mathbf{F}(\tau) - Y)^\top H(\tau) \text{vec}(\mathbf{F}(\tau) - Y)$$

□

Below, we prove all other terms are small when  $m$  is large enough compared to the progressive term.

**Lemma I.3** (Minor effects on non-progress term). *If the following*

- Let  $m \geq \Omega(\lambda^{-2} n^2 d^2 \exp(30B) \sqrt{\log(nd/\delta)})$ .
- Let  $r \in [m]$ , let  $i, j \in [n]$ , let  $k, k_2 \in [d]$
- Let scalar  $v_{0,k,i} \in \mathbb{R}$  be defined as follows

$$v_{0,k,i} := m \sum_{r \in [m]} (\theta_{k,i,r}(\tau+1) - \theta_{k,i,r}(\tau)) \cdot \mathbf{u}_{i,r}(\tau+1)$$

- Let scalar  $v_{1,2,k,i} \in \mathbb{R}$  be defined as follows (we omit  $(\tau)$  in the following terms)

$$v_{1,2,k,i} = m^2 \sum_{r \in [m]} \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \cdot a_r \mathbf{S}_{j,r}(\tau) e_{k_2}^\top \right) x_i$$

- Let scalar  $v_{2,k,i} \in \mathbb{R}$  be defined as follows

$$v_{2,k,i} := m \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \eta^2 \cdot \Theta(1) \cdot \langle \Delta w_r(\tau), x_i \rangle^2$$

- Let  $C_0 := 2 \langle \text{vec}(\mathbf{F}(\tau) - Y), \text{vec}(v_0) \rangle$
- Let  $C_{1,2} := 2 \langle \text{vec}(\mathbf{F}(\tau) - Y), \text{vec}(v_{1,2}) \rangle$
- Let  $C_2 := 2 \langle \text{vec}(\mathbf{F}(\tau) - Y), \text{vec}(v_2) \rangle$
- Let  $C_3 := \|\mathbf{F}(\tau+1) - \mathbf{F}(\tau)\|_F^2$

Then, we have

- $|C_0| \leq 0.1m\eta\lambda \cdot \|\mathbf{F}(\tau) - Y\|_F^2$
- $|C_{1,2}| \leq 0.1m\eta\lambda \cdot \|\mathbf{F}(\tau) - Y\|_F^2$
- $|C_2| \leq \eta^2 m \cdot n^2 d^2 \exp(16B) \cdot \|\mathbf{F}(\tau) - Y\|_F^2$
- $|C_3| \leq \eta^2 m^2 \cdot \|\mathbf{F}(\tau) - Y\|_F^2$

*Proof.* This proof follows from Lemma I.4, Lemma I.5, Lemma I.6 and Lemma I.7. □

## I.1 BOUNDING $C_0$

**Lemma I.4.** *If the following conditions hold*

- Let  $\lambda = \lambda_{\min}(H^*)$
- Let  $C > 10$  denote a sufficiently large constant

- 2268 • Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .  
 2269  
 2270 • Let  $\delta \in (0, 0.1)$ .  
 2271  
 2272 • Let  $m \geq \Omega(\lambda^{-2}n^2d^2 \exp(30B)\sqrt{\log(nd/\delta)})$ .  
 2273  
 2274 • Let  $r \in [m]$ , let  $i, j \in [n]$ , let  $k, k_1 \in [d]$ .  
 2275  
 2276 • Let  $\beta_k(\tau) \in \mathbb{R}^m$  be defined as Definition F.5.  
 2277  
 2278 • Let  $\alpha_i(\tau) \in \mathbb{R}$  be defined as Definition F.3.  
 2279  
 2280 • Let  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  be defined as Definition G.6.  
 2281  
 2282 • Let  $\mathbf{u}_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.2.  
 2283  
 2284 • Let  $\mathbf{v}_k := \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau) \in \mathbb{R}^m$ .  
 2285  
 2286 • Denote  $\mathbf{F}(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.  
 2287  
 2288 • Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.  
 2289  
 2290 • Let  $\eta \in (0, 1/m)$  denote the learning rate.  
 2291  
 2292 • Let scalar  $v_{0,k,i} \in \mathbb{R}$  be defined as follows (we omit  $(\tau)$  in the following terms)

$$v_{0,k,i} = m \sum_{r \in [m]} (\theta_{k,i,r}(\tau+1) - \theta_{k,i,r}(\tau)) \cdot \mathbf{u}_{i,r}(\tau+1)$$

- 2294 • Let  $C_0 := 2\langle \text{vec}(\mathbf{F}(\tau) - Y), \text{vec}(v_0) \rangle$

2296 Then, with a probability at least  $1 - \delta/\text{poly}(nd)$ , we have

$$|C_0| \leq 0.1\eta m \lambda \|\mathbf{F}(\tau) - Y\|_{\mathbb{F}}^2.$$

2299 *Proof.* By Claim G.12, we have

$$\Delta w_r(\tau) = m \sum_{i=1}^n \sum_{k=1}^d (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( \langle v_{k,r}(\tau), \mathbf{S}_i(\tau) \rangle \cdot \mathbf{S}_{i,r}(\tau) \cdot x_i + a_r \mathbf{S}_{i,r}(\tau) e_k \right)$$

2304 Then the  $k_1$ -th entry  $\Delta w_{r,k_1}(\tau)$  for  $k_1 \in [d]$  should be

$$\Delta w_{r,k_1}(\tau) = m \sum_{i=1}^n \sum_{k=1}^d (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( \langle v_{k,r}(\tau), \mathbf{S}_i(\tau) \rangle \cdot \mathbf{S}_{i,r}(\tau) \cdot x_{i,k_1} + a_r \mathbf{S}_{i,r}(\tau) e_{k,k_1} \right) \quad (14)$$

2310 We have

$$\begin{aligned} v_{0,k,i} &= m \sum_{r \in [m]} (\theta_{k,i,r}(\tau+1) - \theta_{k,i,r}(\tau)) \cdot \mathbf{u}_{i,r}(\tau+1) \\ &= m \sum_{r \in [m]} (\beta_{k,r}(\tau+1)\alpha_i(\tau+1)^{-1} - \beta_{k,r}(\tau)\alpha_i(\tau)^{-1}) \cdot \mathbf{u}_{i,r}(\tau+1) \\ &= m \sum_{r \in [m]} (\beta_{k,r}(\tau+1)\alpha_i(\tau+1)^{-1} - \beta_{k,r}(\tau+1)\alpha_i(\tau)^{-1} \\ &\quad + \beta_{k,r}(\tau+1)\alpha_i(\tau)^{-1} - \beta_{k,r}(\tau)\alpha_i(\tau)^{-1}) \cdot \mathbf{u}_{i,r}(\tau+1) \\ &= m \sum_{r \in [m]} (\beta_{k,r}(\tau+1) \cdot (\alpha_i(\tau+1)^{-1} - \alpha_i(\tau)^{-1}) \end{aligned}$$

$$\begin{aligned}
& + (\beta_{k,r}(\tau + 1) - \beta_{k,r}(\tau)) \cdot \alpha_i(\tau)^{-1} \cdot \mathbf{u}_{i,r}(\tau + 1) \\
& = m(Q_{0,1,k,i} + Q_{0,2,k,i})
\end{aligned}$$

where the first step follows from the definition of  $v_{0,k,i}$ , the second step follows from Definition G.6, the third and fourth steps follow from simple algebras, the last step hold since we define

$$\begin{aligned}
Q_{0,1,k,i} & := \sum_{r \in [m]} \beta_{k,r}(\tau + 1) \cdot (\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1}) \cdot \mathbf{u}_{i,r}(\tau + 1), \\
Q_{0,2,k,i} & := \sum_{r \in [m]} (\beta_{k,r}(\tau + 1) - \beta_{k,r}(\tau)) \cdot \alpha_i(\tau)^{-1} \cdot \mathbf{u}_{i,r}(\tau + 1).
\end{aligned}$$

**Bounding first term.** For the first term  $Q_{0,1,k,i}$ , we have its quantity

$$\begin{aligned}
& \left| \sum_{i=1}^n \sum_{k=1}^d Q_{0,1,k,i}(\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m \beta_{k,r}(\tau + 1) \cdot (\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1}) \cdot \mathbf{u}_{i,r}(\tau + 1) (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \exp(B) \cdot \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m \beta_{k,r}(\tau + 1) \cdot (\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1}) (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq B \exp(B) \cdot \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m a_r (\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1}) \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq B \exp(B) \cdot \left| \sum_{r=1}^m a_r (\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1}) \right| \cdot \sqrt{nd} \|\mathbf{F}(\tau) - Y\|_F \tag{15}
\end{aligned}$$

where the first step follows from the definition of  $Q_{0,1,k,i}$ , the second step follows from Part 4 of Lemma L.1 and Definition G.2, the third step follows from Part 1 of Lemma L.1 and  $\|U\|_1 \leq \sqrt{nd} \|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ .

By Part 2 of Lemma I.9, we have

$$\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1} \leq \eta \frac{\sqrt{nd} \exp(15B)}{m^3} \cdot \|\mathbf{F}(\tau) - Y\|_F + \eta^2 \frac{nd \exp(27B)}{\sqrt{m}} \cdot \|\mathbf{F}(\tau) - Y\|_F.$$

Then we apply Hoeffding inequality (Lemma E.4) to random variables  $a_r (\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1})$  for  $r \in [m]$ , and by  $\mathbb{E}[\sum_{r=1}^m a_r (\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1})] = 0$ , we have

$$\begin{aligned}
& \left| \sum_{r=1}^m a_r (\alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1}) \right| \\
& \leq \left( \eta \frac{\sqrt{nd} \exp(15B)}{m^3} + \eta^2 \frac{nd \exp(27B)}{\sqrt{m}} \right) \cdot \|\mathbf{F}(\tau) - Y\|_F \cdot \sqrt{m \log(nd/\delta)}. \tag{16}
\end{aligned}$$

with probability at least  $1 - \delta / \text{poly}(nd)$ .

Through combining Eq. (16) and Eq.(15), we can show that

$$\begin{aligned}
& \left| \sum_{i=1}^n \sum_{k=1}^d Q_{0,1,k,i}(\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \left( \eta \frac{nd \exp(17B)}{m^3} \cdot \|\mathbf{F}(\tau) - Y\|_F^2 + \eta^2 \frac{nd \sqrt{nd} \exp(29B)}{\sqrt{m}} \cdot \|\mathbf{F}(\tau) - Y\|_F^2 \right) \cdot \sqrt{m \log(nd/\delta)}
\end{aligned}$$

with a probability at least  $1 - \delta / \text{poly}(nd)$ .

Thus, by Lemma condition, we can show

$$\eta \frac{nd \exp(17B)}{m^3} \cdot \sqrt{m \log(nd/\delta)} \leq 0.01\eta\lambda,$$

$$\eta^2 \frac{nd\sqrt{nd} \exp(29B)}{\sqrt{m}} \cdot \sqrt{m \log(nd/\delta)} \leq \eta \frac{nd\sqrt{nd} \exp(29B)}{m} \cdot \sqrt{\log(nd/\delta)} \leq 0.01\eta\lambda.$$

**Bounding second term.** On the other hand, for the second term  $Q_{0,2,k,i}$ , we have its quantity

$$\begin{aligned} & \left| \sum_{i=1}^n \sum_{k=1}^d Q_{0,2,k,i}(\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\ & \leq \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m (\beta_{k,r}(\tau+1) - \beta_{k,r}(\tau)) \cdot \alpha_i(\tau)^{-1} \cdot \mathbf{u}_{i,r}(\tau+1) \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\ & \leq \exp(B) \cdot \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m (\beta_{k,r}(\tau+1) - \beta_{k,r}(\tau)) \cdot \alpha_i(\tau)^{-1} \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\ & \leq \frac{\exp(2B)}{m} \cdot \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m (\beta_{k,r}(\tau+1) - \beta_{k,r}(\tau)) \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\ & \leq \frac{\exp(2B)}{m} \cdot \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m (W_{k,r}(\tau+1) \cdot a_r - W_{k,r}(\tau) \cdot a_r) \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\ & \leq \eta \frac{\exp(2B)}{m} \cdot \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m a_r \cdot m \cdot \sum_{j=1}^d (\mathbf{F}_{k_1,j}(\tau) - y_{k_1,j}) \right. \\ & \quad \cdot \left. \left( \langle v_{k_1,r}(\tau), \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \cdot x_{j,k} + a_r \mathbf{S}_{j,r}(\tau) e_{k_1,k} \right) \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\ & \leq \eta \frac{\exp(5B)}{m} \cdot \left| \sum_{r=1}^m \sigma_r \max_{j,k,k_1 \in [d]} C_{j,k,k_1,r} \cdot \|(\mathbf{F}(\tau) - Y) \otimes (\mathbf{F}(\tau) - Y)\|_1 \right| \\ & \leq \eta \frac{\exp(5B)}{m} \cdot \left| \sum_{r=1}^m \sigma_r \max_{j,k,k_1 \in [d]} C_{j,k,k_1,r} \cdot \|\mathbf{F}(\tau) - Y\|_1^2 \right| \\ & \leq \eta \frac{nd \exp(5B)}{m} \cdot \left| \sum_{r=1}^m \sigma_r \max_{j,k,k_1 \in [d]} C_{j,k,k_1,r} \cdot \|\mathbf{F}(\tau) - Y\|_F^2 \right| \end{aligned}$$

where the first step follows from the definition of  $Q_{0,2,k,i}$ , the second and third steps follow from Part 4 of Lemma L.1, the fourth step follows from Definition F.5, the fifth step follows from Eq.(14), the sixth step follows from the definition of Kronecker product,  $1 \leq \mathbf{S}_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1,  $\|x_i\|_2 \leq 1$  and defining

$$C_{j,k,k_1,r} := \langle \mathbf{S}_j, v_{k_1,r} \rangle + e_{k_1,k}, \sigma_r \in \{+1, -1\},$$

the seventh step follows from the definition of  $\ell_1$  norm, the last step follows from  $\|U\|_1 \leq \sqrt{nd}\|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ .

Thus, by following Part 6 of Lemma L.2, we have

$$\begin{aligned} C_{j,k,k_1,r} &= \langle \mathbf{S}_j, v_{k_1,r} \rangle + e_{k_1,k} \\ &\leq \exp(6B) + 1 \\ &\leq \exp(7B) \end{aligned}$$

where the last step follows from simple algebras.

We apply Hoeffding inequality (Lemma E.4) to  $\sigma_r \max_{j,k,k_1 \in [d]} C_{j,k,k_1,r}$  for  $r \in [m]$ .

By  $\mathbb{E}[\sum_{r=1}^m \sigma_r \max_{j,k,k_1 \in [d]} C_{j,k,k_1,r}] = 0$ , we have

$$\left| \sum_{i=1}^n \sum_{k=1}^d Q_{0,2,k,i}(\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \leq \eta \frac{nd \exp(5B)}{m} \cdot \|\mathbf{F}(\tau) - Y\|_F^2 \cdot \exp(6B) \sqrt{m \log(nd/\delta)}.$$

with probability at least  $1 - \delta / \text{poly}(nd)$ .

Then, by Lemma condition, we have

$$\eta \frac{nd \exp(5B)}{m} \cdot \exp(7B) \sqrt{m \log(nd/\delta)} \leq 0.01\eta\lambda.$$

Now we can complete the proof by combining all terms, we have

$$|C_0| \leq 0.1\eta m \lambda \|F(\tau) - Y\|_F^2.$$

□

## I.2 BOUNDING $C_{1,2}$

**Lemma I.5.** *If the following conditions hold*

- Let  $\lambda = \lambda_{\min}(H^*)$
- Let  $C > 10$  denote a sufficiently large constant
- Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- Let  $\delta \in (0, 0.1)$ .
- Let  $m \geq \Omega(\lambda^{-2}n^2d^2 \exp(30B)\sqrt{\log(nd/\delta)})$ .
- Let  $r \in [m]$ , let  $i, j \in [n]$ , let  $k, k_1 \in [d]$ .
- Let  $\beta_k(\tau) \in \mathbb{R}^m$  be defined as Definition F.5.
- Let  $\alpha_i(\tau) \in \mathbb{R}$  be defined as Definition F.3.
- Let  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  be defined as Definition G.6.
- Let  $\mathbf{u}_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.2.
- Let  $\mathbf{S}_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.7.
- Let  $\mathbf{v}_k := \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau) \in \mathbb{R}^m$
- Denote  $F(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.
- Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.
- Let  $\eta > 0$  denote the learning rate.
- Let scalar  $v_{1,2,k,i} \in \mathbb{R}$  be defined as follows (we omit  $(\tau)$  in the following terms)

$$v_{1,2,k,i} = m^2 \sum_{r \in [m]} \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \cdot \mathbf{a}_\tau \mathbf{S}_{j,r}(\tau) e_{k_2}^\top \right) x_i$$

- Let  $C_{1,2} := 2 \langle \text{vec}(F(\tau) - Y), \text{vec}(v_{1,2}) \rangle$

Then, with a probability at least  $1 - \delta / \text{poly}(nd)$ , we have

$$|C_{1,2}| \leq 0.1\eta m \lambda \|F(\tau) - Y\|_F^2$$

*Proof.* We have the quantity of  $v_{1,2,k,i}$

$$\left| \sum_{i=1}^n \sum_{k=1}^d v_{1,2,k,i} (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right|$$

$$\begin{aligned}
&\leq \left| \sum_{i=1}^n \sum_{k=1}^d m^2 \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \right. \\
&\quad \cdot \left. \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \cdot a_r \mathbf{S}_{j,r}(\tau) e_{k_2}^\top \right) x_i \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
&\leq \left| \sum_{i=1}^n \sum_{k=1}^d m^2 \sum_{r=1}^m \beta_{k,r}(\tau) \alpha_i(\tau)^{-1} \cdot \mathbf{u}_{i,r}(\tau) \right. \\
&\quad \cdot \left. \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \cdot a_r \mathbf{S}_{j,r}(\tau) e_{k_2}^\top \right) x_i \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
&\leq \left| \sum_{i=1}^n \sum_{k=1}^d m^2 \sum_{r=1}^m \beta_{k,r}(\tau) \mathbf{S}_{i,r}(\tau) \right. \\
&\quad \cdot \left. \left( -\eta \sum_{j=1}^n \sum_{k_2=1}^d (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \cdot a_r \mathbf{S}_{j,r}(\tau) e_{k_2}^\top \right) x_i \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
&\leq \eta m^2 \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m \beta_{k,r}(\tau) \mathbf{S}_{i,r}(\tau) \right. \\
&\quad \cdot \left. \left( -\sum_{j=1}^n \sum_{k_2=1}^d (\mathbf{F}_{k_2,j}(\tau) - y_{k_2,j}) \cdot a_r \mathbf{S}_{j,r}(\tau) e_{k_2}^\top \right) x_i \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
&\leq \eta \exp(6B) \sum_{r=1}^m |a_r \cdot \max_{k \in [d]} \beta_{k,r}(\tau)| \cdot \|(\mathbf{F}(\tau) - Y) \otimes (\mathbf{F}(\tau) - Y)\|_1 \\
&\leq \eta \exp(6B) \sum_{r=1}^m |a_r \cdot \max_{k \in [d]} \beta_{k,r}(\tau)| \cdot \|\mathbf{F}(\tau) - Y\|_1^2 \\
&\leq \eta n d \exp(6B) \sum_{r=1}^m |a_r \cdot \max_{k \in [d]} \beta_{k,r}(\tau)| \cdot \|\mathbf{F}(\tau) - Y\|_F^2
\end{aligned}$$

where the first step follows from the definition of  $v_{1,2,k,i}$ , the second step follows from Definition G.6, the third step follows from Definition F.5, the fourth step follows from Definition G.7, the fifth step follows from simple algebras, the sixth step follows from  $0 \leq \mathbf{S}_{j,r} \leq \frac{\exp(3B)}{m}$ ,  $\|x_i\|_2 \leq 1$  and the definition of Kronecker product, the seventh step follows from the definition of  $\ell_1$  norm, the last step follows from  $\|U\|_1 \leq \sqrt{nd} \|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ .

Then by Part 1 of Lemma L.1, we have

$$|\max_{k \in [d]} \beta_{k,r}(\tau)| \leq B$$

We apply Hoeffding inequality (Lemma E.4) to random variables  $a_r \cdot \max_{k \in [d]} \beta_{k,r}(\tau)$  for  $r \in [m]$ . By  $\mathbb{E}[\sum_{r=1}^m a_r \cdot \max_{k \in [d]} \beta_{k,r}(\tau)] = 0$ , we have

$$\left| \sum_{i=1}^n \sum_{k=1}^d v_{1,2,k,i} (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \leq \eta n d \exp(6B) B \|\mathbf{F}(\tau) - Y\|_F^2$$

with a probability at least  $1 - \delta / \text{poly}(nd)$ .

By the Lemma condition, we have

$$n d \exp(6B) B \leq 0.1 m \lambda$$

□

2538 I.3 BOUNDING  $C_2$   
25392540 **Lemma I.6.** *If the following conditions hold*

- 2541 • Let  $\lambda = \lambda_{\min}(H^*)$
- 2542 • Let  $C > 10$  denote a sufficiently large constant
- 2543 • Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- 2544 • Let  $\delta \in (0, 0.1)$ .
- 2545 • Let  $m \geq \Omega(\lambda^{-2}n^2d^2 \exp(30B)\sqrt{\log(nd/\delta)})$ .
- 2546 • Let  $r \in [m]$ , let  $i, j \in [n]$ , let  $k, k_1 \in [d]$ .
- 2547 • Let  $\beta_k(\tau) \in \mathbb{R}^m$  be defined as Definition F.5.
- 2548 • Let  $\alpha_i(\tau) \in \mathbb{R}$  be defined as Definition F.3.
- 2549 • Let  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  be defined as Definition G.6.
- 2550 • Let  $\mathbf{u}_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.2.
- 2551 • Let  $\mathbf{S}_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.7.
- 2552 • Let  $v_k := \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau) \in \mathbb{R}^m$
- 2553 • Denote  $\mathbf{F}(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.
- 2554 • Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.
- 2555 • Let  $\eta > 0$  denote the learning rate.
- 2556 • Let scalar  $v_{2,k,i} \in \mathbb{R}$  be defined as follows (we omit  $(\tau)$  in the following terms)

$$2567 v_{2,k,i} := m \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \eta^2 \cdot \Theta(1) \cdot \langle \Delta w_r(\tau), x_i \rangle^2$$

- 2570 • Let  $C_2 := 2\langle \text{vec}(\mathbf{F}(\tau) - Y), \text{vec}(v_2) \rangle$

2571 *Then, with a probability at least  $1 - \delta/\text{poly}(nd)$ , we have*

$$2572 |C_2| \leq \eta^2 m \cdot n^2 d^2 \exp(16B) \|\mathbf{F}(\tau) - Y\|_F^2$$

2573 *Proof.* We have

$$2574 \begin{aligned} & \langle \Delta w_r(\tau), x_i \rangle^2 \\ & \leq \left( m \sum_{j=1}^n \sum_{k=1}^d (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( \langle v_{k,r}(\tau), \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \cdot x_j^\top + a_r \mathbf{S}_{j,r}(\tau) e_k^\top \right) x_i \right)^2 \\ & \leq \exp(12B) \cdot \|\mathbf{F}(\tau) - Y\|_1^2 \\ & \leq nd \exp(12B) \cdot \|\mathbf{F}(\tau) - Y\|_F^2 \end{aligned} \tag{17}$$

2585 where the first step follows from Claim G.12, the second step follows from the definition of  $\ell_1$  norm,  
 2586  $0 \leq \mathbf{S}_{j,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1 and Part 6 of Lemma L.2, last step follows from  
 2587  $\|U\|_1 \leq \sqrt{nd}\|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ .

2588 Then, we can show that

$$2589 \left| \sum_{i=1}^n \sum_{k=1}^d v_{2,k,i} (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right|$$



$$\begin{aligned}
& \leq \left| \sum_{i=1}^n \sum_{k=1}^d m \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \eta^2 \cdot \Theta(1) \cdot \langle \Delta w_r(\tau), x_i \rangle^2 \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \eta^2 \left| \sum_{i=1}^n \sum_{k=1}^d m \sum_{r=1}^m \theta_{k,i,r}(\tau) \cdot \mathbf{u}_{i,r}(\tau) \cdot \langle \Delta w_r(\tau), x_i \rangle^2 \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \eta^2 \left| \sum_{i=1}^n \sum_{k=1}^d m \sum_{r=1}^m \beta_{k,r}(\tau) \cdot \alpha_i(\tau)^{-1} \cdot \mathbf{u}_{i,r}(\tau) \cdot \langle \Delta w_r(\tau), x_i \rangle^2 \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \eta^2 \left| \sum_{i=1}^n \sum_{k=1}^d m \sum_{r=1}^m \beta_{k,r}(\tau) \cdot \mathbf{S}_{i,r}(\tau) \cdot \langle \Delta w_r(\tau), x_i \rangle^2 \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \eta^2 \exp(3B) \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m \beta_{k,r}(\tau) \cdot \langle \Delta w_r(\tau), x_i \rangle^2 \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \eta^2 \exp(4B) \left| \sum_{i=1}^n \sum_{k=1}^d \sum_{r=1}^m a_r \langle \Delta w_r(\tau), x_i \rangle^2 \cdot (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \\
& \leq \eta^2 \exp(4B) \left| \sum_{r=1}^m a_r \max_{i \in [n]} \langle \Delta w_r(\tau), x_i \rangle^2 \cdot \sqrt{nd} \|\mathbf{F}(\tau) - Y\|_F \right| \\
& \leq \eta^2 \sqrt{mnd} \exp(4B) \left| \sum_{r=1}^m a_r \max_{i \in [n]} \langle \Delta w_r(\tau), x_i \rangle^2 \right|
\end{aligned}$$

where the first step follows from the definition of  $v_{2,k,i}$ , the second step follows from simple algebras, the third step follows from Definition G.6, the fourth step follows from Definition G.7, the fifth step follows from  $0 \leq \mathbf{S}_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1, the sixth step follows from Part 1 of Lemma L.1 and Definition F.5, the seventh step follows from definition of  $\ell_1$  norm and  $\|U\|_1 \leq \sqrt{nd} \|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ , the last step follows from Lemma I.8.

Next, by Eq.(17), applying Hoeffding inequality (Lemma E.4) to  $a_r \max_{i \in [n]} \langle \Delta w_r(\tau), x_i \rangle^2$  for  $r \in [m]$  and  $\mathbb{E}[\sum_{r=1}^m a_r \max_{i \in [n]} \langle \Delta w_r(\tau), x_i \rangle^2] = 0$ , we have

$$\left| \sum_{i=1}^n \sum_{k=1}^d v_{2,k,i} (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \right| \leq \eta^2 \sqrt{mn} d^2 \exp(16B) \cdot \|\mathbf{F}(\tau) - Y\|_F^2 \cdot \sqrt{m \log(nd/\delta)}$$

with a probability at least  $1 - \delta / \text{poly}(nd)$ .

By the Lemma condition, we have

$$\eta^2 \sqrt{mn} d^2 \exp(16B) \cdot \sqrt{m \log(nd/\delta)} \leq \eta^2 m \cdot n^2 d^2 \exp(16B)$$

Then we complete the proof.  $\square$

#### I.4 BOUNDING $C_3$

**Lemma I.7.** *If the following conditions hold*

- Let  $\lambda = \lambda_{\min}(H^*)$
- Let  $C > 10$  denote a sufficiently large constant
- Let  $B := \max\{C\sigma \sqrt{\log(nd/\delta)}, 1\}$ .
- Let  $\delta \in (0, 0.1)$ .
- Let  $m \geq \Omega(\lambda^{-2} n^2 d^2 \exp(30B) \sqrt{\log(nd/\delta)})$ .
- Let  $r \in [m]$ , let  $i, j \in [n]$ , let  $k, k_1 \in [d]$ .

- Let  $\beta_k(\tau) \in \mathbb{R}^m$  be defined as Definition F.5.
- Let  $\alpha_i(\tau) \in \mathbb{R}$  be defined as Definition F.3.
- Let  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  be defined as Definition G.6.
- Let  $\mathbf{u}_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.2.
- Let  $\mathbf{S}_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.7.
- Let  $\mathbf{v}_k := \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau) \in \mathbb{R}^m$
- Denote  $\mathbf{F}(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.
- Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.
- Let  $\eta > 0$  denote the learning rate.
- Let  $C_3 := \|\mathbf{F}(\tau + 1) - \mathbf{F}(\tau)\|_F^2$

Then, with a probability at least  $1 - \delta / \text{poly}(nd)$ , we have

$$|C_3| \leq \eta^2 m^2 \|\mathbf{F}(\tau) - Y\|_F^2$$

*Proof.* We have

$$\begin{aligned}
|C_3| &= \|\mathbf{F}(\tau + 1) - \mathbf{F}(\tau)\|_F^2 \\
&= \sum_{i=1}^n \sum_{k=1}^d (\mathbf{F}_{k,i}(\tau + 1) - \mathbf{F}_{k,i}(\tau))^2 \\
&= \sum_{i=1}^n \sum_{k=1}^d m^2 (\langle \beta_k(\tau + 1), \mathbf{S}_i(\tau + 1) \rangle - \langle \beta_k(\tau), \mathbf{S}_i(\tau) \rangle)^2 \\
&= \sum_{i=1}^n \sum_{k=1}^d m^2 \left( \sum_{r=1}^m (\beta_{k,r}(\tau + 1) \cdot \mathbf{S}_{i,r}(\tau + 1) - \beta_{k,r}(\tau) \cdot \mathbf{S}_{i,r}(\tau)) \right)^2 \\
&= \sum_{i=1}^n \sum_{k=1}^d m^2 \left( \sum_{r=1}^m (\beta_{k,r}(\tau + 1) \cdot \mathbf{S}_{i,r}(\tau + 1) - \beta_{k,r}(\tau + 1) \cdot \mathbf{S}_{i,r}(\tau)) \right. \\
&\quad \left. + \beta_{k,r}(\tau + 1) \cdot \mathbf{S}_{i,r}(\tau) - \beta_{k,r}(\tau) \cdot \mathbf{S}_{i,r}(\tau)) \right)^2 \\
&= \sum_{i=1}^n \sum_{k=1}^d m^2 \left( \sum_{r=1}^m (\beta_{k,r}(\tau + 1) \cdot (\mathbf{S}_{i,r}(\tau + 1) - \mathbf{S}_{i,r}(\tau))) \right. \\
&\quad \left. + (\beta_{k,r}(\tau + 1) - \beta_{k,r}(\tau)) \cdot \mathbf{S}_{i,r}(\tau) \right)^2 \\
&= \sum_{i=1}^n \sum_{k=1}^d m^2 (Q_{3,1,i,k} + Q_{3,2,i,k})^2
\end{aligned}$$

where the first step follows from the definition  $C_2$ , the second step follows from the definition of Frobenius norm, the third step follows from Definition G.8, the fourth, fifth and sixth steps follow from simple algebras, the last step follows from defining

$$\begin{aligned}
Q_{3,1,i,k} &= \sum_{r=1}^m \beta_{k,r}(\tau + 1) \cdot (\mathbf{S}_{i,r}(\tau + 1) - \mathbf{S}_{i,r}(\tau)), \\
Q_{3,2,i,k} &= \sum_{r=1}^m (\beta_{k,r}(\tau + 1) - \beta_{k,r}(\tau)) \cdot \mathbf{S}_{i,r}(\tau).
\end{aligned}$$

**Bounding first term.** For the first term, we have

$$\begin{aligned}
|Q_{3,1,i,k}| &= \left| \sum_{r=1}^m \beta_{k,r}(\tau+1) \cdot (S_{i,r}(\tau+1) - S_{i,r}(\tau)) \right| \\
&= \left| \sum_{r=1}^m a_r \cdot w_{r,k}(\tau+1) \cdot (S_{i,r}(\tau+1) - S_{i,r}(\tau)) \right| \\
&\leq |B \cdot \sum_{r=1}^m a_r \cdot (S_{i,r}(\tau+1) - S_{i,r}(\tau))| \\
&\leq |\exp(3B) \cdot \sum_{r=1}^m a_r \cdot \max_{i \in [n]} (\alpha_i(\tau+1)^{-1} - \alpha_i(\tau)^{-1})|
\end{aligned}$$

where the first step follows from the definition of  $Q_{3,1,i,k}$ , the second step follows from Definition F.5, the third step follows from Part 1 of Lemma L.1, last step follows from Part 4 of Lemma L.1, Definition G.7 and  $B \leq \exp(B)$ .

Then by Part 2 of Lemma I.9, applying Hoeffding inequality (Lemma E.4) to the random variables  $a_r \cdot \max_{i \in [n]} (\alpha_i(\tau+1)^{-1} - \alpha_i(\tau)^{-1})$  for  $r \in [m]$  and  $\mathbb{E}[\sum_{r=1}^m a_r \cdot \max_{i \in [n]} (\alpha_i(\tau+1)^{-1} - \alpha_i(\tau)^{-1})] = 0$ , we have

$$|Q_{3,1,i,k}| \leq \left( \eta \frac{\sqrt{nd} \exp(18B)}{m^3} \cdot \|\mathbf{F}(\tau) - Y\|_F + \eta^2 \frac{nd \exp(30B)}{\sqrt{m}} \cdot \|\mathbf{F}(\tau) - Y\|_F \right) \cdot \sqrt{m \log(nd/\delta)}$$

with a probability of at least  $1 - \delta / \text{poly}(nd)$ .

By the Lemma condition, we have

$$\left( \eta \frac{\sqrt{nd} \exp(18B)}{m^3} + \eta^2 \frac{nd \exp(30B)}{\sqrt{m}} \right) \cdot \sqrt{m \log(nd/\delta)} \leq \frac{1}{2\sqrt{nd}} \eta$$

**Bounding second term.** On the other hand, for the second term  $Q_{3,2,k,i}$ , we have

$$\begin{aligned}
|Q_{3,2,k,i}| &= \left| \sum_{r=1}^m (\beta_{k,r}(\tau+1) - \beta_{k,r}(\tau)) \cdot S_{i,r}(\tau) \right| \\
&= \eta \left| \sum_{r=1}^m a_r \Delta w_{r,k}(\tau) \cdot S_{i,r}(\tau) \right| \\
&\leq \eta \frac{\exp(3B)}{m} \left| \sum_{r=1}^m a_r \Delta w_{r,k}(\tau) \right| \\
&\leq \eta \exp(3B) \left| \sum_{r=1}^m a_r \sum_{j=1}^n \sum_{k_1=1}^d (\mathbf{F}_{k_1,j}(\tau) - y_{k_1,j}) \right. \\
&\quad \left. \cdot \left( \langle v_{k_1,r}(\tau), \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \cdot x_{i,k} + a_r \mathbf{S}_{j,r}(\tau) e_{k,k_1} \right) \right| \\
&\leq \eta \frac{\exp(6B)}{m} \left| \sum_{r=1}^m a_r \max_{j \in [n], k, k_1 \in [d]} C_{j,k,k_1,r} \right| \cdot \|\mathbf{F}(\tau) - Y\|_1 \\
&\leq \eta \frac{\sqrt{nd} \exp(6B)}{m} \left| \sum_{r=1}^m a_r \max_{j \in [n], k, k_1 \in [d]} C_{j,k,k_1,r} \right| \cdot \|\mathbf{F}(\tau) - Y\|_F
\end{aligned}$$

where the first step follows from the definition of  $Q_{3,2,k,i}$ , the second step follows from Definition G.13, the third step follows from  $0 \leq S_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1, the fourth step follows from Claim G.12, the fifth step follows from  $0 \leq S_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1,  $\|x_i\|_2 \leq 1$  and defining

$$C_{j,k,k_1,r} := \langle v_{k_1,r}(\tau), \mathbf{S}_j(\tau) \rangle + e_{k,k_1},$$

2754 the last step follows from  $\|U\|_1 \leq \sqrt{nd}\|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ .

2755 Now we follow from Part 6 of Lemma L.2, applying Hoeffding inequality (Lemma E.4) to random  
2756 variables  $a_r \max_{j \in [n], k, k_1 \in [d]} C_{j,k,k_1,r}$  for  $r \in [m]$  and  $\mathbb{E}[\sum_{r=1}^m a_r \max_{j \in [n], k, k_1 \in [d]} C_{j,k,k_1,r}] = 0$ ,  
2757 we have

$$2758 |Q_{3,2,k,i}| \leq \eta \frac{\sqrt{nd} \exp(13B)}{m} \cdot \|F(\tau) - Y\|_F \cdot \sqrt{m \log(nd/\delta)} \leq \frac{1}{2\sqrt{nd}} \eta$$

2762 Finally, we combine all terms, we have

$$2763 |C_3| = \sum_{i=1}^n \sum_{k=1}^d m^2 \left( \frac{1}{2\sqrt{nd}} \eta + \frac{1}{2\sqrt{nd}} \eta \right) \cdot \|F(\tau) - Y\|_F^2$$

$$2764 \leq \eta^2 m^2 \|F(\tau) - Y\|_F^2$$

□

## 2770 I.5 BOUNDING LOSS DURING TRAINING PROCESS

2772 **Lemma I.8.** *If the following conditions hold*

- 2773 • Denote  $F(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.
- 2774 • Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.

2777 Then we have

$$2778 \|F(\tau) - Y\|_F \leq O(\sqrt{nmd})$$

2781 *Proof.* This proof follows from  $\|y_i\| \leq 1$  for  $i \in [n]$  and Definition G.8.

□

## 2783 I.6 HELPFUL LEMMA

2785 **Lemma I.9.** *If the following conditions hold*

- 2786 • Let  $\lambda = \lambda_{\min}(H^*)$ .
- 2787 • Let  $C > 10$  denote a sufficiently large constant.
- 2788 • Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- 2789 • Let  $\delta \in (0, 0.1)$ .
- 2790 • Let  $m \geq \Omega(\lambda^{-2} n^2 d^2 \exp(30B) \sqrt{\log(nd/\delta)})$ .
- 2791 • Let  $r \in [m]$ , let  $i, j \in [n]$ , let  $k, k_1 \in [d]$ .
- 2792 • Let  $\alpha_i(\tau) \in \mathbb{R}$  be defined as Definition F.3.
- 2793 • Let  $\beta_k(\tau) \in \mathbb{R}^m$  be defined as Definition F.5.
- 2794 • Let  $\theta_{k,i}(\tau) \in \mathbb{R}^m$  be defined as Definition G.6.
- 2795 • Let  $u_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.2.
- 2796 • Let  $S_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.7.
- 2797 • Let  $v_k := \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau) \in \mathbb{R}^m$ .
- 2798 • Denote  $F(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.
- 2799 • Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.

2808 Then with a probability at least  $1 - \delta / \text{poly}(nd)$ , we have

- 2810 • *Part 1.*

$$2811 \alpha_i(\tau + 1) - \alpha_i(\tau) \leq \eta \frac{\sqrt{nd} \exp(9B)}{m} \cdot \|F(\tau) - Y\|_F + \eta^2 m^{1.5} \cdot nd \exp(21B) \cdot \|F(\tau) - Y\|_F$$

- 2814 • *Part 2.*

$$2815 \alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1} \leq \eta \frac{\sqrt{nd} \exp(15B)}{m^3} \cdot \|F(\tau) - Y\|_F + \eta^2 \frac{nd \exp(27B)}{\sqrt{m}} \cdot \|F(\tau) - Y\|_F$$

2818 **Proof. Proof of Part 1.**

2819 We have

$$\begin{aligned} 2820 & \alpha_i(\tau + 1) - \alpha_i(\tau) \\ 2821 &= \langle \mathbf{u}_i(\tau + 1), \mathbf{1}_m \rangle - \langle \mathbf{u}_i(\tau), \mathbf{1}_m \rangle \\ 2822 &= \langle \mathbf{u}_i(\tau + 1) - \mathbf{u}_i(\tau), \mathbf{1}_m \rangle \\ 2823 &= \langle \exp(W(\tau + 1)^\top x_i) - \exp(W(\tau)^\top x_i), \mathbf{1}_m \rangle \\ 2824 &= \langle \exp(W(\tau)^\top x_i) \circ (\exp(-\eta \Delta W(\tau)^\top x_i) - \mathbf{1}_m), \mathbf{1}_m \rangle \\ 2825 &= \langle \exp(W(\tau)^\top x_i) \circ (-\eta \Delta W(\tau)^\top x_i + \Theta(1)\eta^2 \cdot (\Delta W(\tau)^\top x_i)^2), \mathbf{1}_m \rangle \\ 2826 &= \langle -\eta \Delta W(\tau)^\top x_i + \Theta(1)\eta^2 \cdot (\Delta W(\tau)^\top x_i)^2, \exp(W(\tau)^\top x_i) \rangle \\ 2827 &\leq \exp(B) \cdot \langle -\eta \Delta W(\tau)^\top x_i + \Theta(1)\eta^2 \cdot (\Delta W(\tau)^\top x_i)^2, \mathbf{1}_m \rangle \\ 2828 &\leq \eta \frac{\sqrt{nd} \exp(9B)}{m} \cdot \|F(\tau) - Y\|_F + \eta^2 m^{1.5} \cdot nd \exp(21B) \cdot \|F(\tau) - Y\|_F \end{aligned}$$

2834 where the first step follows from Definition F.3, the second step follows from simple algebras, the  
2835 third step follows from Definition G.2, the fourth step follows from simple algebra, the fifth step  
2836 follows from Fact E.1, the sixth step follows from simple algebras, the seventh step follows from Part  
2837 4 of Lemma L.1, last step follows from Part 1 and Part 2 of Lemma I.10.

2838 **Proof of Part 2.** We have

$$\begin{aligned} 2839 \alpha_i(\tau + 1)^{-1} - \alpha_i(\tau)^{-1} &= \alpha_i(\tau + 1)^{-1} \alpha_i(\tau)^{-1} \cdot (\alpha_i(\tau + 1) - \alpha_i(\tau)) \\ 2840 &\leq \frac{\exp(6B)}{m^2} \cdot (\alpha_i(\tau + 1) - \alpha_i(\tau)) \\ 2841 &\leq \eta \frac{\sqrt{nd} \exp(15B)}{m^3} \cdot \|F(\tau) - Y\|_F + \eta^2 \frac{nd \exp(27B)}{\sqrt{m}} \cdot \|F(\tau) - Y\|_F \end{aligned}$$

2842 where the first step follows from simple algebras, the second step follows from Part 4 of Lemma L.2,  
2843 the last step follows from Part 1 of this Lemma.  $\square$

2844 **Lemma I.10.** *If the following conditions hold*

- 2850 • Let  $\lambda = \lambda_{\min}(H^*)$ .
- 2851 • Let  $W(\tau) \in \mathbb{R}^{m \times d}$  be defined as Definition G.13, let  $a \in \mathbb{R}^m$  be defined as Definition F.1.
- 2852 • Let  $C > 10$  denote a sufficiently large constant.
- 2853 • Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- 2854 • Let  $\delta \in (0, 0.1)$ .
- 2855 • Let  $m \geq \Omega(\lambda^{-2} n^2 d^2 \exp(30B) \sqrt{\log(nd/\delta)})$ .
- 2856 • Let  $r \in [m]$ , let  $i, j \in [n]$ , let  $k, k_2 \in [d]$ .
- 2857 • Let  $S_i(\tau) \in \mathbb{R}^m$  be defined as Definition G.7.

- Let  $v_{k,r} := \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau) \in \mathbb{R}^m$ .
- Denote  $F(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.
- Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.
- Let  $\eta = \lambda / (m \cdot \text{poly}(n, d, \exp(B)))$  denote the learning rate.

Then with a probability at least  $1 - \delta / \text{poly}(nd)$ , we have

- **Part 1.**

$$|\langle \eta \Delta W(\tau)^\top x_i, \mathbf{1}_m \rangle| \leq \eta \frac{\sqrt{nd} \exp(8B)}{m} \cdot \|F(\tau) - Y\|_F$$

- **Part 2.**

$$|\langle \eta^2 (\Delta W(\tau)^\top x_i)^2, \mathbf{1}_m \rangle| \leq \eta^2 m^{1.5} \cdot nd \exp(20B) \cdot \|F(\tau) - Y\|_F$$

**Proof. Proof of Part 1.** We have

$$\begin{aligned} & |\langle \eta \Delta W(\tau)^\top x_i, \mathbf{1}_m \rangle| \\ &= \eta \left| \sum_{r=1}^m \langle \Delta w_r(\tau), x_i \rangle \right| \\ &\leq \eta \left| \sum_{r=1}^m m \sum_{j=1}^n \sum_{k=1}^d (F_{k,i}(\tau) - y_{k,i}) \cdot \left( \langle v_{k,r}(\tau), S_j(\tau) \rangle \cdot S_{j,r}(\tau) \cdot x_j^\top + a_r S_{j,r}(\tau) e_k^\top \right) x_i \right| \\ &\leq \eta \left| \sum_{r=1}^m m \sum_{j=1}^n \sum_{k=1}^d (F_{k,i}(\tau) - y_{k,i}) \cdot \left( \langle \beta_{k,r}(\tau) \cdot \mathbf{1}_m - \beta_k(\tau), S_j(\tau) \rangle \cdot S_{j,r}(\tau) \cdot x_j^\top + a_r S_{j,r}(\tau) e_k^\top \right) x_i \right| \\ &\leq \eta \left| \sum_{r=1}^m m \sum_{j=1}^n \sum_{k=1}^d (F_{k,i}(\tau) - y_{k,i}) \cdot \left( a_r w_{r,k} + \langle -a \circ W_{k,*}(\tau), S_j(\tau) \rangle \cdot S_{j,r}(\tau) \cdot x_j^\top + a_r S_{j,r}(\tau) e_k^\top \right) x_i \right| \\ &\leq \eta \frac{\exp(3B)}{m} \sum_{r=1}^m \sigma_r \max_{j \in [n], k \in [d]} C_{j,k,r} \|F(\tau) - Y\|_1 \\ &\leq \eta \frac{\sqrt{nd} \exp(3B)}{m} \sum_{r=1}^m \sigma_r \max_{j \in [n], k \in [d]} C_{j,k,r} \|F(\tau) - Y\|_F \end{aligned}$$

where the first step follows from simple algebras, the second step follows from Claim G.12, the third step follows from the definition of  $v_{k,r}$ , the fourth step follows from Definition F.5 and simple algebras, the fifth step follows from  $\|x_i\|_2 \leq 1$ ,  $1 \leq S_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1, definition of  $\ell_1$  norm and defining

$$C_{j,k,r} := |w_{r,k}| + |\langle -W_{k,*}(\tau), S_j(\tau) \rangle| + \|e_k\|, \sigma_r \in \{+1, -1\},$$

the last step follows from  $\|U\|_1 \leq \sqrt{nd} \|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ .

Thus, by following Part 1 and Part 11 of Lemma L.2 and Hoeffding inequality (Lemma E.4), we have

$$|\langle \eta \Delta W(\tau)^\top x_i, \mathbf{1}_m \rangle| \leq \eta \frac{\sqrt{nd} \exp(8B)}{m} \cdot \|F(\tau) - Y\|_F$$

with a probability at least  $1 - \delta / \text{poly}(nd)$ .

**Proof of Part 2.** We have

$$\begin{aligned} & |\langle \eta^2 (\Delta W(\tau)^\top x_i)^2, \mathbf{1}_m \rangle| \\ &\leq \eta^2 \sum_{r=1}^m (\langle \Delta w_r(\tau), x_i \rangle)^2 \end{aligned}$$

$$\begin{aligned}
&\leq \eta^2 \sum_{r=1}^m \left( m \sum_{j=1}^n \sum_{k=1}^d (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( \langle v_{k,r}(\tau), \mathbf{S}_j(\tau) \rangle \cdot \mathbf{S}_{j,r}(\tau) \cdot x_j^\top + a_r \mathbf{S}_{j,r}(\tau) e_k^\top \right) x_i \right)^2 \\
&\leq \eta^2 \exp(6B) \sum_{r=1}^m \left( \sum_{j=1}^n \sum_{k=1}^d (\mathbf{F}_{k,i}(\tau) - y_{k,i}) \cdot \left( \langle v_{k,r}(\tau), \mathbf{S}_j(\tau) \rangle \cdot x_j^\top + a_r e_k^\top \right) x_i \right)^2 \\
&\leq \eta^2 m \exp(20B) \cdot \|\mathbf{F}(\tau) - Y\|_1^2 \\
&\leq \eta^2 m \sqrt{nm} \exp(20B) \cdot \|\mathbf{F}(\tau) - Y\|_1 \\
&\leq \eta^2 m^{1.5} \cdot nd \exp(20B) \cdot \|\mathbf{F}(\tau) - Y\|_F
\end{aligned}$$

where the first step follows from simple algebras, the second step follows from Claim G.12, the third step follows from  $0 \leq S_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of Lemma L.1, the fourth step follows from  $\langle v_{k,r}(\tau), \mathbf{S}_j(\tau) \rangle \leq \exp(6B)$  by Part 6 of Lemma L.2,  $\|x_i\|_2 \leq 1$ ,  $\exp(6B) + 1 \leq \exp(7B)$  and the definition of  $\ell_1$  norm, the fifth step follows from Lemma I.8, the last step follows from  $\|U\|_1 \leq \|U\|_F$  for  $U \in \mathbb{R}^{n \times d}$ .  $\square$

## J CONVERGENCE OF PREFIX LEARNING

Here, we provide all the properties we need for math induction for NTK happening.

**Definition J.1** (Properties). *We state the following properties*

- *General Condition 1.* Let  $\lambda = \lambda_{\min}(H^*) > 0$
- *General Condition 2.* Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- *General Condition 3.* Let  $\eta$  be defined as
$$\eta := \lambda / (m \text{poly}(n, d, \exp(B))).$$
- *General Condition 4.* Let  $D := 2\lambda^{-1} \cdot \exp(20B) \frac{\sqrt{nd}}{m} \|Y - \mathbf{F}(0)\|_F$
- *General Condition 5.* Let  $w_r$  and  $a_r$  be defined as Definition F.1.
- *General Condition 6.*  $D < R = \lambda / \text{poly}(n, d, \exp(B))$
- *General Condition 7.*  $m = \lambda^{-2} \text{poly}(n, d, \exp(B))$
- **Weight Condition.**  $\|w_r(t) - w_r(0)\|_2 \leq D < R, \forall r \in [m]$
- **Loss Condition.**  $\|\text{vec}(\mathbf{F}(i) - Y)\|_2^2 \leq \|\text{vec}(\mathbf{F}(0) - Y)\|_2^2 \cdot (1 - m\eta\lambda/2)^i, \forall i \in [t]$
- **Gradient Condition.**  $\eta \|\Delta w_r(i)\|_2 \leq 0.01 \forall r \in [m], \forall i \in [t]$

### J.1 MAIN RESULT

Our main result is presented as follows.

**Theorem J.2** (Main result, formal version of Theorem 3.2). *For any  $\epsilon, \delta \in (0, 0.1)$ , if the following conditions hold*

- *Let  $\lambda = \lambda_{\min}(H^*) > 0$*
- *Let  $B = \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$*
- *Let  $m = \lambda^{-2} \text{poly}(n, d, \exp(B))$*
- *Let  $\eta = \lambda / (m \text{poly}(n, d, \exp(B)))$*
- *Let  $\hat{T} = \Omega((m\eta\lambda)^{-1} \log(nd/\epsilon))$*

2970 Then, after  $\widehat{T}$  iterations, with probability at least  $1 - \delta$ , we have

$$2971 \quad \|F(\widehat{T}) - Y\|_F^2 \leq \epsilon.$$

2972  
2973  
2974 *Proof.* We have  $\|F(0) - Y\|_F^2 \leq nd$  as Lemma J.6. Using the choice of  $\widehat{T}$ , it follows directly from  
2975 the alternative application of Lemma J.3 and Lemma J.4.  $\square$   
2976

## 2977 J.2 INDUCTION PART 1. FOR WEIGHTS

2978 In this section, we introduce the induction lemma for weights.

2979 **Lemma J.3** (Induction Part 1 for weights). *If the following conditions hold*

- 2980 • *Suppose properties in Definition J.1 are true*

2981 For  $t + 1$  and  $\forall r \in [m]$ , it holds that:

$$2982 \quad \|w_r(t + 1) - w_r(0)\|_2 \leq D.$$

2983 *Proof.* We have

$$2984 \quad \eta \sum_{i=0}^{\infty} (1 - m\eta\lambda/2)^i \leq \eta \frac{4}{m\lambda} \quad (18)$$

2985 where this step follows from Fact E.2.

$$2986 \quad \begin{aligned} 2987 \quad \|w_r(t + 1) - w_r(0)\|_2 &\leq \eta \sum_{\tau=0}^t \|\Delta w_r(\tau)\|_2 \\ 2988 &\leq \eta \sum_{\tau=0}^t \sqrt{nd} \exp(11B) \cdot \|F(\tau) - Y\|_F \\ 2989 &\leq \eta \sqrt{nd} \exp(11B) \cdot \sum_{\tau=0}^t (1 - m\eta\lambda/2)^i \cdot \|F(0) - Y\|_F \\ 2990 &\leq 2\eta \frac{1}{m\lambda} \sqrt{nd} \exp(11B) \cdot \|F(0) - Y\|_F \\ 2991 &\leq D \end{aligned}$$

2992 where the third step follows from the triangle inequality, the second step follows from Eq. (22), the  
2993 third step follows from Lemma J.4, the fourth step follows from Eq. (18), the last step follows from  
2994 *General Condition 4.* in Definition J.1.  $\square$   
2995

## 3014 J.3 INDUCTION PART 2. FOR LOSS

3015 Now, we present our next induction lemma.

3016 **Lemma J.4** (Induction Part 2 for loss). *Let  $t$  be a fixed integer.*

3017 *If the following conditions hold*

- 3018 • *Suppose properties in Definition J.1 are true*

3019 Then we have

$$3020 \quad \|F(t + 1) - y\|_F^2 \leq (1 - m\eta\lambda/2)^{t+1} \cdot \|F(0) - y\|_F^2.$$



3024 *Proof.* We have

$$\begin{aligned}
3025 & \quad \|F(t+1) - y\|_F^2 \\
3026 & \leq \|F(t) - y\|_F^2 + C_0 + C_1 + C_2 + C_3 \\
3027 & = \|F(t) - y\|_F^2 + C_0 + C_{1,1} + C_{1,2} + C_2 + C_3 \\
3028 & \leq \|F(t) - y\|_F^2 \cdot (1 + 0.1\eta m\lambda - 1.6\eta m\lambda + 0.1\eta m\lambda + \eta^2 m \cdot n^2 d^2 \exp(16B) + \eta^2 m^2) \\
3029 & \leq \|F(t) - y\|_F^2 \cdot (1 - 1.4\eta m\lambda + \eta^2 m \cdot n^2 d^2 \exp(16B) + \eta^2 m^2) \tag{19}
\end{aligned}$$

3030 where the first step follows from Lemma I.1, the second step follows from the definitions of  $C_1$ ,  $C_{1,1}$   
3031 and  $C_{1,2}$ , the third step follows from Lemma I.2 and Lemma I.3.

3032 **Choice of parameter.** Here, we explain the condition setting in Definition J.1:

- 3033 • To get our results in Lemma I.2 and Lemma I.3, we have to let  $m \geq \Omega(\lambda^{-2} n^2 d^2 \cdot \exp(30B) \cdot$   
3034  $\sqrt{\log(nd/\delta)})$ .
- 3035 • If we let  $\eta \leq O(\lambda/(mn^2 d^2 \exp(16B)))$ , we can have

$$3036 \quad \eta^2 m \cdot n^2 d^2 \exp(16B) + \eta^2 m^2 \leq 0.9\eta m\lambda. \tag{20}$$

3037 Thus, combining Eq. (19) and Eq. (20), we have

$$3038 \quad \|F(t+1) - y\|_F^2 \leq (1 - m\eta\lambda/2) \cdot \|F(t) - y\|_F^2 \tag{21}$$

3039 Then by Eq. (21), we conclude all  $\|F(\tau) - y\|_F^2$  for  $\tau \in [t]$ , we have

$$3040 \quad \|F(t+1) - y\|_F^2 \leq (1 - m\eta\lambda/2)^{t+1} \cdot \|F(0) - y\|_F^2$$

3041 □

#### 3042 J.4 INDUCTION PART 3. FOR GRADIENT

3043 In this section, we present the induction lemma for gradients.

3044 **Lemma J.5** (Induction Part 3 for gradient). *Let  $t$  be a fixed integer.*

3045 *If the following conditions hold*

- 3046 • *Suppose properties in Definition J.1 are true*

3047 *Then we have*

$$3048 \quad \eta \|\Delta w_r(t)\|_2 \leq 0.01, \forall r \in [m]$$

3049 *Proof.* Firstly, we have

$$\begin{aligned}
3050 & \quad \|\Delta w_r(t)\|_2 \leq \|\Delta w_r(t)\|_1 \\
3051 & \leq \sum_{k_1=1}^d \left| m \sum_{i=1}^n \sum_{k=1}^d (F_{k,i}(t) - y_{k,i}) \cdot \left( \langle v_{k,r}(t), S_i(t) \rangle \cdot S_{i,r}(t) \cdot x_{i,k_1} + a_r S_{i,r}(t) e_{k,k_1} \right) \right| \\
3052 & \leq \sqrt{nd} \exp(11B) \|F(t) - Y\|_F \tag{22}
\end{aligned}$$

3053 where the first step follows from  $\|U\|_F \leq \|U\|_1$  for  $U \in \mathbb{R}^{n \times d}$ , the second step follows from  
3054 Claim G.12, the last step follows from the definition of 4  $\ell_1$  norm,  $0 \leq S_{i,r} \leq \frac{\exp(3B)}{m}$  by Part 11 of  
3055 Lemma L.1,  $\|x_i\|_2 \leq 1$  and Part 6 of Lemma L.2.

3056 Then by the property of  $\eta$  in Definition J.1, we have

$$3057 \quad \eta \|\Delta w_r(t)\|_2 \leq 0.01, \forall r \in [m]$$

3058 □

3078 J.5 BOUNDING LOSS AT INITIALIZATION  
3079

3080 **Lemma J.6.** *If the following conditions hold*

- 3081
- 3082 • Denote  $F(\tau) \in \mathbb{R}^{n \times d}$  as Definition G.8.
  - 3083 • Let  $Y \in \mathbb{R}^{n \times d}$  denote the labels.
- 3084

3085 *Then we have*

$$3086 \quad \|F(0) - Y\|_F \leq O(\sqrt{nd})$$

3087

3088 *Proof.* This proof follows from  $\|y_i\| \leq 1$  for  $i \in [n]$  and Definition G.8. □

3089

3090 K NTK-ATTENTION

3091 In this section, we compute the error bound of our NTK-Attention in approximating prefix matrix  $P \in \mathbb{R}^{m \times d}$ . In Appendix K.1, we provide the formal definition of our NTK-Attention. In Appendix K.2, we give our main theorem of error bound. In Appendix K.3, we state tools from (Alman & Song, 2023).

3092

3093 K.1 DEFINITIONS

3094 **Definition K.1.** *If the following conditions hold:*

- 3095
- 3101 • Given input  $X \in \mathbb{R}^{L \times d}$ , prefix matrix  $P \in \mathbb{R}^{m \times d}$ .
  - 3102
  - 3103 • Let  $S := \begin{bmatrix} P \\ X \end{bmatrix} \in \mathbb{R}^{(m+L) \times d}$ .
  - 3104
  - 3105 • Given projections  $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$
  - 3106
  - 3107 • Let  $Q := XW_Q \in \mathbb{R}^{L \times d}$ .
  - 3108
  - 3109 • Let  $K_P := SW_Q \in \mathbb{R}^{(m+L) \times d}$
  - 3110
  - 3111 • Let  $V_P := SW_V \in \mathbb{R}^{(m+L) \times d}$
  - 3112
  - 3113 • Let  $A := \exp(QK_P^\top) \in \mathbb{R}^{L \times (m+L)}$ .
  - 3114
  - 3115 • Let  $D := \text{diag}(A\mathbf{1}_{(m+L)}) \in \mathbb{R}^{L \times L}$ .

3116 *We define:*

$$3117 \quad \text{Attn}(Q, K, V) := D^{-1}AV_P.$$

3118

3119 K.2 ERROR BOUND

3120 Here, we provide our two statements about error bound.

3121 **Theorem K.2** (Formal version of Theorem 4.1). *Given an input matrix  $X \in \mathbb{R}^{L \times d}$  and prefix matrix  $P \in \mathbb{R}^{m \times d}$ , we denote  $Q = XW_Q$ ,  $K_C = PW_K$  and  $V_C = PW_V$ . If the condition Eq. (7),  $\|Q\|_\infty \leq o(\sqrt{\log m})$ ,  $\|K_C\|_\infty \leq o(\sqrt{\log m})$ ,  $\|V_C\|_\infty \leq o(\sqrt{\log m})$  and  $d = O(\log m)$  holds, then Algorithm 2 outputs a matrix  $T \in \mathbb{R}^{L \times d}$  within time complexity of  $O(L^2d)$  that satisfies:*

3122

$$3123 \quad \|T - \text{PrefixAttn}(X, P)\|_\infty \leq 1/\text{poly}(m).$$

3124

3125 *Proof.* Following Definition K.1, we can have matrix  $A \in \mathbb{R}^{L \times (m+L)}$  as follows:

$$3126 \quad A = QK^\top$$

$$3127 \quad = [\exp(XW_QW_K^\top X^\top) \quad \exp(XW_QW_K^\top P^\top)]$$

3128

3129

3130

3131

where the second step follows from  $K = SW_K$  and  $S = \begin{bmatrix} P \\ X \end{bmatrix}$ .

Our Algorithm 2 actually implement on using  $Q = XW_Q$  and  $PW_K$  to approximate  $\exp(XW_QW_K^\top P^\top)$  by Lemma K.7.

Trivially, this proof follows from Theorem K.5 and Lemma K.7.  $\square$

**Corollary K.3.** *Given an input matrix  $X \in \mathbb{R}^{L \times d}$  and prefix matrix  $P \in \mathbb{R}^{m \times d}$ , we denote  $Q = XW_Q$ ,  $K_C = PW_K$  and  $V_C = PW_V$ . If the condition Eq. (7),  $\|Q\|_\infty \leq o(\sqrt{\log m})$ ,  $\|K_C\|_\infty \leq o(\sqrt{\log m})$ ,  $\|V_C\|_\infty \leq o(\sqrt{\log m})$  and  $d = O(\log m)$  holds, then there exists an algorithm that outputs a matrix  $T \in \mathbb{R}^{L \times d}$  within time complexity of  $O(L^{1+o(1)}d)$  that satisfies:*

$$\|T - \text{PrefixAttn}(X, P)\|_\infty \leq 1/\text{poly}(m).$$

*Proof.* The algorithm and proof can trivially follow from Algorithm 1, 2, 3 and Theorem 1 in HyperAttention (Han et al., 2024).  $\square$

### K.3 TOOLS FROM FAST ATTENTION

In this section, we introduce some tools from previous work which we have used.

**Definition K.4** (Approximate Attention Computation AAttC( $n, d, B, \epsilon_a$ ), Definition 1.2 in (Alman & Song, 2023)). *Let  $\epsilon_a > 0$  and  $B > 0$  be parameters. Given three matrices  $Q, K, V \in \mathbb{R}^{n \times d}$ , with the guarantees that  $\|Q\|_\infty \leq B$ ,  $\|K\|_\infty \leq B$ , and  $\|V\|_\infty \leq B$ , output a matrix  $T \in \mathbb{R}^{n \times d}$  which is approximately equal to  $D^{-1}AV$ , meaning,*

$$\|T - D^{-1}AV\|_\infty \leq \epsilon_a.$$

Here, for a matrix  $M \in \mathbb{R}^{n \times n}$ , we write  $\|M\|_\infty := \max_{i,j} |M_{i,j}|$ .

**Theorem K.5** (Upper bound, Theorem 1.4 in (Alman & Song, 2023)). *There is an algorithm that solves AAttC( $n, d = O(\log n), B = o(\sqrt{\log n}), \epsilon_a = 1/\text{poly}(n)$ ) in time  $n^{1+o(1)}$ .*

**Definition K.6** (Definition 3.1 in (Alman & Song, 2023)). *Let  $r \geq 1$  denote a positive integer. Let  $\epsilon \in (0, 0.1)$  denote an accuracy parameter. Given a matrix  $A \in \mathbb{R}_{\geq 0}^{n \times n}$ , we say  $\tilde{A} \in \mathbb{R}_{\geq 0}^{n \times n}$  is an  $(\epsilon, r)$ -approximation of  $A$  if*

- $\tilde{A} = U_1 \cdot U_2^\top$  for some matrices  $U_1, U_2 \in \mathbb{R}^{n \times r}$  (i.e.,  $\tilde{A}$  has rank at most  $r$ ), and
- $|\tilde{A}_{i,j} - A_{i,j}| \leq \epsilon \cdot A_{i,j}$  for all  $(i, j) \in [n]^2$ .

**Lemma K.7** (Lemma 3.4 in (Alman & Song, 2023)). *Suppose  $Q, K \in \mathbb{R}^{n \times d}$ , with  $\|Q\|_\infty \leq B$ , and  $\|K\|_\infty \leq B$ . Let  $A := \exp(QK^\top/d) \in \mathbb{R}^{n \times n}$ . For accuracy parameter  $\epsilon \in (0, 1)$ , there is a positive integer  $g$  bounded above by*

$$g = O\left(\max\left\{\frac{\log(1/\epsilon)}{\log(\log(1/\epsilon)/B^2)}, B^2\right\}\right),$$

and a positive integer  $r$  bounded above by

$$r \leq \binom{2(g+d)}{2g}$$

such that: There is a matrix  $\tilde{A} \in \mathbb{R}^{n \times n}$  that is an  $(\epsilon, r)$ -approximation (Definition K.6) of  $A \in \mathbb{R}^{n \times n}$ . Furthermore, we can construct the matrices  $U_1 := \phi(Q)$  and  $U_2 := \phi(K)$  through a function  $\phi(\cdot)$  defining  $\tilde{A} = U_1U_2^\top$  can be computed in  $O(n \cdot r)$  time.

### L TAYLOR SERIES

In this section, we provide some perturbation analysis for NTK analysis.

**Lemma L.1** (Lemma B.1 in (Li et al., 2024a)). *If the following conditions hold*

- 3186 • Let  $C > 10$  denote a sufficiently large constant
- 3187
- 3188 • Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- 3189
- 3190 • Let  $W = [w_1, \dots, w_m]$  and  $w_r$  be random Gaussian vectors from  $\mathcal{N}(0, \sigma^2 I_d)$ .
- 3191 • Let  $V = [v_1, \dots, v_m]$  and  $v_r$  denote the vector where  $\|v_r - w_r\|_2 \leq R, \forall r \in [m]$ .
- 3192
- 3193 • Let  $x_i \in \mathbb{R}^d$  and  $\|x_i\|_2 \leq 1, \forall i \in [n]$ .
- 3194
- 3195 • Let  $R \in (0, 0.01)$ .
- 3196 • Let  $S_i$  and  $\tilde{S}_i$  be the softmax function corresponding to  $W$  and  $V$  respectively.
- 3197
- 3198 • Let  $\alpha_i = \langle \mathbf{1}_m, \exp(W^\top x_i) \rangle$  and  $\tilde{\alpha}_i = \langle \mathbf{1}_m, \exp(V^\top x_i) \rangle, \forall i \in [n]$ .

3199 Then, with probability at least  $1 - \delta / \text{poly}(nd)$ , we have

- 3200
- 3201 • Standard inner product
- 3202
- 3203 – Part 1.  $|\langle w_r, x_i \rangle| \leq B, \forall i \in [n], \forall r \in [m]$
- 3204 – Part 2.  $|\langle v_r, x_i \rangle| \leq B + R, \forall i \in [n], \forall r \in [m]$
- 3205 – Part 3.  $|\langle w_r - v_r, x_i + x_j \rangle| \leq 2R, \forall i, j \in [n], \forall r \in [m]$
- 3206
- 3207 • exp function
- 3208 – Part 4.  $\exp(-B) \leq \exp(\langle w_r, x_i \rangle) \leq \exp(B), \forall i \in [n], \forall r \in [m]$
- 3209 – Part 5.  $\exp(-B - R) \leq \exp(\langle v_r, x_i \rangle) \leq \exp(B + R), \forall i \in [n], \forall r \in [m]$
- 3210 – Part 6.  $|\exp(\langle w_r - v_r, x_i + x_j \rangle) - 1| \leq 4R, \forall i, j \in [n], \forall r \in [m]$
- 3211 – Part 7.  $|\exp(\langle w_r, x_i \rangle) - \exp(\langle v_r, x_i \rangle)| \leq R \exp(B + R), \forall i \in [n], \forall r \in [m]$
- 3212
- 3213 • softmax S function
- 3214 – Part 8.  $|\alpha_i - \tilde{\alpha}_i| \leq mR \exp(B + R), \forall i \in [n]$
- 3215 – Part 9.  $|\alpha_i^{-1} - \tilde{\alpha}_i^{-1}| \leq \frac{R}{m} \exp(3B + 2R), \forall i \in [n]$
- 3216 – Part 10.  $|S_{i,r}| \leq \exp(2B)/m, \forall i \in [n], \forall r \in [m]$
- 3217 – Part 11.  $|\tilde{S}_{i,r}| \leq \exp(2B + 2R)/m, \forall i \in [n], \forall r \in [m]$
- 3218 – Part 12.  $|S_{i,r} - \tilde{S}_{i,r}| \leq \frac{R}{m} \exp(4B + 3R), \forall i \in [n], \forall r \in [m]$
- 3219 – Part 13. for any  $z \in \mathbb{R}^m$  and  $\|z\|_\infty \leq 1$ , we have  $|\langle z, S_i \rangle - \langle z, \tilde{S}_i \rangle| \leq R \exp(4B + 3R), \forall i \in [n]$

3222 **Lemma L.2.** If the following conditions hold

- 3223
- 3224 • Let  $C > 10$  denote a sufficiently large constant
- 3225
- 3226 • Let  $B := \max\{C\sigma\sqrt{\log(nd/\delta)}, 1\}$ .
- 3227
- 3228 • Let  $W = [w_1, \dots, w_m]$  and  $w_r$  be random Gaussian vectors from  $\mathcal{N}(0, \sigma^2 I_d)$ .
- 3229 •  $w_r$  for  $r \in [m]$  satisfies  $\|w_r\|_2 \leq B$  with probability at least  $1 - \delta / \text{poly}(nd)$  as in
- 3230 Lemma L.1.
- 3231 • Let  $a \in \mathbb{R}^m$  be defined as Definition F.1.
- 3232
- 3233 • Define  $\beta_k := W_{k,*} \circ a \in \mathbb{R}^m$  for  $k \in [d]$  as Definition F.5.
- 3234
- 3235 • Define  $v_{k,r} := \beta_{k,r} \cdot \mathbf{1}_m - \beta_k \in \mathbb{R}^m$  for  $k \in [d]$  and  $r \in [m]$  as Definition H.1.
- 3236
- 3237 • Define  $\alpha_i$  for  $i \in [n]$  as Definition F.3.

3238 Then, with probability at least  $1 - \delta / \text{poly}(nd)$ , we have

- 3239 • Part 1.  $|\beta_{k,r}| \leq B$

- 3240 • *Part 2.*  $\|\beta_k\|_2 \leq B\sqrt{m}$   
 3241  
 3242 • *Part 3.*  $\|v_{k,r}\|_2 \leq 2\sqrt{m}B$   
 3243  
 3244 • *Part 4.*  $|\alpha^{-1}| \leq \exp(B)/m$   
 3245  
 3246 • *Part 5.*  $\langle \beta_k, S_i \rangle \leq \exp(4B)$   
 3247  
 3248 • *Part 6.*  $\langle v_{k,r}, S_i \rangle \leq \exp(6B)$

3249 *Proof. Proof of Part 1.* We can get the proof by Gaussian tail bound.

3250 **Proof of Part 2.** We have

$$\begin{aligned} 3251 \|\beta_k\|_2 &= \sqrt{\sum_{r=1}^m \beta_{k,r}^2} \\ 3252 &\leq \sqrt{\sum_{r=1}^m B^2} \\ 3253 &\leq \sqrt{m} \cdot B \end{aligned}$$

3254 where the first step follows from the definition of  $\ell_2$  norm, the second step follows from Part 1 of this  
 3255 Lemma, the last step follows from simple algebras.

3256 **Proof of Part 3.** We have

$$\begin{aligned} 3257 \|v_{k,r}\|_2 &= \sqrt{\sum_{r_1=1}^m (\beta_{k,r} - \beta_{k,r_1})^2} \\ 3258 &\leq \sqrt{\sum_{r_1=1}^m \beta_{k,r}^2 + \beta_{k,r_1}^2 + |2\beta_{k,r}\beta_{k,r_1}|} \\ 3259 &\leq \sqrt{\sum_{r_1=1}^m 4B^2} \\ 3260 &\leq 2\sqrt{m} \cdot B \end{aligned}$$

3261 where the first step follows from the definition of  $\ell_2$  norm, the second step follows from simple  
 3262 algebras, the third step follows from Part 1 of this Lemma, the last step follows from simple algebras.

3263 **Proof of Part 4.** This proof follows from Part 4 of Lemma L.1 and Definition F.3.

3264 **Proof of Part 5.** We have

$$\begin{aligned} 3265 \langle \beta_k, S_i \rangle &\leq \|\beta_k\|_2 \cdot \|S_i\|_2 \\ 3266 &\leq \sqrt{m}B \cdot \|S_i\|_2 \\ 3267 &\leq \sqrt{m}B \cdot \sqrt{\sum_{r=1}^m S_{i,r}^2} \\ 3268 &\leq \sqrt{m}B \cdot \sqrt{\sum_{r=1}^m \frac{\exp(6B)}{m^2}} \\ 3269 &\leq \sqrt{m}B \cdot \sqrt{\frac{\exp(6B)}{m}} \\ 3270 &\leq B \exp(3B) \\ 3271 &\leq \exp(4B) \end{aligned}$$

3294 where the first step follows from Cauchy-Schwarz inequality, the second step follows from Part 2  
3295 of this Lemma, the third step follows from the definition of  $\ell_2$  norm, the fourth step follows from  
3296 Part 11 of Lemma L.1, the fifth step follows from triangle inequality, the sixth step follows from  
3297  $B \leq \exp(B)$ , last step follows from simple algebras.

3298 **Proof of Part 6.** This proof follows from Part 3 of this Lemma,  $B \leq \exp(B)$  and Part 11 of  
3299 Lemma L.1. □

3300  
3301  
3302  
3303  
3304  
3305  
3306  
3307  
3308  
3309  
3310  
3311  
3312  
3313  
3314  
3315  
3316  
3317  
3318  
3319  
3320  
3321  
3322  
3323  
3324  
3325  
3326  
3327  
3328  
3329  
3330  
3331  
3332  
3333  
3334  
3335  
3336  
3337  
3338  
3339  
3340  
3341  
3342  
3343  
3344  
3345  
3346  
3347