
Transformer Efficiently Learns Low-dimensional Target Functions In-context

Anonymous Authors¹

Abstract

Transformers can efficiently learn in-context from example demonstrations. We study ICL of a nonlinear function class via transformer with a nonlinear MLP layer: given a class of *single-index* target functions $f_*(\mathbf{x}) = \sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle)$, where the index features $\boldsymbol{\beta} \in \mathbb{R}^d$ are drawn from a rank- $r \ll d$ subspace, we show that a nonlinear transformer optimized by gradient descent learns f_* in-context with a prompt length that only depends on the dimension of function class r . In contrast, an algorithm that directly learns f_* on the test prompt yields a statistical complexity that scales with the ambient dimension d . Our result highlights the adaptivity of ICL to low-dimensional structures of the function class.

1. Introduction

Transformers (Vaswani et al., 2017) possess the remarkable ability of *in-context learning (ICL)* (Brown et al., 2020), whereby the model constructs a predictor from a prompt consisting of pairs of labeled examples without updating any parameters. A common explanation is that the trained transformer can implement a learning algorithm, such as gradient descent on the in-context examples, in its forward pass (Dai et al., 2022; Von Oswald et al., 2023).

Many recent theoretical works focus on learning *linear functions* using linear transformers, and it can be shown that minima of the pretraining loss implements one (pre-conditioned) gradient descent step on the least squares objective computed on the test prompt (Zhang et al., 2023; Ahn et al., 2023; Mahankali et al., 2023).

The motivation of our work is the observation that the simple setting of learning linear models with linear transformers does not fully capture the statistical efficiency and adaptivity of ICL. Specifically,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

- A linear transformer has limited expressivity: specifically, the pretrained linear transformer cannot outperform directly solving linear regression on the test prompt. Thus, we ask the following question: *With the aid of MLP layer, can a pretrained transformer learn a nonlinear function class in-context, and outperform simple baselines such as one gradient step on the test prompt?*
- A key feature of ICL is the *adaptivity* to structure of the function class; for example, prior empirical results have shown that transformers may match the performance of ridge regression or LASSO, depending on sparsity of the pretrained task distribution (Garg et al., 2022). Such adaptivity cannot be fully explained by the one gradient step algorithm on the test prompt, which does not take into account the “prior” distribution of target functions. Hence a natural question to ask is that, *“Can a pretrained transformer adapt to certain structures of the target function class, and how does such adaptivity contribute to the statistical efficiency of ICL?”*

1.1. Our Contributions

Gaussian single-index models. To address the above questions, we study the in-context learning of the *single-index* function class, where the t -th pretraining task is constructed as $\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$, $y_i = \sigma_*^t(\langle \mathbf{x}_i, \boldsymbol{\beta}_t \rangle) + \varsigma_i$, where $\sigma_*^t : \mathbb{R} \rightarrow \mathbb{R}$ is the link function, and $\boldsymbol{\beta}_t \in \mathbb{R}^d$ is the index feature vector which is randomly drawn from some fixed *rank- r subspace* for some $r \leq d$. Thus, outputs only depend on the direction of $\boldsymbol{\beta}_t$ in the d -dimensional input space (See Section 2 for details). Due to the nonlinear link function, single-index targets cannot be learned by linear transformers.

For this function class, the statistical efficiency of simple algorithms that can be implemented on the in-context examples has been extensively studied: given a link function with degree P and information exponent k (defined as the index of the smallest non-zero coefficient in the Hermite expansion of σ_*), we know that kernel methods can learn the function with $n \gtrsim d^P$ samples (Ghorbani et al., 2021; Donhauser et al., 2021), whereas two-layer neural network trained by gradient descent can achieve a better sample complexity $n \gtrsim d^{\Theta(k)}$ (Ben Arous et al., 2021; Bietti et al., 2022). These serve as a baseline for comparing the statisti-

cal efficiency of ICL.

Moreover, in our problem setting, there is a low-dimensional structure: the subspace from which β_t is drawn is low-dimensional. In particular, when $r \ll d$, we expect the transformer to efficiently extract the rank- r subspace during pretraining, and hence outperform baseline algorithms that directly learn the target function from the test prompt, which cannot make use of such information.

Transformer learns single-index models in-context.

We characterize the sample complexity of learning the single-index model in-context, using a transformer with a nonlinear MLP block and linear self-attention module, optimized by gradient descent.

Informally speaking, our main theorem states that the length of the test prompt N^* required to achieve sufficiently small generalization error (See (2.1) for the definition) is $r^{\Theta(P)}$, where P is the highest degree of the link function. Most importantly, it does not depend on the ambient dimension d (up to polylogarithmic term), but only the dimension r of the feature subspace. When $r \ll d$, we see a separation between ICL and algorithms that directly learn the single-index function from the test prompt such as gradient descent (where the sample complexity scales with d). This highlights the benefit of ICL in adapting to low-dimensional structures of the target function class, by feature extraction via pretraining.

1.2. Related Works

Recent works (Zhang et al., 2023; Ahn et al., 2023; Mahankali et al., 2023; Wu et al., 2023; Zhang et al., 2024) studied the training of linear transformer to learn linear target functions in-context. Similar theoretical works are also established for transformers with SoftMax attention (Huang et al., 2023; Nichani et al., 2024; Chen et al., 2024). Our setting closely resembles (Kim & Suzuki, 2024), where a nonlinear MLP block is followed by a linear attention layer; the main difference is that we establish learnability for a concrete nonlinear function class, whereas (Kim & Suzuki, 2024) focused on global convergence of optimization. Finally, (Cheng et al., 2023) showed that transformers learn nonlinear functions in-context via a functional gradient update, but no statistical guarantees or optimization complexity were given.

The statistical and computational complexity of learning low-dimensional functions has been extensively studied. Typical target functions include single-index models (Ben Arous et al., 2021; Ba et al., 2022; Bietti et al., 2022; Mousavi-Hosseini et al., 2023; Damian et al., 2023; Ba et al., 2023) and multi-index models (Damian et al., 2022; Abbe et al., 2022; 2023; Bietti et al., 2023).

2. Problem Setting

Notations. We use boldface to represent vectors, matrices, and tensors. Let N be a nonnegative integer. Then, $[N]$ denotes the set $\{n \in \mathbb{Z} \mid 1 \leq n \leq N\}$. For a nonnegative integer i , the i -th Hermite polynomial is defined as $\text{He}_i(z) = e^{\frac{z^2}{2}} \frac{d^i}{dz^i} e^{-\frac{z^2}{2}}$. For a set S , $\text{Unif}(S)$ denotes the uniform distribution over S . We denote the unit sphere $\{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\| = 1\}$ by \mathbb{S}^{d-1} . $\tilde{O}(\cdot), \tilde{\Omega}(\cdot)$ represent $O(\cdot)$ and $\Omega(\cdot)$ notations where polylogarithmic terms are hidden. If necessary, we specify the targeted variables in O, Ω, \tilde{O} and $\tilde{\Omega}$, as $O_d(\cdot)$ for example. We write $a \lesssim b$ when there exists a constant c such that $a \leq cb$ holds.

2.1. Data Generating Process

First, we introduce the basic setting of in-context learning (Brown et al., 2020) of simple function classes as investigated in (Garg et al., 2022; Akyürek et al., 2022). In each inference (test) task, learners are fed a sequence of inputs and outputs $(\mathbf{x}_1, y_1, \dots, \mathbf{x}_N, y_N, \mathbf{x})$ referred to as *prompt*, where $\mathbf{x}_i, \mathbf{x} \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$. The labeled examples $\mathbf{X} = (\mathbf{x}_1 \ \dots \ \mathbf{x}_N) \in \mathbb{R}^{d \times N}$, $\mathbf{y} = (y_1 \ \dots \ y_N)^\top \in \mathbb{R}^N$ are called *context*, and \mathbf{x} is the *query*. We assume that the output y_i can be expressed as $y_i = f_*(\mathbf{x}_i) + \varsigma_i$, where f_* is the true function describing input-output relation and ς_i is label noise. The task is to predict the *response* $y = f_*(\mathbf{x}) + \varsigma$ corresponding to the query \mathbf{x} given the context, without updating model parameter. We specify the distribution of inputs and outputs as follows:

Assumption 1. *The prompt $(\mathbf{x}_1, y_1, \dots, \mathbf{x}_N, y_N, \mathbf{x})$ and the response y is generated as $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, \mathbf{x} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$, $y_i = f_*(\mathbf{x}_i) + \varsigma_i$, $y = f_*(\mathbf{x}) + \varsigma$, where $\varsigma_i, \varsigma \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(\{-\tau, \tau\})$, and the true function f_* is generated from the following distribution.*

- Let \mathcal{S} be an $r \leq d$ -dimensional linear subspace of \mathbb{R}^d . Draw a vector β uniformly from the unit sphere in \mathcal{S} , i.e., from $\{\beta \mid \beta \in \mathcal{S}, \|\beta\| = 1\}$.
- Draw Hermite coefficients $\{c_i\}_{i=2}^P$ from a distribution satisfying $\mathbb{E}[c_2] \neq 0$, $\sum_{i=2}^P c_i^2 = \Theta_{d,r}(1)$ (a.s.), and $(c_2, \dots, c_P) \neq (0, \dots, 0)$ (a.s.). Then, we define $f_*(\mathbf{x}_i) = \sigma_*(\langle \mathbf{x}_i, \beta \rangle)$, where $\sigma_*(z) = \sum_{i=2}^P c_i \text{He}_i(z)$.

Throughout the paper, we assume that $P \ll d, r$ and $r \ll d$: specifically, we take $P = O_{d,r}(1)$. This entails that the class of target functions is *low-dimensional*, and such structure should be adapted by the transformer via pretraining. As mentioned in the introduction section, the difficulty to learn a single-index model is governed by the *information exponent* of the link function σ_* : when we conduct the Hermite expansion as $\sigma_*(z) = \sum_{i \geq 0} c_i \frac{\text{He}_i(z)}{i!}$, then the information exponent is defined by $\min\{i \mid c_i \neq 0\}$. In

the case of Assumption 1, this is equal to the minimal i such that $c_i \neq 0$. Therefore, our case allows that exponent changes across the tasks, and models the situation where the difficulty to learn f_* varies across tasks.

Let $f(\mathbf{X}, \mathbf{y}, \mathbf{x})$ be an estimator for y . We evaluate the model performance by the *expected ICL risk* defined as

$$\mathcal{R}_N(f) := \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, y} [|f(\mathbf{X}, \mathbf{y}, \mathbf{x}) - y|], \quad (2.1)$$

where the expectation is taken over prompts with length N and responses.

2.2. Student Model: transformer with Nonlinear MLP Layer

As a learning model capable of in-context learning, we consider a transformer composed of a single-layer self-attention module preceded by an embedding module using a nonlinear two-layer perceptron. Let $\mathbf{E} \in \mathbb{R}^{d_e \times d_N}$ be an embedding matrix constructed using a prompt $(\mathbf{x}_1, y_1, \dots, \mathbf{x}_N, y_N, \mathbf{x})$. Using single-layer SoftMax-based self-attention module (Vaswani et al., 2017), the prediction of y is constructed as the right-bottom entry of $\mathbf{E} + \mathbf{W}^P \mathbf{W}^V \mathbf{E} \cdot \text{softmax}\left(\frac{(\mathbf{W}^K \mathbf{E})^\top \mathbf{W}^Q \mathbf{E}}{\rho}\right)$, where ρ is a normalization constant and $\mathbf{W}^K, \mathbf{W}^Q \in \mathbb{R}^{d_k \times d_e}$, $\mathbf{W}^V \in \mathbb{R}^{d_v \times d_e}$ and $\mathbf{W}^P \in \mathbb{R}^{d_e \times d_v}$ are parameters called key, query, value and projection matrix, respectively. In this paper, we take \mathbf{E} as

$$\mathbf{E} = \begin{bmatrix} \sigma(\mathbf{W}\mathbf{X} + \mathbf{b}) & \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) \\ \mathbf{y}^\top & 0 \end{bmatrix}. \quad (2.2)$$

Now $\sigma(\mathbf{W}\mathbf{X} + \mathbf{b})$ is a $m \times N$ matrix whose (i, j) -th element is $\sigma(\mathbf{w}_i^\top \mathbf{x}_j + b_i)$ and $\sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$ is a m -dimensional vector whose i -th element is $\sigma(\mathbf{w}_i^\top \mathbf{x} + b_i)$, where $\mathbf{w}_1, \dots, \mathbf{w}_m \in \mathbb{R}^d$ and $b_1, \dots, b_m \in \mathbb{R}$ are parameters and $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ is a nonlinear activation function. In this paper, we use $\sigma(z) = \text{ReLU}(z) = \max\{z, 0\}$.

In other words, we consider a two-layer neural network whose width is m , and take the output of each neuron $\sigma(\mathbf{w}^\top \mathbf{x} + b)$ at the hidden layer as the embedding. Using the output of a neural network as an embedding is adopted in some recent works (Guo et al., 2023; Kim & Suzuki, 2024). This MLP layer can extract features of the ground truth efficiently.

We further simplify the original self-attention module following the same line as (Wu et al., 2023; Zhang et al., 2023): we omit the softmax activation, set $\rho = d_N - 1 = N$, merge some parameter matrices and let some entries in the merged matrices as zero. We can show that the prediction of the output for \mathbf{x} by the simplified transformer can

Algorithm 1 Pretraining of transformer with MLP layer

- 1: **Input:** Learning rate η_1 , weight decay rate λ_1, λ_2 , prompt length N_1, N_2 and number of tasks T_1, T_2 .
- 2: Draw data $\{(\mathbf{x}_{t,1}, y_{t,1}, \dots, \mathbf{x}_{t,N_1}, y_{t,N_1}, \mathbf{x}_t, y_t)\}_{t=1}^{T_1}$ with prompt length N_1 and $\{(\mathbf{x}_{t,1}, y_{t,1}, \dots, \mathbf{x}_{t,N_2}, y_{t,N_2}, \mathbf{x}_t, y_t)\}_{t=T_1+1}^{T_1+T_2}$ with prompt length N_2 .
- 3: Initialize MLP weights as $\mathbf{w}_j^{(0)} \sim \text{Unif}(\mathbb{S}^{d-1})$ ($j \leq m/2$) and $\mathbf{w}_j^{(0)} = \mathbf{w}_{m-j}^{(0)}$ ($j > m/2$), biases as $b_j = 0$ ($j \in [m]$), and the attention matrix diagonally as $\Gamma_{j,j}^{(0)} \sim \text{Unif}(\{\pm 1\})$ ($j \leq m/2$) and $\Gamma_{j,j}^{(0)} = -\Gamma_{m-j, m-j}^{(0)}$ ($j > m/2$).
- 4: $\mathbf{w}_j^{(1)} \leftarrow \mathbf{w}_j^{(0)} - \eta_1 \left[\nabla_{\mathbf{w}_j} \hat{\mathcal{R}}(f) + \lambda_1 \mathbf{w}_j^{(0)} \right]$
- 5: Re-initialize \mathbf{b} as $b_j \sim \text{Unif}([-1, 1])$ ($j \in [m]$)
- 6: Find the minimizer $\mathbf{\Gamma}^*$ of $\min_{\mathbf{\Gamma}} \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} (y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \mathbf{\Gamma}, \mathbf{b}))^2 + \lambda_2 \|\mathbf{\Gamma}\|_F^2$
- 7: **Output:** parameters $(\mathbf{W}^{(1)}, \mathbf{\Gamma}^*, \mathbf{b})$

be written as

$$f(\mathbf{X}, \mathbf{y}, \mathbf{x}; \mathbf{W}, \mathbf{\Gamma}, \mathbf{b}) = \left\langle \frac{\mathbf{\Gamma} \sigma(\mathbf{W}\mathbf{X} + \mathbf{b}) \mathbf{y}}{N} \varphi(\mathbf{W}\mathbf{x} + \mathbf{b}) \right\rangle \quad (2.3)$$

where $\mathbf{\Gamma}$ is a parameter matrix. See Appendix E for the derivation of equation (2.3). We call $\mathbf{\Gamma}$ the *attention matrix*.

2.3. Pretraining: Empirical Risk Minimization via Gradient Descent

We pretrain parameters of the transformer (2.3) by the gradient-based algorithm, written in Algorithm 1. Throughout the paper, we assume that the width m is even.

In Algorithm 1, we first conduct one-step gradient descent for regularized empirical risk $\hat{\mathcal{R}}(f) := \frac{1}{T_1} \sum_{t=1}^{T_1} (y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}, \mathbf{\Gamma}, \mathbf{b}))^2 + \lambda_1 \|\mathbf{w}\|^2$ and update the MLP weight \mathbf{w} . Secondly, we conduct standard ridge regression with respect to the attention matrix $\mathbf{\Gamma}$. Note that the minimizer can be efficiently found because the optimization problem is convex with respect to $\mathbf{\Gamma}$. The symmetric initialization (line 3 in Algorithm 1) ensures that the output of the transformer is zero at initialization and removes the interaction between neurons: it is used in some recent works (Chizat et al., 2019; Damian et al., 2022).

3. Main Result: Transformer Learns Single-index Models In-context

3.1. Main Theorem

We state our main theorem to describe ICL ability of transformers. See Appendix D.3 for the proof.

Theorem 1. Assume that the data distribution is specified as Assumption 1. Consider pretraining transformer (2.3) via Algorithm 1 with $m \gtrsim r^P, T_1 = \tilde{\Omega}(d^4), N_1 = \tilde{\Omega}(d^2), \eta_1 = \Omega\left(\sqrt{\frac{d^3}{r}} \frac{1}{\text{poly log } d}\right)$ and $\lambda_1 = \eta_1^{-1}$. Then, with probability at least 0.99 over the data distribution and the random initialization, there exists $\lambda_2 > 0$ such that the ICL risk with prompt length N^* is upper bounded as

$$\mathcal{R}_{N^*}(f) - \tau \lesssim \frac{r^{3P/2}}{\sqrt{m}} + \text{polylog}(T_2) \frac{r^{5P/2}}{\sqrt{T_2}} d\sqrt{d} + \sqrt{r^{5P}} \sqrt{\frac{r^2}{N_2} + \frac{r^2}{N^*}},$$

where $f = f(\mathbf{X}, \mathbf{y}, \mathbf{x}; \mathbf{W}^{(1)}, \mathbf{\Gamma}^*, \mathbf{b})$ and dependence on $\text{poly log}(d)$ in the right-hand side is ignored.

We can evaluate the sample and prompt-length complexity of Algorithm 1 using Theorem 1: to achieve $\mathcal{R}_{N^*}(f) \leq \tau + \epsilon$ for given $\epsilon > 0$, it is sufficient to set $m = \tilde{\Omega}(r^{3P}), T_2 = \tilde{\Omega}(r^{5P}d^3)$ and $N_2, N^* = \tilde{\Omega}(r^{5P+2})$. Most importantly, at the test time, the required prompt length N^* only depends on the inner dimension r , up to polylogarithmic terms. We emphasize that the nonlinear link function σ_* and the true direction $\beta \in \mathcal{S}$ varies across tasks, and then the difficulty to learn the input-output relation varies. Nevertheless, Theorem 1 shows that transformers can learn the relation on a short prompt which does not scale with d .

As we have nonlinearity in our model, empirical risk minimization problem becomes nonconvex and establishing optimization guarantee via gradient descent becomes more difficult than the linear setting. We established the guarantee by making use of one-step gradient descent, which has been considered in literatures of feature learning (Ba et al., 2022; Damian et al., 2022) and which has successfully yielded end-to-end optimization guarantee. Now we have to extend these results utilizing discussions specific to ICL as we have an attention layer.

Comparison against baseline methods. Another aspect specific to our result is that we can say transformers can outperform learning algorithms that directly act on test prompts: a lot of works (Zhang et al., 2023; Ahn et al., 2023; Mahankali et al., 2023; Wu et al., 2023; Zhang et al., 2024) discussed linear transformers, but due to the linearity of the studied transformer, ICL cannot outperform linear estimators on the test prompt.

Let us concretely compare our result with algorithms acting on the test prompt: these algorithms read each context $(\mathbf{x}_1, y_1, \dots, \mathbf{x}_N, y_N)$ and update their parameters, then make a prediction of the response y for the query \mathbf{x} . This is simply a regression problem to estimate a single-index model $f_*(\mathbf{x}) = \sigma_*(\langle \mathbf{x}, \beta \rangle)$ using samples

$(\mathbf{x}_1, y_1, \dots, \mathbf{x}_N, y_N)$. Thus, the required prompt length is equal to the number of samples needed to learn the single-index model. Sample complexities for various algorithms to learn single-index models were shown: for linear methods such as kernel methods, $d^{\Omega(P)}$ samples are necessary to achieve ϵ -error, that is, to achieve $\mathbb{E}_{\mathbf{x}}[|f(\mathbf{x}) - y|] \leq \tau + \epsilon$ for given $\epsilon > 0$ where f is the estimator (Ghorbani et al., 2021; Donhauser et al., 2021). On the other hand, neural networks can learn single-index models with $d^{\Omega(k^*)}$ samples (Ben Arous et al., 2021; Bietti et al., 2022) by gradient descent, where k^* is the information exponent of σ_* , i.e., the minimal i such that $c_i \neq 0$. However, for the easiest case where $c_2 \neq 0$ and thus $k^* = 2$, the sample complexity is at least linear in the ambient dimension d . Moreover, a lower bound is known for a general framework of algorithms called CSQ algorithm, which includes stochastic gradient descent on neural networks: any CSQ algorithm needs $\Omega(d^{k^*/2})$ samples to achieve ϵ -error (Damian et al., 2022). Therefore, if we run these learning algorithms on each test prompt, they require the prompt length which depends on $\text{poly}(d)$. Thus when $r \ll d$, pretrained transformers can capture the input-output relation in-context with a significantly shorter prompt length.

Discussion on the Mechanism. The proof of Theorem 1 is composed of several parts. First, we show that after one-step gradient descent, $\mathbf{w}^{(1)}$ aligns with \mathcal{S} , i.e., $\mathbf{w}^{(1)}$ is almost contained in \mathcal{S} (See Appendix B for details). Secondly, we show in Appendix C that there is an attention matrix $\mathbf{\Gamma}$ such that the entire transformer approximates the true function well, which is crucial in the generalization error analysis (Appendix D). Thus, we can say that the MLP layer succeeds in “memorizing” the low-dimensional feature space even with the single step gradient descent, and the attention matrix works to approximate the link function correctly.

4. Conclusion and Future Direction

In this work, we studied the ICL ability of transformers and showed that they can adapt to the intrinsic low-dimensional structure of nonlinear true functions, and then outperformed algorithms working directly on the test prompt, in that the required prompt length only scaled with the inner dimension $r \ll d$.

There are several important future challenges. It is intriguing to explore the regime $r \approx d$, where memorizing the feature space \mathcal{S} by pretraining is no longer meaningful. Extending our result to multi-index models is also an interesting future direction. Finally, we used a nonlinear MLP layer with a linear self-attention module, but it is important to study the ICL on nonlinear functions via nonlinear self-attention such as softmax-based self-attention modules.

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Abbe, E., Adsera, E. B., and Misiakiewicz, T. The merged-staircase property: a necessary and nearly sufficient condition for sgd learning of sparse functions on two-layer neural networks. In *Conference on Learning Theory*, pp. 4782–4887. PMLR, 2022.
- Abbe, E., Adsera, E. B., and Misiakiewicz, T. SGD learning on neural networks: leap complexity and saddle-to-saddle dynamics. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 2552–2623. PMLR, 2023.
- Ahn, K., Cheng, X., Daneshmand, H., and Sra, S. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2023.
- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., Wu, D., and Yang, G. High-dimensional asymptotics of feature learning: How one gradient step improves the representation. *Advances in Neural Information Processing Systems*, 35: 37932–37946, 2022.
- Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., and Wu, D. Learning in the presence of low-dimensional structure: A spiked random matrix perspective. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=HlIAoCHDWW>.
- Ben Arous, G., Gheissari, R., and Jagannath, A. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *The Journal of Machine Learning Research*, 22(1):4788–4838, 2021.
- Bietti, A., Bruna, J., Sanford, C., and Song, M. J. Learning single-index models with shallow neural networks. *Advances in Neural Information Processing Systems*, 35: 9768–9783, 2022.
- Bietti, A., Bruna, J., and Pillaud-Vivien, L. On learning Gaussian multi-index models with gradient flow. *arXiv preprint arXiv:2310.19793*, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Chen, S., Sheen, H., Wang, T., and Yang, Z. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv preprint arXiv:2402.19442*, 2024.
- Cheng, X., Chen, Y., and Sra, S. Transformers implement functional gradient descent to learn non-linear functions in context. *arXiv preprint arXiv:2312.06528*, 2023.
- Chizat, L., Oyallon, E., and Bach, F. On lazy training in differentiable programming. *Advances in neural information processing systems*, 32, 2019.
- Dai, D., Sun, Y., Dong, L., Hao, Y., Ma, S., Sui, Z., and Wei, F. Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers. *arXiv preprint arXiv:2212.10559*, 2022.
- Damian, A., Lee, J., and Soltanolkotabi, M. Neural networks can learn representations with gradient descent. In *Conference on Learning Theory*, pp. 5413–5452. PMLR, 2022.
- Damian, A., Nichani, E., Ge, R., and Lee, J. D. Smoothing the landscape boosts the signal for SGD: Optimal sample complexity for learning single index models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=73XPopmbXH>.
- Donhauser, K., Wu, M., and Yang, F. How rotational invariance of common kernels prevents generalization in high dimensions. In *International Conference on Machine Learning*, pp. 2804–2814. PMLR, 2021.
- Garg, S., Tsipras, D., Liang, P. S., and Valiant, G. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Ghorbani, B., Mei, S., Misiakiewicz, T., and Montanari, A. Linearized two-layers neural networks in high dimension. *The Annals of Statistics*, 49(2):1029–1054, 2021.
- Götze, F., Sambale, H., and Sinulis, A. Concentration inequalities for polynomials in α -sub-exponential random variables. 2021.
- Guo, T., Hu, W., Mei, S., Wang, H., Xiong, C., Savarese, S., and Bai, Y. How do transformers learn in-context beyond simple functions? a case study on learning with representations. *arXiv preprint arXiv:2310.10616*, 2023.
- Huang, Y., Cheng, Y., and Liang, Y. In-context convergence of transformers. *arXiv preprint arXiv:2310.05249*, 2023.

- 275 Kim, J. and Suzuki, T. Transformers learn nonlinear fea-
276 tures in context: Nonconvex mean-field dynamics on the
277 attention landscape. *arXiv preprint arXiv:2402.01258*,
278 2024.
- 279 Mahankali, A., Hashimoto, T. B., and Ma, T. One step
280 of gradient descent is provably the optimal in-context
281 learner with one layer of linear self-attention. *arXiv*
282 *preprint arXiv:2307.03576*, 2023.
- 284 Maurer, A. A vector-contraction inequality for rademacher
285 complexities. In *Algorithmic Learning Theory: 27th In-*
286 *ternational Conference, ALT 2016, Bari, Italy, October*
287 *19-21, 2016, Proceedings 27*, pp. 3–17. Springer, 2016.
- 288 Mousavi-Hosseini, A., Park, S., Girotti, M., Mitliagkas,
289 I., and Erdogdu, M. A. Neural networks efficiently
290 learn low-dimensional representations with SGD. In *The*
291 *Eleventh International Conference on Learning Repre-*
292 *sentations*, 2023.
- 294 Nichani, E., Damian, A., and Lee, J. D. How transform-
295 ers learn causal structure with gradient descent. *arXiv*
296 *preprint arXiv:2402.14735*, 2024.
- 297 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
298 L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I.
299 Attention is all you need. In Guyon, I., Luxburg, U. V.,
300 Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S.,
301 and Garnett, R. (eds.), *Advances in Neural Information*
302 *Processing Systems*, volume 30. Curran Associates, Inc.,
303 2017. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
304 [cc/paper_files/paper/2017/file/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
305 [3f5ee243547dee91fbd053c1c4a845aa-Paper.](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
306 [pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 308 Von Oswald, J., Niklasson, E., Randazzo, E., Sacramento,
309 J., Mordvintsev, A., Zhmoginov, A., and Vladymyrov,
310 M. Transformers learn in-context by gradient descent.
311 In *International Conference on Machine Learning*, pp.
312 35151–35174. PMLR, 2023.
- 314 Wainwright, M. J. *High-dimensional statistics: A non-*
315 *asymptotic viewpoint*, volume 48. Cambridge university
316 press, 2019.
- 317 Wu, J., Zou, D., Chen, Z., Braverman, V., Gu, Q., and
318 Bartlett, P. L. How many pretraining tasks are needed for
319 in-context learning of linear regression? *arXiv preprint*
320 *arXiv:2310.08391*, 2023.
- 322 Zhang, R., Frei, S., and Bartlett, P. L. Trained trans-
323 formers learn linear models in-context. *arXiv preprint*
324 *arXiv:2306.09927*, 2023.
- 325 Zhang, R., Wu, J., and Bartlett, P. L. In-context learn-
326 ing of a linear transformer block: Benefits of the mlp
327 component and one-step gd initialization. *arXiv preprint*
328 *arXiv:2402.14951*, 2024.
- 329

A. Preliminaries

By coordinate transformation, without loss of generality we can assume that $\mathcal{S} = \{(x_1, \dots, x_r, 0, \dots, 0) \mid x_1, \dots, x_r \in \mathbb{R}\}$ and $\beta \sim \text{Unif}(\mathbb{S}^{r-1})$, i.e., $\beta \sim \text{Unif}(\{(\beta_1, \dots, \beta_r) \mid \beta_1^2 + \dots + \beta_r^2 = 1\})$. Therefore, we assume this in the entire proof. For a vector w , we use $w_{a:b}$ for $a \leq b$ to denote the vector $[w_a, w_{a+1}, \dots, w_b]^\top$.

A.1. Definition of the Term ‘‘With High Probability’’

We sometimes use the term ‘‘an event A occurs with high probability’’. Now we explain the definition of this.

Definition 2. We say that an event A occurs with high probability when there exists a sufficiently large constant C^* which does not depend on the ambient dimension d and

$$\Pr[A] \geq 1 - \text{poly}(d)d^{-C^*}$$

holds.

Note that if A_1, \dots, A_M occurs with high probability where $M = \text{poly}(d)$, then $A_1 \cap \dots \cap A_M$ occurs with high probability.

A.2. Tensors

In this paper, a k -tensor is a multidimensional array which has k indices: for example, matrices are 2-tensors. Let A be a k -tensor. A_{i_1, \dots, i_k} denotes (i_1, \dots, i_k) -th entry of A . Let A be a k -tensor and B be a l -tensor where $k \geq l$. $A(B)$ denotes $k - l$ tensor whose (i_1, \dots, i_{k-l}) -th entry is

$$A(B)_{i_1, \dots, i_{k-l}} = \sum_{j_1, \dots, j_l} A_{i_1, \dots, i_{k-l}, j_1, \dots, j_l} B_{j_1, \dots, j_l},$$

and is defined only when sizes are compatible. If $k = l$, we sometimes write $A(B)$ as $A \circ B$. Let $v \in \mathbb{R}^d$ be a vector and k be a positive integer. Then, $v^{\otimes k} \in \mathbb{R}^{d \times \dots \times d}$ denotes a k -tensor whose (i_1, \dots, i_k) -th entry is $v_{i_1} \dots v_{i_k}$.

Let $f(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ be a d -variable function. A k -tensor $\nabla^k f(x)$ is defined as

$$(\nabla^k f(x))_{i_1, \dots, i_k} = \frac{\partial}{\partial x_{i_1}} \dots \frac{\partial}{\partial x_{i_k}} f(x).$$

A.3. Hermite Polynomials

We frequently use (probabilists’) Hermite polynomials, which is defined by $\text{He}_i(z) = e^{\frac{z^2}{2}} \frac{d^i}{dz^i} e^{-\frac{z^2}{2}}$, where i is a non-negative integer. Hermite polynomials have orthogonality, in that $\mathbb{E}_{z \sim \mathcal{N}(0,1)}[\text{He}_i(z)\text{He}_j(z)] = i! \delta_{i,j}$. The Hermite expansion for $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is defined as $\sigma(z) = \sum_{i \geq 0} \frac{a_i}{i!} \text{He}_i(z)$ where $a_i = \mathbb{E}_{z \sim \mathcal{N}(0,1)}[\sigma(z)\text{He}_i(z)]$. Similarly, the multivariate Hermite expansion for $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as $f(z) = \sum_{i_1 \geq 0, \dots, i_d \geq 0} \frac{a_{i_1, \dots, i_d}}{(i_1)! \dots (i_d)!} \text{He}_{i_1}(z_1) \dots \text{He}_{i_d}(z_d)$, where $a_{i_1, \dots, i_d} = \mathbb{E}_{z_i, \dots, z_d \sim \mathcal{N}(0,1)}[f(z)\text{He}_{i_1}(z_1) \dots \text{He}_{i_d}(z_d)]$. The coefficient a_{i_1, \dots, i_d} can also be obtained by $a_{i_1, \dots, i_d} =$

$$\mathbb{E}_{z_i, \dots, z_d \sim \mathcal{N}(0,1)} \left[\frac{\partial^{i_1}}{\partial z_1^{i_1}} \dots \frac{\partial^{i_d}}{\partial z_d^{i_d}} f(z) \right].$$

The lemma below is useful to find a basis of the set of true functions.

Lemma 3. Suppose $\beta \in \mathcal{S}$. Then,

$$\text{He}_p(\langle x, \beta \rangle) = \sum_{p_1 \geq 0, \dots, p_r \geq 0}^{p_1 + \dots + p_r = p} \frac{(p_1 + \dots + p_r)!}{p_1! \dots p_r!} \cdot \beta_1^{p_1} \dots \beta_r^{p_r} \cdot \text{He}_{p_1}(x_1) \dots \text{He}_{p_r}(x_r)$$

holds.

Proof. Note that $\mathbb{E}_{z_i, \dots, z_r \sim \mathcal{N}(0,1)} \left[\frac{\partial^{i_1}}{\partial z_1^{i_1}} \dots \frac{\partial^{i_r}}{\partial z_r^{i_r}} \text{He}_p(\langle x, \beta \rangle) \right]$ is nonzero only when $i_1 + \dots + i_r = p$. When $i_1 + \dots + i_r = p$, then $\mathbb{E}_{z_1, \dots, z_r \sim \mathcal{N}(0,1)} \left[\frac{\partial^{i_1}}{\partial z_1^{i_1}} \dots \frac{\partial^{i_r}}{\partial z_r^{i_r}} \text{He}_p(\langle x, \beta \rangle) \right] = p! \beta_1^{i_1} \dots \beta_r^{i_r}$ holds. Then, from the multivariate Hermite expansion we obtain the claim. \square

B. Proofs for MLP Layer

First, note that $f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(0)}, \mathbf{\Gamma}^{(0)}, \mathbf{b}) = 0$ at initialization. Then, from line 4 of Algorithm 1, for each $j \in [m]$,

$$\begin{aligned} \mathbf{w}_j^{(1)} &= 2\eta_1 \frac{1}{T} \sum_{t=1}^T y_t \nabla_{\mathbf{w}_j^{(0)}} f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(0)}, \mathbf{\Gamma}^{(0)}, \mathbf{b}) \\ &= 2\eta_1 \mathbf{\Gamma}_{j,j}^{(0)} \left(\sum_{t=1}^T \frac{1}{T} y_t \sigma(\mathbf{w}_j^{(0)\top} \mathbf{x}_t) \frac{1}{N} \sum_{i=1}^N y_{t,i} \sigma'(\mathbf{w}_j^{(0)\top} \mathbf{x}_{t,i}) \mathbf{x}_{t,i} \right. \\ &\quad \left. + \sum_{t=1}^T \frac{1}{T} y_t \sigma'(\mathbf{w}_j^{(0)\top} \mathbf{x}_t) \mathbf{x}_t \frac{1}{N} \sum_{i=1}^N y_{t,i} \sigma(\mathbf{w}_j^{(0)\top} \mathbf{x}_{t,i}) \right) \end{aligned}$$

holds because we assumed $\eta_1 = \lambda_1^{-1}$ (for simplicity let $T := T_1$ and $N := N_1$ in this section). Here $\mathbf{x}_{t,i}$ is the i -th column of \mathbf{X}_t and $y_{t,i}$ is the i -th element of \mathbf{y}_t . Now let

$$\begin{aligned} \mathbf{g}_T(\mathbf{w}) &:= \left(\sum_{i=1}^T \frac{1}{T} y_i \sigma(\mathbf{w}^\top \mathbf{x}_i) \frac{1}{N} \sum_{i=1}^N y_{t,i} \sigma'(\mathbf{w}^\top \mathbf{x}_{t,i}) \mathbf{x}_{t,i} + \sum_{i=1}^T \frac{1}{T} y_i \sigma'(\mathbf{w}^\top \mathbf{x}_i) \mathbf{x}_i \frac{1}{N} \sum_{i=1}^N y_{t,i} \sigma(\mathbf{w}^\top \mathbf{x}_{t,i}) \right) \end{aligned}$$

and $\mathbf{g}(\mathbf{w}) = \mathbb{E}[\mathbf{g}_T(\mathbf{w})]$, where the expectation is taken with respect to the data distribution. Note that $\mathbf{w}_j^{(1)} = 2\eta_1 \mathbf{\Gamma}_{j,j}^{(0)} \mathbf{g}_T(\mathbf{w}_j^{(0)})$.

In this section we make

- asymptotic expansion of $\mathbf{g}(\mathbf{w})$, and
- uniform upper bound for the difference $\|\mathbf{g}(\mathbf{w}) - \mathbf{g}_T(\mathbf{w})\|$.

Asymptotic Expansion of $\mathbf{g}(\mathbf{w})$. First, note that

$$\mathbf{g}(\mathbf{w}) = 2\mathbb{E}_y[\mathbb{E}_x[y\sigma'(\mathbf{w}^\top \mathbf{x})\mathbf{x}]\mathbb{E}_x[y\sigma(\mathbf{w}^\top \mathbf{x})]],$$

where $\mathbb{E}_y[\cdot]$ means the expectation with respect to the distribution of β and $\{c_i\}$.

Now let $\sigma(z) = \sum_{i \geq 0} a_i \frac{\text{He}_i(z)}{i!}$ be the Hermite expansion of student activation. The asymptotic expansions of $\mathbb{E}_x[y\sigma'(\mathbf{w}^\top \mathbf{x})\mathbf{x}]$ and $\mathbb{E}_x[y\sigma(\mathbf{w}^\top \mathbf{x})]$ are known as follows:

Lemma 4. *It holds that*

$$\begin{aligned} \mathbb{E}_x[y\sigma'(\mathbf{w}^\top \mathbf{x})\mathbf{x}] &= \sum_{k=1}^P \frac{a_{k+1} \mathbb{E}_x[\nabla^{k+1} f_*(\mathbf{x})](\mathbf{w}^{\otimes k})}{k!} + \mathbf{w} \sum_{k=2}^P \frac{a_{k+2} \mathbb{E}_x[\nabla^k f_*(\mathbf{x})](\mathbf{w}^{\otimes k})}{k!}, \\ \mathbb{E}_x[y\sigma(\mathbf{w}^\top \mathbf{x})] &= \sum_{k=2}^P \frac{a_k \mathbb{E}_x[\nabla^k f_*(\mathbf{x})](\mathbf{w}^{\otimes k})}{k!}. \end{aligned}$$

Proof. It is obtained from the proof of Lemma 7 in (Damian et al., 2022). □

Then, $\mathbf{g}(\mathbf{w})$ can be expanded as

$$\begin{aligned} \mathbf{g}(\mathbf{w}) &= 2\mathbb{E}_y \left[\left(\sum_{k=1}^P \frac{a_{k+1} \mathbb{E}_x[\nabla^{k+1} f_*(\mathbf{x})](\mathbf{w}^{\otimes k})}{k!} \right. \right. \\ &\quad \left. \left. + \mathbf{w} \sum_{k=2}^P \frac{a_{k+2} \mathbb{E}_x[\nabla^k f_*(\mathbf{x})](\mathbf{w}^{\otimes k})}{k!} \right) \left(\sum_{k=2}^P \frac{a_k \mathbb{E}_x[\nabla^k f_*(\mathbf{x})](\mathbf{w}^{\otimes k})}{k!} \right) \right] \end{aligned}$$

$$\begin{aligned}
 &=: 2\mathbb{E}_y \left[a_2 \mathbb{E}_{\mathbf{x}} [\nabla^2 f_*(\mathbf{x})](\mathbf{w}) \frac{a_2 \mathbb{E}_{\mathbf{x}} [\nabla^2 f_*(\mathbf{x})](\mathbf{w}^{\otimes 2})}{2} + \mathbf{s}(y, \mathbf{w}) \right] \\
 &= a_2^2 \mathbb{E}_y [\mathbb{E}_{\mathbf{x}} [\nabla^2 f_*(\mathbf{x})](\mathbf{w}) \mathbb{E}_{\mathbf{x}} [\nabla^2 f_*(\mathbf{x})](\mathbf{w}^{\otimes 2})] + 2\mathbb{E}_y [\mathbf{s}(y, \mathbf{w})].
 \end{aligned}$$

The main term is proportional to $\mathbb{E}_y [(\mathbf{H}_{f_*} \mathbf{w})(\mathbf{H}_{f_*} \circ \mathbf{w}^{\otimes 2})]$, where

$$\mathbf{H}_{f_*} = \mathbb{E}_{\mathbf{x}} [\nabla^2 f_*(\mathbf{x})] = 2c_2 \boldsymbol{\beta} \boldsymbol{\beta}^\top$$

be the expected Hessian of f_* . Let us calculate this main term explicitly. Recall that we assumed that $\boldsymbol{\beta} \sim \text{Unif}(\mathbb{S}^{r-1})$. Note that by letting $\boldsymbol{\beta}' \sim \mathcal{N}(0, \boldsymbol{\Sigma}_{\boldsymbol{\beta}})$ where

$$\boldsymbol{\Sigma}_{\boldsymbol{\beta}} = \begin{pmatrix} \mathbf{I}_r & 0 \\ 0^\top & \mathbf{O}_{d-r} \end{pmatrix},$$

and $z \sim \chi_r$ independent from $\boldsymbol{\beta}$, then $\boldsymbol{\beta}' \sim \boldsymbol{\beta}z$ holds.

Now it holds that

$$\begin{aligned}
 \mathbb{E}_y [(\mathbf{H}_{f_*} \mathbf{w})(\mathbf{H}_{f_*} \circ \mathbf{w}^{\otimes 2})] &= \mathbb{E}_{\boldsymbol{\beta}, c_2} [4c_2^2 (\boldsymbol{\beta} \boldsymbol{\beta}^\top \mathbf{w})(\boldsymbol{\beta} \boldsymbol{\beta}^\top \circ \mathbf{w}^{\otimes 2})] \\
 &= 4\mathbb{E}[c_2^2] \mathbb{E}_{\boldsymbol{\beta}} [\boldsymbol{\beta}^{\otimes 4}] \mathbf{w}^{\otimes 3}.
 \end{aligned}$$

Furthermore, as

$$\mathbb{E}_{\boldsymbol{\beta}} [\beta_i \beta_j \beta_k \beta_l] = \begin{cases} \frac{3}{\mathbb{E}_{z \sim \chi_r} [z^4]} & (i = j = k = l \leq r) \\ \frac{1}{\mathbb{E}_{z \sim \chi_r} [z^4]} & (i = j \leq r, k = l \leq r, i \neq k \text{ or } i = k \leq r, j = l \leq r, i \neq j, \\ & \text{or } i = l \leq r, j = k \leq r, i \neq j) \\ 0 & \text{(otherwise)} \end{cases}$$

it follows that

$$(\mathbb{E}[c_2^2] \mathbb{E}_{\boldsymbol{\beta}} [\boldsymbol{\beta}^{\otimes 4}] \mathbf{w}^{\otimes 3})_i = \begin{cases} \mathbb{E}[c_2^2] \frac{3}{r(r+2)} w_i \sum_{j=1}^r w_j^2 & (i \leq r) \\ 0 & (i > r) \end{cases}.$$

Then, we arrive at

$$\mathbb{E}_y [(\mathbf{H}_{f_*} \mathbf{w})(\mathbf{H}_{f_*} \circ \mathbf{w}^{\otimes 2})] = \mathbb{E}[c_2^2] \frac{12}{r(r+2)} \|\mathbf{w}_{1:r}\|_2^2 \mathbf{w}_{1:r}.$$

Next, we upper bound the residual term $\mathbb{E}_y [\mathbf{s}(y, \mathbf{w})]$.

Lemma 5.

$$\sup_{\mathbf{w} \sim \mathbb{S}^{d-1}} \|\mathbb{E}_y [\mathbf{s}(y, \mathbf{w})]\| = O\left(\sqrt{\frac{r^2}{d^4}}\right)$$

holds.

Proof. Note that

$$\sup_{\mathbf{w} \sim \mathbb{S}^{d-1}} \|\mathbb{E}_y [\mathbf{s}(y, \mathbf{w})]\| \leq \sup_{\mathbf{w} \sim \mathbb{S}^{d-1}} \mathbb{E}_y [\|\mathbf{s}(y, \mathbf{w})\|^2]^{1/2}.$$

Then, by Minkowski's inequality, it suffices to show that

$$\mathbb{E}_y \left[\left\| \frac{a_{k+1} \mathbb{E}[\nabla^{k+1} f^t](\mathbf{w}^{\otimes k})}{k!} \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right\|^2 \right]^{1/2} \lesssim \sqrt{\frac{r^2}{d^4}} \quad (k \geq 1, l \geq 2, (k, l) \neq (1, 2))$$

and

$$\mathbb{E}_y \left[\left\| \mathbf{w} \frac{a_{k+2} \mathbb{E}[\nabla^k f^t](\mathbf{w}^{\otimes k})}{k!} \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right\|^2 \right]^{1/2} \lesssim \sqrt{\frac{r^2}{d^4}} \quad (k \geq 2, l \geq 2).$$

The former is obtained by

$$\begin{aligned} & \mathbb{E}_y \left[\left\| \frac{a_{k+1} \mathbb{E}[\nabla^{k+1} f^t](\mathbf{w}^{\otimes k})}{k!} \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right\|^2 \right]^{1/2} \\ & \leq \mathbb{E}_y \left[\left\| \frac{a_{k+1} \mathbb{E}[\nabla^{k+1} f^t](\mathbf{w}^{\otimes k})}{k!} \right\|^4 \right]^{1/4} \left[\left\| \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right\|^4 \right]^{1/4} \\ & \lesssim \sqrt{\frac{r^{\lfloor k/2 \rfloor}}{d^k}} \sqrt{\frac{r^{\lfloor l/2 \rfloor}}{d^l}}. \end{aligned}$$

Now, for deriving the last line, we used Corollary 9 and Lemma 24 in (Damian et al., 2022). The latter can be derived by following the same line. \square

Bounding the Difference between Empirical and Population Gradient. We upper bound $\|g_T(\mathbf{w}) - g(\mathbf{w})\|$ by extending Lemma 19 in (Damian et al., 2022). In the paper they bound the difference between empirical and population gradient of a two-layer fully-connected neural network. However, in our case we have an attention module and nonlinear activation appears twice in the gradient. This yields the need for using concentration argument in a ‘‘nested’’ way.

Lemma 6. *Let*

$$\mathbf{z}_1(\mathbf{w}) = \sum_{t=1}^T \frac{1}{T} y_t \sigma(\mathbf{w}^\top \mathbf{x}_t) \frac{1}{N} \sum_{i=1}^N y_{t,i} \sigma'(\mathbf{w}^\top \mathbf{x}_{t,i}) \mathbf{x}_{t,i}$$

and

$$\mathbf{z}_2(\mathbf{w}) = \sum_{t=1}^T \frac{1}{T} y_t \sigma'(\mathbf{w}^\top \mathbf{x}_t) \mathbf{x}_t \frac{1}{N} \sum_{i=1}^N y_{t,i} \sigma(\mathbf{w}^\top \mathbf{x}_{t,i}).$$

Then, with high probability over the data distribution,

$$\begin{aligned} \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \|\mathbf{z}_1(\mathbf{w}) - \mathbb{E}[\mathbf{z}_1(\mathbf{w})]\| &= \tilde{O} \left(\left(\sqrt{\frac{1}{d}} + \sqrt{\frac{d}{N}} \right) \sqrt{\frac{d}{T}} \right), \\ \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \|\mathbf{z}_2(\mathbf{w}) - \mathbb{E}[\mathbf{z}_2(\mathbf{w})]\| &= \tilde{O} \left(\left(\frac{1}{d} + \sqrt{\frac{d}{N}} \right) \sqrt{\frac{d}{T}} \right) \end{aligned}$$

holds.

Proof. First, consider $\mathbf{z}_1(\mathbf{w})$. Let $\iota = C \log N$, where C is a sufficiently large constant. From Lemma 19 in (Damian et al., 2022), with probability at least $1 - 2N e^{-\iota}$ it holds that

$$\frac{1}{N} \sum_{i=1}^N y_{t,i} \sigma'(\mathbf{w}^\top \mathbf{x}_{t,i}) \mathbf{x}_{t,i} \leq \mathbb{E}[y_{t,i} \sigma'(\mathbf{w}^\top \mathbf{x}_{t,i}) \mathbf{x}_{t,i}] + C_t \sqrt{\frac{d \iota^{P+1}}{N}} \quad \text{for all } \mathbf{w} \in \mathbb{S}^{d-1} \quad (\text{B.1})$$

with fixed t , where C_t is a constant which only depends on t and P . Let R_t be the right hand side of (B.1). Note that $R_t = \tilde{O} \left(\sqrt{\frac{1}{d}} + \sqrt{\frac{d \iota^{P+1}}{N}} \right)$ with high probability. Moreover, from Lemma 17 in (Damian et al., 2022), $y_t \leq (C'_t \iota)^{P/2} =: R'_t$ holds with high probability. In the following, we assume that

- (B.1) holds,
- $y_t \leq R'_t$ holds with all t , and
- $\|\mathbf{x}_{t,i}\|, \|\mathbf{x}_t\| \leq a\sqrt{d}$ holds for all t and i , where a is a sufficiently large constant.

which occurs with high probability. Let

$$\tilde{z}_1(\mathbf{w}) = \frac{1}{T} \sum_{i=1}^T y_t \sigma(\mathbf{w}^\top \mathbf{x}_t) \mathbf{v}_t(\mathbf{w}) \mathbf{1}[v_t(\mathbf{w}) \leq R_t \wedge y_t \leq R'_t],$$

where $\mathbf{v}_t(\mathbf{w}) := \frac{1}{N} \sum_{i=1}^N y_{t,i} \sigma'(\mathbf{w}^\top \mathbf{x}_{t,i}) \mathbf{x}_{t,i}$. From the assumption above, $\tilde{z}_1 = z_1$ holds uniformly over all $\mathbf{w} \in \mathbb{S}^{d-1}$. Let us bound $\sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \|\mathbf{z}_1(\mathbf{w}) - \mathbb{E}[\mathbf{z}_1(\mathbf{w})]\|$ following the same line as the proof of Lemma 19 in (Damian et al., 2022). As

$$\mathbf{z}_1(\mathbf{w}) - \mathbb{E}[\mathbf{z}_1(\mathbf{w})] = (\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]) - (\mathbb{E}[\mathbf{z}_1(\mathbf{w})] - \mathbb{E}[\tilde{z}_1(\mathbf{w})]),$$

it suffices to bound $\sup_{\mathbf{w}} \|\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\|$ and $\sup_{\mathbf{w}} \|\mathbb{E}[\mathbf{z}_1(\mathbf{w})] - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\|$. First,

$$\begin{aligned} & \sup_{\mathbf{w}} \|\mathbb{E}[\mathbf{z}_1(\mathbf{w})] - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\| \\ &= \sup_{\mathbf{w}} \|\mathbb{E}[y_t \sigma(\mathbf{w}^\top \mathbf{x}_t) \mathbf{v}_t(\mathbf{w}) \mathbf{1}[v_t(\mathbf{w}) > R_t \vee y_t > R'_t]]\| \\ &\leq \sup_{\mathbf{w}} \sqrt{\mathbb{E}[\|\mathbf{v}_t(\mathbf{w})\|^2]} \cdot \sup_{\mathbf{w}} \sqrt{\mathbb{E}[y_t^2 \sigma(\mathbf{w}^\top \mathbf{x}_t)^2 \mathbf{1}[v_t(\mathbf{w}) > R_t \vee y_t > R'_t]^2]} \\ &\leq \sup_{\mathbf{w}} \mathbb{E}[\|\mathbf{v}_t(\mathbf{w})\|^2]^{1/2} \cdot \sup_{\mathbf{w}} \mathbb{E}[y_t^4]^{1/4} \mathbb{E}[\sigma(\mathbf{w}^\top \mathbf{x}_t)^8]^{1/8} \mathbb{P}[v_t(\mathbf{w}) > R_t \vee y_t > R'_t]^{1/8}. \end{aligned}$$

Then, as $\sup_{\mathbf{w}} \mathbb{E}[\|\mathbf{v}_t(\mathbf{w})\|^2]^{1/2}$, $\sup_{\mathbf{w}} \mathbb{E}[y_t^4]^{1/4}$ and $\sup_{\mathbf{w}} \mathbb{E}[\sigma(\mathbf{w}^\top \mathbf{x}_t)^8]^{1/8}$ are polynomial in d , by setting C sufficiently large, it follows that $\sup_{\mathbf{w}} \|\mathbb{E}[\mathbf{z}_1(\mathbf{w})] - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\| = O(d^{-C^*})$ for any $C^* \geq 1$ with high probability.

Second, let us bound $\sup_{\mathbf{w}} \|\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\|$. From Lemmas 18 and 19 in (Damian et al., 2022), there exists ϵ -covering \mathcal{N}_ϵ of \mathbb{S}^{d-1} with $|\mathcal{N}_\epsilon| \leq e^{C_1 d \log(NT/\epsilon)}$ and $1/4$ -covering $\mathcal{N}_{1/4}$ of \mathbb{S}^{d-1} with $|\mathcal{N}_{1/4}| \leq e^{C_2 d}$ such that for all $\mathbf{w} \in \mathbb{S}^{d-1}$, there exists $\boldsymbol{\pi}(\mathbf{w}) \in \mathcal{N}_\epsilon$ such that $\|\mathbf{w} - \boldsymbol{\pi}(\mathbf{w})\| \leq \epsilon$ and $\mathbf{v}_t(\mathbf{w}) = \mathbf{v}_t(\boldsymbol{\pi}(\mathbf{w}))$ holds for all t . Then

$$\begin{aligned} & \sup_{\mathbf{w}} \|\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\| \\ &\leq \sup_{\mathbf{w} \in \mathcal{N}_\epsilon} \|\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\| + \sup_{\mathbf{w}} \|\tilde{z}_1(\mathbf{w}) - \tilde{z}_1(\boldsymbol{\pi}(\mathbf{w}))\| + \sup_{\mathbf{w}} \|\mathbb{E}[\tilde{z}_1(\mathbf{w})] - \mathbb{E}[\tilde{z}_1(\boldsymbol{\pi}(\mathbf{w}))]\| \\ &\leq \sup_{\mathbf{w} \in \mathcal{N}_\epsilon} \|\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\| + \sup_{\mathbf{w}} \|\tilde{z}_1(\mathbf{w}) - \tilde{z}_1(\boldsymbol{\pi}(\mathbf{w}))\| + O(d^{-C^*} + \epsilon d) \\ &\leq \sup_{\mathbf{w} \in \mathcal{N}_\epsilon} \|\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\| + O(RR'\epsilon\sqrt{d}) + O(d^{-C^*} + \epsilon d) \end{aligned}$$

Let us bound the first term in the right-hand side. Notice that

$$\sup_{\mathbf{w} \in \mathcal{N}_\epsilon} \|\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]\| \leq 2 \sup_{\mathbf{w} \in \mathcal{N}_\epsilon} \sup_{\mathbf{u} \in \mathcal{N}_{1/4}} \mathbf{u} \cdot [\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]].$$

Now, $\tilde{z}_1(\mathbf{w})$ is RR' -sub Gaussian and then with probability $1 - 2e^{-z}$, it holds that

$$\mathbf{u} \cdot [\tilde{z}_1(\mathbf{w}) - \mathbb{E}[\tilde{z}_1(\mathbf{w})]] \leq RR' \sqrt{2z/T},$$

where $R = \max_t R_t$ and $R' = \max_t R'_t$. Therefore, with probability at least $1 - 2e^{C_3 d \log(NT/\epsilon) - z}$ we have

$$\sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \|\mathbf{z}_1(\mathbf{w}) - \mathbb{E}[\mathbf{z}_1(\mathbf{w})]\| = O\left(RR' \sqrt{\frac{2z}{T}} + RR'\epsilon\sqrt{d} + \epsilon d + d^{-C^*}\right).$$

By taking $z = C_3 d \log(NT/\epsilon) + \iota$ and $\epsilon = O(d^{-C^*})$, we arrive at

$$\begin{aligned} \sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \|\mathbf{z}_1(\mathbf{w}) - \mathbb{E}[\mathbf{z}_1(\mathbf{w})]\| &= O\left(\iota^{P/2} \left(\sqrt{\frac{1}{d}} + \sqrt{\frac{d\iota^{P+1}}{N}} \right) \sqrt{\frac{2z}{T}} + \epsilon d + O(d^{-C^*})\right) \\ &= \tilde{O}\left(\left(\sqrt{\frac{1}{d}} + \sqrt{\frac{d}{N}}\right) \sqrt{\frac{d}{T}}\right) \end{aligned}$$

with high probability.

Bounds for $\sup_{\mathbf{w} \in \mathbb{S}^{d-1}} \|\mathbf{z}_2(\mathbf{w}) - \mathbb{E}[\mathbf{z}_2(\mathbf{w})]\|$ can be obtained by the same procedure, noting that $\mathbb{E}[y_{t,i} \sigma(\mathbf{w}^\top \mathbf{x}_{t,i})] = O(1/d)$. \square

Combining Lemmas 5 and 6, we immediately obtain the following:

Corollary 7. *Assume $T = \tilde{\Omega}(d^4)$ and $N = \tilde{\Omega}(d^2)$. Then, $\|\mathbf{g}_T(\mathbf{w})\| = \tilde{O}(\sqrt{\frac{r}{d^3}})$ and $\|\mathbf{g}_T(\mathbf{w})_{r+1:d}\| = \tilde{O}(\sqrt{\frac{r^2}{d^4}})$ holds with high probability.*

It is used to derive prompt length-free generalization result (Appendix D.2).

C. Construction of Attention Matrix

In this section, we construct $\mathbf{\Gamma}$ which satisfies the following approximation property:

Theorem 8. *With high probability (see Section A.1 for the definition), there exists $\mathbf{\Gamma}$ such that*

$$\left| \left\langle \frac{\mathbf{\Gamma} \sigma(\mathbf{W}^{(1)} \mathbf{X}^t + \mathbf{b}) \mathbf{y}}{\dim(\mathbf{y}^t)}, \begin{bmatrix} \sigma(\mathbf{w}_1^{(1)\top} \mathbf{x}^t + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^{(1)\top} \mathbf{x}^t + b_m) \end{bmatrix} \right\rangle - y^t \right| - \tau \lesssim \text{poly log}(d) \left(\frac{r^{5P/2}}{\sqrt{N_2}} + \frac{r^{3P/2}}{\sqrt{m}} \right)$$

for all $t \in \{T_1 + 1, \dots, T_2\}$, where $\{\mathbf{w}_j^{(1)}\}_j$ are neurons obtained by line 4 of Algorithm 1 with $\eta_1 = \Omega\left(\sqrt{\frac{d^3}{r}} \frac{1}{\text{poly log } d}\right)$, $m \gtrsim r^P$ and \mathbf{b} is re-initialized in Line 5 of Algorithm 1. Moreover, $\|\mathbf{\Gamma}\| = \tilde{O}(\sqrt{r^{5P}/m^2})$ is satisfied.

We give an intuitive explanation of the way to construct such $\mathbf{\Gamma}$. First, recall that the true function satisfies $f_*(\mathbf{x}) = \sum_{i=2}^P c_i \text{He}_i(\langle \mathbf{x}, \boldsymbol{\beta} \rangle)$. Let $\mathcal{F}_* = \{f_* : \sum_{i=2}^P c_i \text{He}_i(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) \mid \boldsymbol{\beta} \in \mathcal{S}, c_i \in \mathbb{R}\}$. We can find an orthogonal basis of \mathcal{F}_* : let $\{\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_r\}$ be an orthonormal basis of \mathcal{S} . Then, we can show that $\mathcal{G} = \{\prod_{j=1}^r \text{He}_{p_j}(\langle \boldsymbol{\beta}_j, \cdot \rangle) \mid 2 \leq p_1 + \dots + p_r \leq P, p_1 \geq 0, \dots, p_r \geq 0\}$ forms a basis for \mathcal{F}_* and $\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, I_d)}[g(\mathbf{x})g'(\mathbf{x})] = 0$ is satisfied if $g \neq g'$.

Let $q = |\mathcal{G}| (= O(r^P))$. Assign numbers to the functions in \mathcal{G} arbitrarily as $\mathcal{G} = \{g_1, \dots, g_q\}$. We can show that there exists $\mathbf{a}^1, \dots, \mathbf{a}^q \in \mathbb{R}^m$ such that $\sum_{j=1}^m \mathbf{a}_j^i \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j) \simeq g_i(\mathbf{x})$ for each $i \in [q]$. This comes from the approximation property of two-layer fully-connected neural networks.

A key fact is that we can construct attention matrix using this $\mathbf{a}^1, \dots, \mathbf{a}^q$. Let $\mathbf{A} = (\mathbf{a}^1 \ \dots \ \mathbf{a}^q) \in \mathbb{R}^{m \times q}$ and $\mathbf{D} = \text{diag}\{\alpha_1, \dots, \alpha_q\}$ where $\alpha_i = \mathbb{E}[g_i(\mathbf{x})^2]^{-1}$. Then, letting $\dim(\mathbf{y}) = N$ and $\mathbf{\Gamma} = \mathbf{A} \mathbf{D} \mathbf{A}^\top$, we notice that

$$\begin{aligned} & \left\langle \frac{\mathbf{\Gamma} \sigma(\mathbf{W} \mathbf{X} + \mathbf{b}) \mathbf{y}}{\dim(\mathbf{y})}, \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\rangle \\ &= \sum_{i=1}^q \alpha_i \left(\frac{1}{N} \sum_{j=1}^N \left(\sum_{k=1}^m \mathbf{a}_k^i \sigma(\mathbf{w}_k^\top \mathbf{x}_j + b_k) \right) y_j \right) \left(\sum_{k=1}^m \mathbf{a}_k^i \sigma(\mathbf{w}_k^\top \mathbf{x} + b_k) \right) \\ &\simeq \sum_{i=1}^q \alpha_i \mathbb{E}[g_i(\mathbf{x}) \sigma_*(\langle \boldsymbol{\beta}, \mathbf{x} \rangle)] g_i(\mathbf{x}) = \sigma_*(\langle \boldsymbol{\beta}, \mathbf{x} \rangle). \end{aligned} \tag{C.1}$$

660 The most essential point is that the self-attention architecture approximately calculates $\mathbb{E}[g_i(\mathbf{x})\sigma_*(\langle \boldsymbol{\beta}, \mathbf{x} \rangle)]$, i.e., the inner
 661 product between the true function and a base g_i . If N is sufficiently large, then we can approximate the true label well.
 662 Recall that we defined $\mathcal{G} = \{\prod_{j=1}^r \text{He}_{p_j}(\langle \boldsymbol{\beta}_j, \cdot \rangle) \mid 2 \leq p_1 + \dots + p_r \leq P, p_1 \geq 0, \dots, p_r \geq 0\}$ where $\{\boldsymbol{\beta}_j\}$ is a
 663 basis of \mathcal{S} and this forms a basis of the set of true functions. We let $|\mathcal{G}| = q$ and number the elements of \mathcal{G} arbitrarily as
 664 $\mathcal{G} = \{g_1, \dots, g_q\}$. Importantly, $q \lesssim r^P$ holds.

665 Note that, to derive (C.1), there are two types of ‘‘approximation error’’: one is the error by a two-layer neural network
 666 $\sum_{j=1}^m \mathbf{a}_j^i \sigma(\mathbf{w}_j^\top \mathbf{x} + b_j)$ to approximate each basis $g_i(\mathbf{x})$. The other is discrepancy between average $\frac{1}{N} \sum_{j=1}^N g_i(\mathbf{x}_j) y_j$ and
 667 its expectation $\mathbb{E}[g_i(\mathbf{x})\sigma_*(\langle \boldsymbol{\beta}, \mathbf{x} \rangle)]$. In Appendix C.1 we bound the former error, and the latter is evaluated in Appendix C.2.
 668 Finally, we prove Theorem 8 in Appendix C.3.

670 C.1. Approximation Result for Two-layer NN

671 We use the following proposition stating approximation ability of two-layer fully-connected neural networks:

672 **Proposition 9.** *Let $\{\mathbf{w}_j^{(1)}\}_j$ be neurons obtained by line 4 of Algorithm 1 with $m \gtrsim r^P$, $\eta_1 = \Omega\left(\sqrt{\frac{d^3}{r}} \frac{1}{\text{poly} \log d}\right)$ and
 673 assume \mathbf{b} is initialized as in Line 5 of Algorithm 1. Then, with high probability over the data distribution, there exists
 674 $\mathbf{a}^1, \dots, \mathbf{a}^q \in \mathbb{R}^m$ such that for each $i \in [q]$,*

$$675 \left(\sum_{j=1}^m \mathbf{a}_j^i \sigma(\mathbf{w}_j^{(1)\top} \mathbf{x} + b_j) - g_i(\mathbf{x}) \right)^2 = \tilde{O}\left(\frac{r^P}{m} + \frac{1}{N_2}\right)$$

676 holds for all $\mathbf{x} = \mathbf{x}_t$ ($t \in \{T_1 + 1, \dots, T_2\}$) and $\mathbf{x} = \mathbf{x}_{t,n}$ ($t \in \{T_1 + 1, \dots, T_2\}, n \in [N_2]$). Moreover, $\|\mathbf{a}^i\|^2 =$
 677 $\tilde{O}\left(\frac{r^P}{m}\right)$ ($i \in [q]$) holds.

678 This is almost a counterpart of Lemma 13 in (Damian et al., 2022) and the proof is almost the same. However, there are
 679 slight differences as the property of $\mathbf{w}_j^{(1)}$ differs: specifically, the statement of Lemmas 10 and 11 in (Damian et al., 2022)
 680 to derive Lemma 13 in the paper should be slightly changed. Here we describe the counterparts of these lemmas.

681 **Lemma 10** (counterpart of Lemma 10 in (Damian et al., 2022)). *Let $\mathbf{r}(w) = \mathbf{g}_T(w) - a_2^2 \mathbb{E}[c_2^2] \frac{12}{r(r+2)} \|\mathbf{w}_{1:r}\|_2^2 \mathbf{w}_{1:r}$. Then,*

$$682 \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} [\|\Pi_S \mathbf{r}(w)\|^j]^{1/j} \lesssim \tilde{O}\left(\left(\sqrt{\frac{1}{d}} + \sqrt{\frac{d}{N_1}}\right) \sqrt{\frac{d}{T_1}}\right) + O\left(\sqrt{\frac{r^2}{d^4}}\right)$$

683 holds with high probability, where Π_S be the orthogonal projection onto \mathcal{S} (recall that we assumed that \mathcal{S} be the subspace
 684 spanned by first r standard basis: that is, $\Pi_S \mathbf{v} = \mathbf{v}_{1:r}$).

685 **Proof.** Since $\mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} [\|\Pi_S [\mathbf{g}_T(w) - \mathbf{g}(w)]\|^j]^{1/j} = \tilde{O}\left(\left(\sqrt{\frac{1}{d}} + \sqrt{\frac{d}{N_1}}\right) \sqrt{\frac{d}{T_1}}\right)$ from Lemma 6, it suffices to upper
 686 bound

$$687 \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \Pi_S \left[\mathbf{g}(w) - a_2^2 \mathbb{E}[c_2^2] \frac{12}{r(r+2)} \|\mathbf{w}_{1:r}\|_2^2 \mathbf{w}_{1:r} \right] \right\|^j \right]^{1/j} = \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} [\|\mathbb{E}_y [\Pi_S \mathbf{s}(y, \mathbf{w})]\|^j]^{1/j}.$$

688 Then it is sufficient to show that

$$689 \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \mathbb{E}_y \left[\frac{a_{k+1} \mathbb{E}[\nabla^{k+1} f^t](\mathbf{w}^{\otimes k})}{k!} \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right] \right\|^j \right]^{1/j}$$

$$690 \lesssim \sqrt{\frac{r^2}{d^4}} \quad (k \geq 1, l \geq 2, (k, l) \neq (1, 2))$$

691 and

$$692 \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \mathbb{E}_y \left[\Pi_S \mathbf{w} \frac{a_{k+2} \mathbb{E}[\nabla^k f^t](\mathbf{w}^{\otimes k})}{k!} \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right] \right\|^j \right]^{1/j} \lesssim \sqrt{\frac{r^2}{d^4}} \quad (k \geq 2, l \geq 2).$$

First,

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \mathbb{E}_y \left[\frac{a_{k+1} \mathbb{E}[\nabla^{k+1} f^t](\mathbf{w}^{\otimes k})}{k!} \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right] \right\|^j \right]^{1/j} \\
 & \leq \mathbb{E}_y \left[\mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \frac{a_{k+1} \mathbb{E}[\nabla^{k+1} f^t](\mathbf{w}^{\otimes k})}{k!} \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right\|^j \right] \right]^{1/j} \\
 & \leq \mathbb{E}_y \left[\mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \frac{a_{k+1} \mathbb{E}[\nabla^{k+1} f^t](\mathbf{w}^{\otimes k})}{k!} \right\|^{2j} \right]^{1/2} \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right\|^{2j} \right]^{1/2} \right]^{1/j} \\
 & \lesssim \sqrt{\frac{r^{\lfloor k/2 \rfloor}}{d^k}} \sqrt{\frac{r^{\lfloor l/2 \rfloor}}{d^l}} \lesssim \sqrt{\frac{r^2}{d^4}}.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \mathbb{E}_y \left[\Pi_{\mathcal{S}} \mathbf{w} \frac{a_{k+2} \mathbb{E}[\nabla^k f^t](\mathbf{w}^{\otimes k})}{k!} \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right] \right\|^j \right]^{1/j} \\
 & \leq \mathbb{E}_y \left[\mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\|\Pi_{\mathcal{S}} \mathbf{w}\|^{2j} \right]^{1/2} \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \frac{a_{k+2} \mathbb{E}[\nabla^k f^t](\mathbf{w}^{\otimes k})}{k!} \right\|^{4j} \right]^{1/4} \right. \\
 & \quad \left. \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} \left[\left\| \frac{a_l \mathbb{E}[\nabla^l f^t](\mathbf{w}^{\otimes l})}{l!} \right\|^{4j} \right]^{1/4} \right]^{1/j} \\
 & \lesssim \sqrt{\frac{r}{d}} \sqrt{\frac{r^{\lfloor k/2 \rfloor}}{d^k}} \sqrt{\frac{r^{\lfloor l/2 \rfloor}}{d^l}} \lesssim \sqrt{\frac{r^2}{d^4}}.
 \end{aligned}$$

Then we obtain $\mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} [\|\Pi_{\mathcal{S}} \mathbf{r}(w)\|^j]^{1/j} \lesssim \tilde{O}\left(\left(\sqrt{\frac{1}{d}} + \sqrt{\frac{d}{N_1}}\right)\sqrt{\frac{d}{T_1}}\right) + O\left(\sqrt{\frac{r^2}{d^4}}\right)$ as desired. \square

To derive the counterpart of Lemma 11 in (Damian et al., 2022) we need the following statement:

Lemma 11 (Tensor expectation lower bound). *Let \mathbf{T} be a $k < P$ -symmetric tensor which has support on \mathcal{S} and let $\bar{\mathbf{w}} = \mathbf{w}_{1:r} \|\mathbf{w}_{1:r}\|^2$. Then*

$$\mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}} [\mathbf{T}(\bar{\mathbf{w}}^{\otimes k})^2] \gtrsim \frac{r^{2i}}{d^{3i}} \mathbb{E}[\|\mathbf{T}(\bar{\mathbf{w}}^{\otimes k-i})\|_F^2].$$

Proof. Let $\mathbf{u} \sim \mathcal{N}(0, I_d)$, $z \sim \chi(d)$ and $\bar{\mathbf{u}} = \mathbf{u}_{1:r} \|\mathbf{u}_{1:r}\|^2$. Then we can decompose as $\bar{\mathbf{u}} = z^3 \bar{\mathbf{w}}$. Therefore

$$\mathbb{E}[\mathbf{T}(\bar{\mathbf{u}}^{\otimes k})^2] = \mathbb{E}[z^{6k}] \mathbb{E}[\mathbf{T}(\bar{\mathbf{w}}^{\otimes k})^2]$$

holds.

On the other hand, let $\mathbf{x} \sim \mathbb{S}^{r-1}$ and $z' \sim \chi(r)$. Then we can decompose as $\bar{\mathbf{u}} = (z')^3 \mathbf{x}$. Thus

$$\mathbb{E}[\mathbf{T}(\bar{\mathbf{u}}^{\otimes k})^2] = \mathbb{E}[(z')^{6k}] \mathbb{E}[\mathbf{T}(\mathbf{x}^{\otimes k})^2]$$

holds. It implies that

$$\mathbb{E}[\mathbf{T}(\bar{\mathbf{w}}^{\otimes k})^2] = \frac{\mathbb{E}[(z')^{6k}]}{\mathbb{E}[z^{6k}]} \mathbb{E}[\mathbf{T}(\mathbf{x}^{\otimes k})^2].$$

Similarly,

$$\mathbb{E}[\|\mathbf{T}(\bar{\mathbf{w}}^{\otimes k-i})\|_F^2] = \frac{\mathbb{E}[(z')^{6(k-i)}]}{\mathbb{E}[z^{6(k-i)}]} \mathbb{E}[\|\mathbf{T}(\mathbf{x}^{\otimes(k-i)})\|_F^2]$$

is satisfied. Now from Corollary 13 in (Damian et al., 2022),

$$\mathbb{E}[\|\mathbf{T}(\mathbf{x}^{\otimes(k-i)})\|_F^2] \lesssim r^i \mathbb{E}[\mathbf{T}(\mathbf{x}^{\otimes k})^2]$$

770 holds. Therefore we obtain

$$\begin{aligned}
 771 \quad \mathbb{E}[\mathbf{T}(\bar{\mathbf{w}}^{\otimes k})^2] &= \frac{\mathbb{E}[(z')^{6k}]}{\mathbb{E}[z^{6k}]} \mathbb{E}[\mathbf{T}(\mathbf{x}^{\otimes k})^2] \\
 772 &\gtrsim r^{-i} \frac{\mathbb{E}[(z')^{6k}]}{\mathbb{E}[z^{6k}]} \mathbb{E}[\|\mathbf{T}(\mathbf{x}^{\otimes(k-i)})\|_F^2] \\
 773 &= r^{-i} \frac{\mathbb{E}[(z')^{6k}]}{\mathbb{E}[z^{6k}]} \frac{\mathbb{E}[z^{6(k-i)}]}{\mathbb{E}[(z')^{6(k-i)}]} \mathbb{E}[\|\mathbf{T}(\bar{\mathbf{w}}^{\otimes(k-i)})\|_F^2] \\
 774 &= \frac{r^{2i}}{d^{3i}} \mathbb{E}[\|\mathbf{T}(\bar{\mathbf{w}}^{\otimes(k-i)})\|_F^2].
 \end{aligned}$$

781 □

782
783
784 **Corollary 12.** Let \mathbf{T} be a $k < P$ -symmetric tensor which has support on \mathcal{S} and let $\hat{\mathbf{w}} = \frac{1}{r(r+2)} \mathbf{w}_{1:r} \|\mathbf{w}_{1:r}\|^2$. Then

$$785 \quad \mathbb{E}_{\mathbf{w} \sim \mathbb{S}^{d-1}}[\mathbf{T}(\hat{\mathbf{w}}^{\otimes k})^2] \gtrsim \frac{1}{d^{3i}} \mathbb{E}[\|\mathbf{T}(\hat{\mathbf{w}}^{\otimes(k-i)})\|_F^2].$$

786
787 Using Lemma 10 and Corollary 12 in the proof of Lemma 11 in (Damian et al., 2022) yields the following Lemma:

788
789 **Lemma 13** (counterpart of Lemma 11 in (Damian et al., 2022)). Suppose $r^2 \lesssim d$, $T_1 \gtrsim d^3$ and $N_1 \gtrsim d^2$. Let \mathbf{T} be

790 $k < P$ -symmetric tensor with $\|\mathbf{T}\|_F = 1$ and assume \mathbf{T} has support on \mathcal{S} . Then,

$$791 \quad \mathbb{E}_{\mathbf{w} \sim \text{Unif}\mathbb{S}^{d-1}}[(\Pi_{\mathcal{S}} \mathbf{g}_{\mathbf{T}}(\mathbf{w}))^{\otimes 2k}](\mathbf{T}, \mathbf{T}) \gtrsim d^{-3k}$$

792 holds.

793
794 Now we have Lemmas 10 and 13. By letting $\eta = \Omega\left(\sqrt{\frac{d^3}{r}} \frac{1}{\text{poly log } d}\right)$ and following the same line towards the proof of

795 Lemma 13 in (Damian et al., 2022), we obtain Proposition 9.

800 C.2. Concentration of Correlation between a Label and a Base

801 In this subsection, we give an upper bound for

$$802 \quad \left| \frac{1}{N} \sum_{j=1}^N y_j g(\mathbf{x}_j) - \mathbb{E}[y g(\mathbf{x})] \right|$$

803 as follows:

804
805 **Proposition 14.** Let $g \in \mathcal{G}$. With high probability,

$$806 \quad \left| \frac{1}{N} \sum_{j=1}^N y_j g(\mathbf{x}_j) - \mathbb{E}[y g(\mathbf{x})] \right| \lesssim \frac{r^{3P/2}}{\sqrt{N}} (\log d)^{P/2}$$

807 holds.

808 Without loss of generality, we assume that $\mathcal{S} = \{(\beta_1, \dots, \beta_r, 0, \dots, 0) \mid \beta_1, \dots, \beta_r \in \mathbb{R}\}$ and $g(\mathbf{x})$ can be written in the

809 form as $g(\mathbf{x}) = \text{He}_{q_1}(x_1) \cdots \text{He}_{q_r}(x_r)$, satisfying $q_1 + \dots + q_r \leq P$. Note that

$$810 \quad \frac{1}{N} \sum_{j=1}^N y_j g(\mathbf{x}_j) - \mathbb{E}[y g(\mathbf{x})] = \frac{1}{N} \sum_{j=1}^N \varsigma_j g(\mathbf{x}_j) + \frac{1}{N} \sum_{j=1}^N \sigma_*(\langle \mathbf{x}_j, \boldsymbol{\beta} \rangle) g(\mathbf{x}_j) - \mathbb{E}[\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x})],$$

811 where $\sigma_*(z) = \sum_{i=2}^P c_p \text{He}_i(z)$. First we give an upper bound for $\frac{1}{N} \sum_{j=1}^N \sigma_*(\langle \mathbf{x}_j, \boldsymbol{\beta} \rangle) g(\mathbf{x}_j) - \mathbb{E}[\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x})]$.

Lemma 15. Let $D = \deg \frac{1}{N} \sum_{j=1}^N \sigma_*(\langle \mathbf{x}_j, \boldsymbol{\beta} \rangle) g(\mathbf{x}_j) \leq 2P$. For all $t > 0$,

$$\begin{aligned} & \mathbb{P}\left(\left|\frac{1}{N} \sum_{j=1}^N \sigma_*(\langle \mathbf{x}_j, \boldsymbol{\beta} \rangle) g(\mathbf{x}_j) - \mathbb{E}[\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x})]\right| \geq t\right) \\ & \leq 2 \exp\left(-\frac{1}{C_D} \min_{1 \leq k \leq D} \left(\frac{t\sqrt{N}}{M^k \|\mathbb{E}[\nabla^k(\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F}\right)^{2/k}\right) \end{aligned}$$

holds, where C_D and M are absolute constants.

Proof. Let $F(\mathbf{x}_1, \dots, \mathbf{x}_N) = \frac{1}{N} \sum_{j=1}^N \sigma_*(\langle \mathbf{x}_j, \boldsymbol{\beta} \rangle) g(\mathbf{x}_j)$. It is a polynomial for $x_{11}, \dots, x_{1d}, \dots, x_{Nd}$, which are standard Gaussian variables. For the standard Gaussian $X \sim \mathcal{N}(0, 1)$, its Orlicz norm $\|X\|_{\Psi_2}$ is bounded; there exists M such that $\|X\|_{\Psi_2} \leq M$. Moreover,

$$\|\mathbb{E}[\nabla^k F(\mathbf{x}_1, \dots, \mathbf{x}_N)]\|_F = \sqrt{N}^{-1} \|\mathbb{E}[\nabla^k(\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F$$

is satisfied. Then, Theorem 1.2 in (Götze et al., 2021) yields the lemma. \square

Then, our goal is to bound $\|\mathbb{E}[\nabla^k(\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F$.

Lemma 16.

$$\|\mathbb{E}[\nabla^k(\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F = O(r^{3P})$$

holds.

Proof. First,

$$\begin{aligned} & \|\mathbb{E}[\nabla^k(\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F^2 \\ & = \left\| \sum_{p=2}^P c_p \mathbb{E}[\nabla^k(\text{He}_p(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\right\|_F^2 \\ & \leq \sum_{p=2}^P c_p \|\mathbb{E}[\nabla^k(\text{He}_p(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F^2 \\ & \leq (c_2^2 + \dots + c_P^2) \left(\sum_{p=2}^P \|\mathbb{E}[\nabla^k(\text{He}_p(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F^2 \right) \\ & \lesssim \sum_{p=2}^P \|\mathbb{E}[\nabla^k(\text{He}_p(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F^2 \\ & \leq \sum_{p=2}^P \left(\sum_{\substack{p_1 + \dots + p_r = p \\ p_1 \geq 0, \dots, p_r \geq 0}} \left(\frac{(p_1 + \dots + p_r)!}{p_1! \dots p_r!} \cdot \beta_1^{p_1} \dots \beta_r^{p_r} \right)^2 \right) \\ & \quad \cdot \sum_{\substack{p_1 + \dots + p_r = p \\ p_1 \geq 0, \dots, p_r \geq 0}} \|\mathbb{E}[\nabla^k(\text{He}_{p_1}(x_1) \dots \text{He}_{p_r}(x_r) g(\mathbf{x}))]\|_F^2 \end{aligned}$$

holds. Noting that $\left(\frac{(p_1 + \dots + p_r)!}{p_1! \dots p_r!} \cdot \beta_1^{p_1} \dots \beta_r^{p_r}\right)^2 \leq p!$ under $p_1 + \dots + p_r = p$ and the number of the combination of $(p_1, \dots, p_r) \geq 0$ such that $p_1 + \dots + p_r = p$ is at most r^p , there exists a constant C_P depending only on P such that

$$\|\mathbb{E}[\nabla^k(\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) g(\mathbf{x}))]\|_F^2 \leq C_P r^P \sum_{p=2}^P \sum_{\substack{p_1 + \dots + p_r = p \\ p_1 \geq 0, \dots, p_r \geq 0}} \|\mathbb{E}[\nabla^k(\text{He}_{p_1}(x_1) \dots \text{He}_{p_r}(x_r) g(\mathbf{x}))]\|_F^2. \quad (\text{C.2})$$

Next, let us bound $\|\mathbb{E}[\nabla^k(\text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r)g(\mathbf{x}))]\|_F^2$, where $g(\mathbf{x}) = \text{He}_{q_1}(x_1) \cdots \text{He}_{q_r}(x_r)$. It holds that

$$\nabla^k(\text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r)g(\mathbf{x})) = \sum_{i=0}^k \binom{k}{i} \nabla^i \text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r) \nabla^{k-i} \text{He}_{q_1}(x_1) \cdots \text{He}_{q_r}(x_r). \quad (\text{C.3})$$

Let $s = |p_1 - q_1| + \cdots + |p_r - q_r|$ and $t = \sum_{i=1}^r |p_i - q_i| \mathbf{1}_{p_i > q_i}$. Let us first consider the case where $k = s$ to simplify the explanation; in this case, each element of $\mathbb{E}[\nabla^k(\text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r)g(\mathbf{x}))]$ can be decomposed into the terms as

$$\mathbb{E}[\partial_{x_1}^{u_1} \cdots \partial_{x_r}^{u_r} \text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r) \partial_{x_1}^{v_1} \cdots \partial_{x_r}^{v_r} \text{He}_{q_1}(x_1) \cdots \text{He}_{q_r}(x_r)], \quad (\text{C.4})$$

where $u_1 + \cdots + u_r + v_1 + \cdots + v_r = s$. However, this expectation becomes nonzero only when $u_i = |p_i - q_i| \mathbf{1}_{p_i > q_i}$ and $v_i = |p_i - q_i| \mathbf{1}_{p_i < q_i}$. Therefore, by seeing (C.3), we can notice that

$$\mathbb{E}[\nabla^i \text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r) \nabla^{k-i} \text{He}_{q_1}(x_1) \cdots \text{He}_{q_r}(x_r)] \quad (\text{C.5})$$

is nonzero tensor only when $i = t$ and $\mathbb{E}[\nabla^t \text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r) \nabla^{k-t} \text{He}_{q_1}(x_1) \cdots \text{He}_{q_r}(x_r)]$ has only one nonzero element. It implies that $\|\mathbb{E}[\nabla^s(\text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r)g(\mathbf{x}))]\|_F^2 \lesssim \binom{s}{t}^2$.

Let us consider the other cases: obviously, $\|\mathbb{E}[\nabla^k(\text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r)g(\mathbf{x}))]\|_F^2 = 0$ if $k < s$. Moreover, if $k - s$ is odd, one can confirm that $\|\mathbb{E}[\nabla^k(\text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r)g(\mathbf{x}))]\|_F^2 = 0$. Consider the case where $k - s = 2l > 0$. The expectation (C.4) under $u_1 + \cdots + u_r + v_1 + \cdots + v_r = l = s + 2l$ is nonzero when $u_i = |p_i - q_i| \mathbf{1}_{p_i > q_i} + l_i$ and $v_i = |p_i - q_i| \mathbf{1}_{p_i < q_i} + l_i$, where $l_1 + \cdots + l_r = l$. It implies that (C.5) is nonzero only when $i = t + l$, and in this case, the tensor (C.5) has at most r^l nonzero entries. As a consequence,

$$\|\mathbb{E}[\nabla^s(\text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r)g(\mathbf{x}))]\|_F^2 \lesssim r^l \binom{s+l}{t+l}^2$$

holds. Overall, $\|\mathbb{E}[\nabla^k(\text{He}_{p_1}(x_1) \cdots \text{He}_{p_r}(x_r)g(\mathbf{x}))]\|_F^2 \leq C'_P r^P$ is satisfied where C'_P depends only on P . Plugging this bound into (C.2), we arrive at

$$\|\mathbb{E}[\nabla^k(\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle)g(\mathbf{x}))]\|_F^2 \leq C_P C'_P r^P \sum_{p=2}^P \sum_{p_1 \geq 0, \dots, p_r \geq 0}^{p_1 + \dots + p_r = p} r^P \lesssim r^{3P},$$

as desired. \square

Plugging the result above and $t = \Theta(\frac{r^{3P/2}}{\sqrt{N}}(\log d)^{P/2})$ into Lemma 15 yields the following corollary.

Corollary 17. *With high probability,*

$$\left| \frac{1}{N} \sum_{j=1}^N \sigma_*(\langle \mathbf{x}_j, \boldsymbol{\beta} \rangle)g(\mathbf{x}_j) - \mathbb{E}[\sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle)g(\mathbf{x})] \right| \lesssim \frac{r^{3P/2}}{\sqrt{N}}(\log d)^{P/2}$$

holds.

Proof. [Proof of Proposition 14] Now we have Corollary 17, then it remains to show that $\frac{1}{N} \sum_{j=1}^N \varsigma_j g(\mathbf{x}_j) \lesssim \frac{r^{3P/2}}{\sqrt{N}}(\log d)^{P/2}$ with high probability. This is obvious from Lemma 18. \square

C.3. Proof of Theorem 8

We prove Theorem 8 using the preparations above. It suffices to show the theorem in the case where $T = 1$, as long as T is polynomial in d . Then, we drop the subscript t .

Note that, from Lemma 3, we can expand $f_*(\mathbf{x})$ as $f_*(\mathbf{x}) = \sum_{i=1}^q \alpha_i \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, I_d)}[y g_i(\mathbf{x})] g_i(\mathbf{x})$, where $\alpha_i = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, I_d)}[g_i(\mathbf{x})^2]^{-1}$. Note that $\alpha_i = O_{d,r}(1)$ and $\alpha_i^{-1} = O_{d,r}(1)$. Now let $\mathbf{A} = (\mathbf{a}^1 \cdots \mathbf{a}^q) \in \mathbb{R}^{m \times q}$ and

$D = \text{diag}\{\alpha_1, \dots, \alpha_q\}$. First,

$$\begin{aligned}
 & \left\langle \frac{ADA^\top \sigma(WX + b)y}{\dim(y)}, \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\rangle - y \\
 &= \sum_{i=1}^q \alpha_i (g_i(\mathbf{x}) + \epsilon_i(\mathbf{x})) \left(\mathbb{E}[yg_i(\mathbf{x})] + \left(\frac{1}{N_2} \sum_{j=1}^{N_2} y_j g_i(\mathbf{x}_j) - \mathbb{E}[yg_i(\mathbf{x})] \right) + \frac{1}{N_2} \sum_{j=1}^{N_2} y_j \epsilon_i(\mathbf{x}_j) \right) \\
 &\quad - \sigma_*(\langle \mathbf{x}, \boldsymbol{\beta} \rangle) - \varsigma \\
 &= \sum_{i=1}^q \alpha_i (g_i(\mathbf{x}) + \epsilon_i(\mathbf{x})) \left(\mathbb{E}[yg_i(\mathbf{x})] + \left(\frac{1}{N_2} \sum_{j=1}^{N_2} y_j g_i(\mathbf{x}_j) - \mathbb{E}[yg_i(\mathbf{x})] \right) + \frac{1}{N_2} \sum_{j=1}^{N_2} y_j \epsilon_i(\mathbf{x}_j) \right) \\
 &\quad - \sum_{i=1}^q g_i(\mathbf{x}) \alpha_i \mathbb{E}[yg_i(\mathbf{x})] - \varsigma
 \end{aligned}$$

holds, where $\epsilon_i(\mathbf{x}) = \sum_{j=1}^m \mathbf{a}_j^i \sigma(\mathbf{w}_j^{(1)\top} \mathbf{x} + b_j) - g_i(\mathbf{x})$.

Then, it follows that

$$\begin{aligned}
 & \left| \left\langle \frac{ADA^\top \sigma(WX + b)y}{\dim(y)}, \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\rangle - y \right| \\
 & \leq \sum_{i=1}^q \alpha_i |\epsilon_i(\mathbf{x}) \mathbb{E}[yg_i(\mathbf{x})]| + \sum_{i=1}^q \alpha_i \left| g_i(\mathbf{x}) \left(\frac{1}{N_2} \sum_{j=1}^{N_2} y_j g_i(\mathbf{x}_j) - \mathbb{E}[yg_i(\mathbf{x})] \right) + g_i(\mathbf{x}) \frac{1}{N_2} \sum_{j=1}^{N_2} y_j \epsilon_i(\mathbf{x}_j) \right| \\
 & \quad + \sum_{i=1}^q \alpha_i \left| \epsilon_i(\mathbf{x}) \left(\frac{1}{N_2} \sum_{j=1}^{N_2} y_j g_i(\mathbf{x}_j) - \mathbb{E}[yg_i(\mathbf{x})] \right) + \epsilon_i(\mathbf{x}) \frac{1}{N_2} \sum_{j=1}^{N_2} y_j \epsilon_i(\mathbf{x}_j) \right| + \tau.
 \end{aligned}$$

To bound these terms, we need to bound $g_i(\mathbf{x})$ and $\frac{1}{N_2} \sum_{j=1}^{N_2} y_j \epsilon_i(\mathbf{x}_j)$, whose expectations are zero.

Lemma 18. *Let $s = \deg g_i$. Then*

$$|g_i(\mathbf{x})| \lesssim (\log d)^{s/2}$$

with high probability for each $i \in [q]$.

Proof. Note that $\mathbb{E}[\nabla^d g_i(\mathbf{x})]$ has only one nonzero element only when $d = s$. Then, from Theorem 1.2 in (Götze et al., 2021),

$$\mathbb{P}(|g_i(\mathbf{x})| \geq t) \leq 2 \exp\left(-\frac{1}{C} \left(\frac{t}{M^s}\right)^{2/s}\right)$$

holds. Plugging $t = \Omega((\log d)^{s/2})$ yields the result. \square

Lemma 19.

$$\left| \frac{1}{N_2} \sum_{j=1}^{N_2} y_j \epsilon_i(\mathbf{x}_j) \right| \lesssim \frac{r^{P/2} (\log d)^{P/2}}{\sqrt{m}} + \frac{(\log d)^{P/2}}{\sqrt{N_2}}$$

with high probability.

Proof. From Lemma 17 in (Damian et al., 2022), $\sigma_*(\langle \boldsymbol{\beta}, \mathbf{x}_j \rangle) \lesssim (\log d)^{P/2}$ holds with high probability. Then, the lemma follows immediately from Proposition 9. \square

In addition to the lemmas above, we know that $\mathbb{E}[yg_i(\mathbf{x})] = O(1)$ from Lemma 3 and $q = O(r^P)$. Note that as we assumed $m \gtrsim r^P$, $\epsilon_i(\mathbf{x}) = O(\text{poly log}(d))$ holds. Then, we arrive at

$$\begin{aligned}
 & \left| \left\langle \frac{\mathbf{A} \mathbf{D} \mathbf{A}^\top \sigma(\mathbf{W} \mathbf{X} + \mathbf{b}) \mathbf{y}}{\dim(\mathbf{y})}, \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\rangle - y \right| \\
 & \leq \sum_{i=1}^q |\epsilon_i(\mathbf{x}) \mathbb{E}[yg_i(\mathbf{x})]| + \sum_{i=1}^q \left| g_i(\mathbf{x}) \left(\frac{1}{N_2} \sum_{j=1}^{N_2} y_j g_i(\mathbf{x}_j) - \mathbb{E}[yg_i(\mathbf{x})] \right) + g_i(\mathbf{x}) \frac{1}{N_2} \sum_{j=1}^{N_2} y_j \epsilon_i(\mathbf{x}_j) \right| \\
 & + \sum_{i=1}^q \left| \epsilon_i(\mathbf{x}) \left(\frac{1}{N_2} \sum_{j=1}^{N_2} y_j g_i(\mathbf{x}_j) - \mathbb{E}[yg_i(\mathbf{x})] \right) + \epsilon_i(\mathbf{x}) \frac{1}{N_2} \sum_{j=1}^{N_2} y_j \epsilon_i(\mathbf{x}_j) \right| + \tau \\
 & \lesssim r^P \left(\frac{r^{P/2}}{\sqrt{m}} + \frac{1}{\sqrt{N_2}} \right) + \text{poly log}(d) \left(r^P \cdot \frac{r^{3P/2}}{\sqrt{N_2}} + r^P \cdot \left(\frac{r^{P/2}}{\sqrt{m}} + \frac{1}{\sqrt{N_2}} \right) \right) + \tau \\
 & \lesssim \text{poly log}(d) \left(\frac{r^{5P/2}}{\sqrt{N_2}} + \frac{r^{3P/2}}{\sqrt{m}} \right) + \tau.
 \end{aligned}$$

For $\|\mathbf{\Gamma}\|$, we obtain

$$\begin{aligned}
 \|\mathbf{\Gamma}\| & \leq \|\mathbf{A}\| \|\mathbf{D}\| \|\mathbf{A}\| \\
 & \leq \sum_{i=1}^q \|\mathbf{a}^i\|^2 \cdot \sqrt{q} \\
 & = \tilde{O} \left(\sqrt{r^{5P}/m^2} \right).
 \end{aligned}$$

D. Generalization Error Analysis and Proof of Theorem 1

D.1. Rademacher Complexity Bound

Let

$$\begin{aligned}
 \mathcal{F}_{N,G,W,B} & = \left\{ (\mathbf{X}, \mathbf{y}, \mathbf{x}) \mapsto \left\langle \frac{\mathbf{\Gamma} \sigma(\mathbf{W} \mathbf{X} + \mathbf{b}) \mathbf{y}}{N}, \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\rangle \right\} \\
 & \quad \left\{ \|\mathbf{\Gamma}\|_F \leq G, \|\mathbf{w}_1\|_2, \dots, \|\mathbf{w}_m\|_2 \leq W, |b_1| \dots |b_m| \leq B \right\}
 \end{aligned}$$

be the set of transformers whose parameter norms are constrained, and let

$$\text{Rad}_T(\mathcal{F}_{N,G,W,B}) = \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, \epsilon} \left[\sup_{f \in \mathcal{F}} \frac{1}{T} \sum_{t=1}^T \epsilon^t f(\mathbf{X}^t, \mathbf{y}^t, \mathbf{x}^t) \right]$$

be its Rademacher complexity, where $\epsilon_i \sim \text{Unif}(\{\pm 1\})$.

We evaluate the Rademacher complexity as follows:

Proposition 20.

$$\text{Rad}_T(\mathcal{F}_{N,G,W,B}) = O \left(\text{polylog}(T) \frac{mG}{\sqrt{T}} \left(r^{P/2} \sqrt{d} + d \sqrt{d} \right) \right).$$

holds, when $B = O_{d,r}(1)$ and $W = O_{d,r}(1)$.

1045 **Proof.**

1046 First,

$$\begin{aligned}
 & \text{Rad}_T(\mathcal{F}_{N,G,W,B}) \\
 &= \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, \epsilon} \left[\sup_{f \in \mathcal{F}} \frac{1}{T} \sum_{t=1}^T \epsilon^t f(\mathbf{X}^t, \mathbf{y}^t, \mathbf{x}^t) \right] \\
 &= \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, \epsilon} \left[\sup_{\Gamma, \mathbf{W}, \mathbf{b}} \frac{1}{T} \sum_{t=1}^T \epsilon^t \left\langle \frac{\Gamma \sigma(\mathbf{W} \mathbf{X}^t + \mathbf{b}) \mathbf{y}^t}{N}, \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x}^t + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x}^t + b_m) \end{bmatrix} \right\rangle \right] \\
 &\leq \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, \epsilon} \left[\sup_{\Gamma, \mathbf{W}, \mathbf{b}} \frac{1}{T} \|\Gamma\|_F \left\| \sum_{t=1}^T \epsilon^t \frac{\sigma(\mathbf{W} \mathbf{X}^t + \mathbf{b}) \mathbf{y}^t}{N} \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x}^t + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x}^t + b_m) \end{bmatrix} \right\|_F \right] \\
 &\leq \frac{Gm}{T} \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, \epsilon} \left[\sup_{\mathbf{W}, \mathbf{b}} \left\| \sum_{t=1}^T \epsilon^t \frac{\sigma(\mathbf{W} \mathbf{X}^t + \mathbf{b}) \mathbf{y}^t}{N} \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x}^t + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x}^t + b_m) \end{bmatrix} \right\|_\infty \right] \\
 &\leq \frac{Gm}{TN} \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, \epsilon} \left[\sup_{\mathbf{w}, \mathbf{b}, \mathbf{w}', \mathbf{b}'} \left| \sum_{t=1}^T \epsilon^t \sum_{i=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_i^t + b) y_i^t \sigma(\mathbf{w}'^\top \mathbf{x}^t + b') \right| \right] \\
 &\leq \frac{Gm}{TN} \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, \epsilon} \left[\sup_{\mathbf{w}, \mathbf{b}, \mathbf{w}', \mathbf{b}'} \sum_{t=1}^T \epsilon^t \sum_{i=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_i^t + b) y_i^t \sigma(\mathbf{w}'^\top \mathbf{x}^t + b') \right. \\
 &\quad \left. + \sup_{\mathbf{w}, \mathbf{b}, \mathbf{w}', \mathbf{b}'} - \sum_{t=1}^T \epsilon^t \sum_{i=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_i^t + b) y_i^t \sigma(\mathbf{w}'^\top \mathbf{x}^t + b') \right] \\
 &\leq \frac{2Gm}{TN} \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, \epsilon} \left[\sup_{\mathbf{w}, \mathbf{b}, \mathbf{w}', \mathbf{b}'} \sum_{t=1}^T \epsilon^t \sum_{i=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_i^t + b) y_i^t \sigma(\mathbf{w}'^\top \mathbf{x}^t + b') \right] \\
 &\leq \frac{2Gm}{T} \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_{\mathbf{w}, \mathbf{b}, \mathbf{w}', \mathbf{b}'} \sum_{t=1}^T \epsilon^t \sigma(\mathbf{w}^\top \mathbf{x}^t + b) y^t \sigma(\mathbf{w}'^\top \mathbf{x}'^t + b') \right] \\
 &\leq \frac{2Gm}{T} \mathbb{E}_{R, R'} \left[\mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_{\mathbf{w}, \mathbf{b}, \mathbf{w}', \mathbf{b}'} \sum_{t=1}^T \epsilon^t \sigma(\mathbf{w}^\top \mathbf{x}^t + b) y^t \sigma(\mathbf{w}'^\top \mathbf{x}'^t + b') \right] \Big|_{A_{R, R'}} \right],
 \end{aligned}$$

1083 where

$$A_{R, R'} := \{\max\{\|\mathbf{x}_{1:r}^t\|^2, \|\mathbf{x}'_{1:r}{}^t\|^2\}_{t=1}^T = R^2, \max\{\|\mathbf{x}_{r+1:d}^t\|^2, \|\mathbf{x}'_{r+1:d}{}^t\|^2\}_{t=1}^T = R'^2\}.$$

1084 We utilize multivariate contraction inequality (Maurer, 2016) to bound the last line.

1085 **Lemma 21.** *Let $f(x, y, z) = \sigma(x)y\sigma(z)$, whose domain is restricted to $|x| \leq R_1, |y| \leq R_2, |z| \leq R_3$. Then, f is*

1086 $\sqrt{R_1^2 R_2^2 + R_2^2 R_3^2 + R_3^2 R_1^2} \leq R_1 R_2 + R_2 R_3 + R_3 R_1$ -Lipschitz continuous.

1087 **Lemma 22.** *There exists a polynomial g such that*

$$\begin{aligned}
 & \frac{2Gm}{T} \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_{\mathbf{w}, \mathbf{b}, \mathbf{w}', \mathbf{b}'} \sum_{t=1}^T \epsilon^t \sigma(\mathbf{w}^\top \mathbf{x}^t + b) y^t \sigma(\mathbf{w}'^\top \mathbf{x}'^t + b') \Big|_{A_{R, R'}} \right] \\
 & \leq \frac{2\sqrt{2}GmL}{T} \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_{\mathbf{w}, \mathbf{b}, \mathbf{w}', \mathbf{b}'} \sum_{t=1}^T \left(\epsilon^{t,1} (\mathbf{w}^\top \mathbf{x}^t + b) + \epsilon^{t,2} y^t + \epsilon^{t,3} (\mathbf{w}'^\top \mathbf{x}'^t + b') \right) \Big|_{A_{R, R'}} \right],
 \end{aligned}$$

1088 where $L = 2(W(R + R') + B)(g(R) + \tau) + (W(R + R') + B)^2$. Moreover, $g(z)$ is at most of degree P , increasing in

1089 the region $z \geq 0$, and its coefficient is $O(\sqrt{\sum_i c_i^2})$.

Proof. First, as $\|\mathbf{x}^t\| \leq \sqrt{R^2 + (R')^2}$, $|\mathbf{w}^\top \mathbf{x}^t + b| \leq W\sqrt{R^2 + (R')^2} + B \leq W(R + R') + B$ holds. Secondly,

$$\begin{aligned} |y^t| &\leq \left| \sum_{i=2}^P c_i \text{He}_i(\boldsymbol{\beta}^\top \mathbf{x}_{1:r}^t) \right| + \tau \\ &\leq \sqrt{\sum_{i=2}^P c_i^2} \sqrt{\sum_{i=2}^P \text{He}_i(\boldsymbol{\beta}^\top \mathbf{x}_{1:r}^t)^2} + \tau \end{aligned}$$

holds. Let $\text{He}_i(z) = \sum_{j=0}^P h_{ij} z^j$. Define $f(z)$ as $f(z) = \sum_{j=0}^P c_j z^j$ where $c_j = \max_{i=2}^P |h_{ij}|$. Then, one can show that $\text{He}_i(z)^2 \leq f(|z|)^2$ for all i and z . Moreover, $f(z)$ is increasing in the region $z \geq 0$. Using this f , we can obtain

$$\begin{aligned} |y^t| &\leq \sqrt{\sum_{i=2}^P c_i^2} \sqrt{\sum_{i=2}^P f(R)^2} + \tau \\ &= \sqrt{\sum_{i=2}^P c_i^2} \sqrt{P-1} f(R) + \tau \end{aligned}$$

Then, the lemma immediately follows from vector-value contraction inequality (Maurer, 2016), by letting $g(z) = \sqrt{\sum_{i=2}^P c_i^2} \sqrt{P-1} f(z)$. \square

Moreover, we can observe that

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_{\mathbf{w}, b, \mathbf{w}', b'} \sum_{t=1}^T \left(\epsilon^{t,1} (\mathbf{w}^\top \mathbf{x}^t + b) + \epsilon^{t,2} y^t + \epsilon^{t,3} (\mathbf{w}'^\top \mathbf{x}'^t + b') \right) \middle| A_{R,R'} \right] \\ &\leq \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_{\mathbf{w}, b, \mathbf{w}', b'} \sum_{t=1}^T \epsilon^{t,1} (\mathbf{w}^\top \mathbf{x}^t + b) + \sum_{t=1}^T \epsilon^{t,2} y^t + \sup_{\mathbf{w}, b, \mathbf{w}', b'} \sum_{t=1}^T \epsilon^{t,3} (\mathbf{w}'^\top \mathbf{x}'^t + b') \middle| A_{R,R'} \right] \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_{\mathbf{w}, b, \mathbf{w}', b'} \sum_{t=1}^T \epsilon^{t,1} (\mathbf{w}^\top \mathbf{x}^t + b) + \sup_{\mathbf{w}, b, \mathbf{w}', b'} \sum_{t=1}^T \epsilon^{t,3} (\mathbf{w}'^\top \mathbf{x}'^t + b') \middle| A_{R,R'} \right] \\ &= 2\mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_{\mathbf{w}} \mathbf{w}^\top \sum_{t=1}^T \epsilon^t \mathbf{x}^t \middle| A_{R,R'} \right] + 2\mathbb{E}_{\mathbf{x}, \mathbf{y}, \mathbf{x}', \epsilon} \left[\sup_b \sum_{t=1}^T \epsilon^t b \middle| A_{R,R'} \right] \\ &\leq 2W\mathbb{E}_{\mathbf{x}} \left[\sqrt{\left\| \sum_{t=1}^T \epsilon^t \mathbf{x}^t \right\|^2} \middle| A_{R,R'} \right] + 2B\mathbb{E}_{\epsilon} \left[\sqrt{\left(\sum_{t=1}^T \epsilon^t \right)^2} \right] \\ &\leq 2W\mathbb{E}_{\mathbf{x}} \sqrt{\left\| \sum_{t=1}^T \epsilon^t \mathbf{x}^t \right\|^2 \middle| A_{R,R'}} + 2B\sqrt{\mathbb{E}_{\epsilon} \left[\left(\sum_{t=1}^T \epsilon^t \right)^2 \right]} \\ &\leq 2W\sqrt{T}(R + R') + 2B\sqrt{T}. \end{aligned}$$

Finally we should evaluate

$$\begin{aligned} &\mathbb{E}_{R,R'} \left[\frac{2\sqrt{2}GmL}{T} \left(2W\sqrt{T}(R + R') + 2B\sqrt{T} \right) \right] \\ &= \frac{mGC'_{W,B}}{\sqrt{T}} \mathbb{E}_{R,R'} \left[2(R + R')(W(R + R') + B)(g(R) + \tau) + (R + R')(W(R + R') + B)^2 \right] \\ &\quad + \frac{mGC'_{W,B}}{\sqrt{T}} \mathbb{E}_{R,R'} \left[2(W(R + R') + B)(g(R) + \tau) + (W(R + R') + B)^2 \right] \end{aligned}$$

$$= \frac{mG}{\sqrt{T}} O(\mathbb{E}_R[R^{P+1}] + \mathbb{E}_{R'}[R']\mathbb{E}_R[R^P] + \mathbb{E}_{R,R'}[(R + R')^3]),$$

where $C'_{W,B}$ is a constant which only depends on W and B . in order to do so, it suffices to evaluate $\mathbb{E}_R[R^k]$ and where $\mathbb{E}_{R'}[(R')^k]$.

Lemma 23.

$$\mathbb{E}_{z_1, \dots, z_{2T} \sim \chi(r)} \left[\max_{t \in [2T]} z_t^k \right] = O_{r,T} \left((r + \sqrt{8r})^{k/2} (1 + \log T)^{k/4} \right)$$

holds. Note that we regard k as $O_{r,T}(1)$.

Proof. From the concentration inequality for chi-squared distribution (Wainwright, 2019)[Example 2.11],

$$z_t^2 \leq r + \sqrt{8r \log(1/\delta)} \quad (\text{D.1})$$

holds with probability at least $1 - \delta$ for each t . Then, equation (D.1) holds uniformly over all t with probability at least $1 - 2T\delta$. Let $A_i = (r + \sqrt{8r \log(2T \cdot 2^i)})^{k/2}$. Then, with probability at least $1 - 2^{-i}$, $\max_{t \in [2T]} z_t^k \leq A_i$ holds. Therefore we can divide the certain event into events with probability $1/2, 1/4, \dots$, where $\max_{t \in [2T]} z_t^k \leq A_1, \max_{t \in [2T]} z_t^k \leq A_2, \dots$ holds. It implies that

$$\begin{aligned} \mathbb{E}_{z_1, \dots, z_{2T} \sim \chi(r)} \left[\max_{t \in [2T]} z_t^k \right] &\leq \sum_{i=1}^{\infty} 2^{-i} (r + \sqrt{8r \log(2T \cdot 2^i)})^{k/2} \\ &\leq 2 \int_0^{1/2} \left(r + \sqrt{8r \log \frac{2T}{t}} \right)^{k/2} dt \\ &\leq (r + \sqrt{8r})^{k/2} \int_0^{1/2} \left(\sqrt{\log \frac{2T}{t}} \right)^{k/2} dt \\ &\leq (r + \sqrt{8r})^{k/2} (1 + \log T)^{k/4} \int_0^{1/2} \left(\sqrt{\log \frac{2}{t}} \right)^{k/2} dt. \end{aligned}$$

□

From Lemma 23, we arrive at

$$\text{Rad}_T(\mathcal{F}_{N,G,W,B}) = O\left(\text{polylog}(T) \frac{mG}{\sqrt{T}} \left(r^{P/2} \sqrt{d} + d\sqrt{d} \right)\right).$$

□

D.2. Prompt Length-free Generalization Bound

Let the ICL risk for prompt length N be

$$\mathcal{R}_N(\mathbf{\Gamma}, \mathbf{W}, \mathbf{b}) = \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, y} [|f(\mathbf{X}_{1:N}, \mathbf{y}_{1:N}, \mathbf{x}; \mathbf{W}, \mathbf{\Gamma}, \mathbf{b}) - y|],$$

where the length of $\mathbf{X}_{1:N}$ and $\mathbf{y}_{1:N}$ is fixed to N . In this section, we upper bound $|\mathcal{R}_N(\mathbf{\Gamma}, \mathbf{W}, \mathbf{b}) - \mathcal{R}_M(\mathbf{\Gamma}, \mathbf{W}, \mathbf{b})|$ under the condition $N, M \gtrsim r^{\Theta(P)}$.

Proposition 24. Assume that $\|\mathbf{w}_{j,1:r}\| = O(1), \|\mathbf{w}_{j,1:r}\| = O(\sqrt{r/d})$ and $|b_j| = O(1)$ for each $j \in [m]$. Then,

$$|\mathcal{R}_N(\mathbf{\Gamma}, \mathbf{W}, \mathbf{b}) - \mathcal{R}_M(\mathbf{\Gamma}, \mathbf{w}, \mathbf{b})| = \tilde{O}\left(\|\mathbf{\Gamma}\|_F \sqrt{r^2 m^2 / N + r^2 m^2 / M}\right)$$

holds.

1210 **Proof.** Note that

$$\begin{aligned}
 & 1211 \quad |\mathcal{R}_N(\Gamma, \mathbf{W}, \mathbf{b}) - \mathcal{R}_M(\Gamma, \mathbf{w}, \mathbf{b})| \\
 & 1212 \quad \leq \mathbb{E}[|f(\mathbf{X}_{1:N}, \mathbf{y}_{1:N}, \mathbf{x}; \mathbf{W}, \Gamma, \mathbf{b}) - f(\mathbf{X}_{1:M}, \mathbf{y}_{1:M}, \mathbf{x}; \mathbf{W}, \Gamma, \mathbf{b})|] \\
 & 1213 \quad = \mathbb{E} \left[\left\| \left\langle \Gamma \left(\frac{\sigma(\mathbf{W}\mathbf{X}_{1:N} + \mathbf{b})\mathbf{y}_{1:N}}{N} - \frac{\sigma(\mathbf{W}\mathbf{X}_{1:M} + \mathbf{b})\mathbf{y}_{1:M}}{M} \right), \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\rangle \right\| \right] \\
 & 1214 \quad \\
 & 1215 \quad \leq \|\Gamma\|_F \mathbb{E} \left[\left\| \left(\frac{\sigma(\mathbf{W}\mathbf{X}_{1:N} + \mathbf{b})\mathbf{y}_{1:N}}{N} - \frac{\sigma(\mathbf{W}\mathbf{X}_{1:M} + \mathbf{b})\mathbf{y}_{1:M}}{M} \right) \right\| \left\| \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\| \right] \\
 & 1216 \quad \\
 & 1217 \quad \\
 & 1218 \quad \\
 & 1219 \quad \\
 & 1220 \quad \\
 & 1221 \quad \\
 & 1222 \quad \\
 & 1223 \quad \leq \|\Gamma\|_F \mathbb{E} \left[\left\| \left(\frac{\sigma(\mathbf{W}\mathbf{X}_{1:N} + \mathbf{b})\mathbf{y}_{1:N}}{N} - \frac{\sigma(\mathbf{W}\mathbf{X}_{1:M} + \mathbf{b})\mathbf{y}_{1:M}}{M} \right) \right\|^2 \right]^{1/2} \mathbb{E} \left[\left\| \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\|^2 \right]^{1/2} \\
 & 1224 \quad \\
 & 1225 \quad \\
 & 1226 \quad \\
 & 1227 \quad
 \end{aligned}$$

1228 First let us bound $\mathbb{E} \left[\left\| \left(\frac{\sigma(\mathbf{W}\mathbf{X}_{1:N} + \mathbf{b})\mathbf{y}_{1:N}}{N} - \frac{\sigma(\mathbf{W}\mathbf{X}_{1:M} + \mathbf{b})\mathbf{y}_{1:M}}{M} \right) \right\|^2 \right]$.

$$\begin{aligned}
 & 1230 \quad \mathbb{E} \left[\left\| \left(\frac{\sigma(\mathbf{W}\mathbf{X}_{1:N} + \mathbf{b})\mathbf{y}_{1:N}}{N} - \frac{\sigma(\mathbf{W}\mathbf{X}_{1:M} + \mathbf{b})\mathbf{y}_{1:M}}{M} \right) \right\|^2 \right] \\
 & 1231 \quad \leq \sum_{j=1}^m \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N \sigma(\mathbf{w}_j^\top \mathbf{x}_i + b_j)y_i - \frac{1}{M} \sum_{i=1}^M \sigma(\mathbf{w}_j^\top \mathbf{x}_i + b_j)y_i \right)^2 \right] \\
 & 1232 \quad \\
 & 1233 \quad \\
 & 1234 \quad \\
 & 1235 \quad \\
 & 1236 \quad
 \end{aligned}$$

1237 For simplicity, we drop the subscript j and obtain

$$\begin{aligned}
 & 1238 \quad \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_i + b)y_i - \frac{1}{M} \sum_{i=1}^M \sigma(\mathbf{w}^\top \mathbf{x}_i + b)y_i \right)^2 \right] \\
 & 1239 \quad \leq 2\mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_i + b)y_i - \mathbb{E}[\sigma(\mathbf{w}^\top \mathbf{x}_i + b)y_i] \right)^2 \right] \\
 & 1240 \quad + 2\mathbb{E} \left[\left(\frac{1}{M} \sum_{i=1}^M \sigma(\mathbf{w}^\top \mathbf{x}_i + b)y_i - \mathbb{E}[\sigma(\mathbf{w}^\top \mathbf{x}_i + b)y_i] \right)^2 \right] \\
 & 1241 \quad \\
 & 1242 \quad \\
 & 1243 \quad \leq \frac{2}{N} \mathbb{E} \left[(\sigma(\mathbf{w}^\top \mathbf{x}_i + b)y_i)^2 \right] + \frac{2}{M} \mathbb{E} \left[(\sigma(\mathbf{w}^\top \mathbf{x}_i + b)y_i)^2 \right] \\
 & 1244 \quad \\
 & 1245 \quad \leq \frac{2}{N} \mathbb{E} \left[((\mathbf{w}^\top \mathbf{x}_i + b)y_i)^2 \right] + \frac{2}{M} \mathbb{E} \left[((\mathbf{w}^\top \mathbf{x}_i + b)y_i)^2 \right] \\
 & 1246 \quad \\
 & 1247 \quad \leq \frac{4}{N} \mathbb{E} \left[(\mathbf{w}^\top \mathbf{x}_i)^2 y_i^2 + b^2 y_i^2 \right] + \frac{4}{M} \mathbb{E} \left[(\mathbf{w}^\top \mathbf{x}_i)^2 y_i^2 + b^2 y_i^2 \right] \\
 & 1248 \quad \\
 & 1249 \quad \leq \frac{4}{N} \mathbb{E} \left[(\mathbf{w}^\top \mathbf{x}_i)^2 y_i^2 \right] + \frac{4}{M} \mathbb{E} \left[(\mathbf{w}^\top \mathbf{x}_i)^2 y_i^2 \right] + O(1/N + 1/M). \\
 & 1250 \quad \\
 & 1251 \quad \\
 & 1252 \quad \\
 & 1253 \quad \\
 & 1254 \quad \\
 & 1255 \quad \\
 & 1256 \quad
 \end{aligned}$$

1257 Then, we need to evaluate $\mathbb{E}[(\mathbf{w}^\top \mathbf{x}_i)^2 y_i^2]$; we need to be careful to obtain a tight bound for this value. Using the fact that

1258 pretrained \mathbf{w} almost aligns to the true subspace \mathcal{S} , we obtain

$$\begin{aligned}
 & 1259 \quad \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_i)^2 y_i^2] \leq \mathbb{E}[(\mathbf{w}_{1:r}^\top \mathbf{x}_{i,1:r} + \mathbf{w}_{r+1:d}^\top \mathbf{x}_{i,r+1:d})^2 y_i^2] \\
 & 1260 \quad \leq 2\|\mathbf{w}_{1:r}\|^2 \mathbb{E}[(\mathbf{x}_{i,1:r})^2 y_i^2] + 2\|\mathbf{w}_{r+1:d}\|^2 \mathbb{E}[(\mathbf{x}_{i,r+1:d})^2 y_i^2] \\
 & 1261 \quad \\
 & 1262 \quad \\
 & 1263 \quad \\
 & 1264 \quad
 \end{aligned}$$

$$= O(1) \cdot \mathbb{E}[(\mathbf{x}_{i,1:r})^2 y_i^2] + \tilde{O}(r/d) \cdot \mathbb{E}[(\mathbf{x}_{i,r+1:d})^2 y_i^2].$$

Moreover,

$$\begin{aligned} & \mathbb{E}[x_1^2 y^2] \\ & \leq 2\mathbb{E}_{c,\beta,s} \left[x_1^2 \left(\sum_{i=2}^P c_i^2 \right) \left(\sum_{i=2}^P \text{He}_i(\boldsymbol{\beta}^\top \mathbf{x})^2 \right) \right] + 2\mathbb{E}_{c,\beta,s} [x_1^2 s^2] \\ & \lesssim \sum_{i=2}^P \mathbb{E}_{c,\beta,s} [x_1^2 \text{He}_i(\boldsymbol{\beta}^\top \mathbf{x})^2] + O(1) \\ & \leq \sum_{i=2}^P \mathbb{E}[x_1^4]^{1/2} \mathbb{E}[\text{He}_i(\boldsymbol{\beta}^\top \mathbf{x})^4]^{1/2} + O(1) = O(1). \end{aligned}$$

Then, we obtain $\mathbb{E}[(\mathbf{w}^\top \mathbf{x}_i)^2 y_i^2] = \tilde{O}(r)$. Thus, an upper bound

$$\mathbb{E} \left[\left\| \left(\frac{\sigma(\mathbf{W} \mathbf{X}_{1:N} + \mathbf{b}) \mathbf{y}_{1:N}}{N} - \frac{\sigma(\mathbf{W} \mathbf{X}_{1:M} + \mathbf{b}) \mathbf{y}_{1:M}}{M} \right) \right\|^2 \right] = \tilde{O}(rm/N + rm/M)$$

is obtained.

Second, we bound $\mathbb{E} \left[\left\| \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\|^2 \right]$ as

$$\begin{aligned} \mathbb{E} \left[\left\| \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\|^2 \right] & \leq 2 \sum_{j=1}^m \mathbb{E}[(\mathbf{w}_j^\top \mathbf{x})^2 + b_j^2] \\ & \leq 2 \sum_{j=1}^m (b_j^2 + 4\|\mathbf{w}_{j,1:r}\|^2 \mathbb{E}[\|\mathbf{x}_{1:r}\|^2] + 4\|\mathbf{w}_{j,r+1:d}\|^2 \mathbb{E}[\|\mathbf{x}_{r+1:d}\|^2]) \\ & = \tilde{O}(mr). \end{aligned}$$

Putting all things together, we arrive at

$$\begin{aligned} & |\mathcal{R}_N(\boldsymbol{\Gamma}, \mathbf{W}, \mathbf{b}) - \mathcal{R}_M(\boldsymbol{\Gamma}, \mathbf{w}, \mathbf{b})| \\ & \leq \|\boldsymbol{\Gamma}\|_F \mathbb{E} \left[\left\| \left(\frac{\sigma(\mathbf{W} \mathbf{X}_{1:N} + \mathbf{b}) \mathbf{y}_{1:N}}{N} - \frac{\sigma(\mathbf{W} \mathbf{X}_{1:M} + \mathbf{b}) \mathbf{y}_{1:M}}{M} \right) \right\|^2 \right]^{1/2} \mathbb{E} \left[\left\| \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\|^2 \right]^{1/2} \\ & = \tilde{O} \left(\|\boldsymbol{\Gamma}\|_F \sqrt{r^2 m^2 / N + r^2 m^2 / M} \right). \end{aligned}$$

□

D.3. Proof of Theorem 1

Finally we are ready to prove our main theorem.

Proof. [Proof of Theorem 1] Let $\bar{\boldsymbol{\Gamma}}$ be the attention matrix constructed in Theorem 8 and let $\boldsymbol{\Gamma}^*$ be the minimizer of the ridge regression problem (line 6 in Algorithm 1). By the equivalence between optimization with L_2 regularization and norm-constrained optimization, there exists $\lambda_2 > 0$ such that

$$\|\boldsymbol{\Gamma}^*\|_F \leq \|\bar{\boldsymbol{\Gamma}}\|_F = O(\sqrt{r^{5P}}/m),$$

$$\begin{aligned}
 \left(\frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| \right)^2 &\leq \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} (y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b}))^2 \\
 &\leq \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} (y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \bar{\Gamma}, \mathbf{b}))^2.
 \end{aligned}$$

Then, from Theorem 8,

$$\frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| - \tau \lesssim \text{poly log}(d) \left(\frac{r^{5P/2}}{\sqrt{N_2}} + \frac{r^{3P/2}}{\sqrt{m}} \right)$$

holds.

We first evaluate $\mathcal{R}_{N_2}(f) - \tau$ where $f = f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})$: First,

$$\begin{aligned}
 &\mathcal{R}_{N_2}(f) - \tau \\
 &= \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| \\
 &\quad + \left(\mathcal{R}_{N_2}(f) - \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| \right) - \tau \\
 &\lesssim \text{poly log}(d) \left(\frac{r^{5P/2}}{\sqrt{N_2}} + \frac{r^{3P/2}}{\sqrt{m}} \right) \\
 &\quad + \sup_{f \in \mathcal{F}_{N,G,W,B}} \left(\mathcal{R}_{N_2}(f) - \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| \right)
 \end{aligned}$$

holds, noting that $f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b}) \in \mathcal{F}_{N,G,W,B}$ with $W = O(1), B = O(1)$ and $G = \tilde{O}(\sqrt{r^{5P}/m^2})$. We can evaluate the expectation value of the second term of the last line using the Rademacher complexity as

$$\begin{aligned}
 &\mathbb{E} \left[\sup_{f \in \mathcal{F}_{N,G,W,B}} \left(\mathcal{R}_{N_2}(f) - \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| \right) \right] \\
 &= \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, y} \left[\sup_{f \in \mathcal{F}_{N,G,W,B}} \left(\mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, y} [|y - f(\mathbf{X}, \mathbf{y}, \mathbf{x}; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})|] \right. \right. \\
 &\quad \left. \left. - \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| \right) \right] \\
 &\leq 2 \mathbb{E}_{\mathbf{X}, \mathbf{y}, \mathbf{x}, y, \epsilon} \left[\sup_{f \in \mathcal{F}_{N,G,W,B}} \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} \epsilon_t |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| \right] (\epsilon_t \stackrel{\text{i.i.d.}}{\sim} \text{Unif}\{\pm 1\}) \\
 &\lesssim \text{Rad}_{T_2}(\mathcal{F}_{N,G,W,B}) + \mathbb{E}_{y, \epsilon} \left[\frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} \epsilon_t y_t \right] (\because \text{Eq.(1) in (Maurer, 2016)}) \\
 &\lesssim \text{Rad}_{T_2}(\mathcal{F}_{N,G,W,B}) + \mathbb{E}_{y, \epsilon} \left[\left(\frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} \epsilon_t y_t \right)^2 \right]^{1/2} \\
 &\leq \text{Rad}_{T_2}(\mathcal{F}_{N,G,W,B}) + \frac{1}{\sqrt{T_2}} \mathbb{E}[y^2]^{1/2} = \tilde{O} \left(\text{polylog}(T_2) \frac{r^{5P/2}}{\sqrt{T_2}} \left(r^{P/2} \sqrt{d} + d \sqrt{d} \right) \right)
 \end{aligned}$$

by Proposition 20 and $\mathbb{E}[y^2] = O(1)$ from Assumption 1. Then, from Markov's inequality, we have $\sup_{f \in \mathcal{F}_{N,G,W,B}} \left(\mathcal{R}_{N_2}(f) - \frac{1}{T_2} \sum_{t=T_1+1}^{T_1+T_2} |y_t - f(\mathbf{X}_t, \mathbf{y}_t, \mathbf{x}_t; \mathbf{W}^{(1)}, \Gamma^*, \mathbf{b})| \right) =$

1375 $\tilde{O}\left(\text{polylog}(T_2) \frac{r^{5P/2}}{\sqrt{T_2}} \left(r^{P/2}\sqrt{d} + d\sqrt{d}\right)\right)$ with probability at least $1 - \delta$ where δ is a sufficiently small constant.
 1376 Then, we obtain
 1377

$$1378 \mathcal{R}_{N_2}(f) - \tau = \text{poly log}(d) \left(\frac{r^{5P/2}}{\sqrt{N_2}} + \frac{r^{3P/2}}{\sqrt{m}} + \text{polylog}(T_2) \frac{r^{5P/2}}{\sqrt{T_2}} \left(r^{P/2}\sqrt{d} + d\sqrt{d}\right) \right).$$

1381 Now we have done the upper bound for $\mathcal{R}_{N_2}(f)$. For $\mathcal{R}_{N^*}(f)$, we can use Proposition 24 because we have Corollary 7
 1382 and can ensure that the assumptions of Proposition 24 are satisfied. \square
 1383
 1384

1385 E. Derivation of Simplified Self-attention Module

1387 We derive equation (2.3), following the same line as (Zhang et al., 2023). Recall that the prediction of y by the original
 1388 self-attention module is defined as the right-bottom entry of
 1389

$$1390 f_{\text{Attn}} = \mathbf{E} + \mathbf{W}^P \mathbf{W}^V \mathbf{E} \cdot \text{softmax} \left(\frac{(\mathbf{W}^K \mathbf{E})^\top \mathbf{W}^Q \mathbf{E}}{\rho} \right),$$

1393 where the embedding matrix is given as (2.2). As mentioned in Section 2.2, we set $\rho = N$, omit softmax and merge
 1394 $\mathbf{W}^P \mathbf{W}^V$ as $\mathbf{W}^{PV} \in \mathbb{R}^{(m+1) \times (m+1)}$ and $(\mathbf{W}^K)^\top \mathbf{W}^Q$ as $\mathbf{W}^{KQ} \in \mathbb{R}^{(m+1) \times (m+1)}$.
 1395

1396 Now, we further assume that \mathbf{W}^{PV} \mathbf{W}^{KQ} are in the form as
 1397

$$1398 \mathbf{W}^{PV} = \begin{bmatrix} * & * \\ 0_{1 \times m} & v \end{bmatrix}, \mathbf{W}^{KQ} = \begin{bmatrix} \mathbf{K} & * \\ 0_{1 \times m} & * \end{bmatrix}.$$

1400 Then, we obtain the simplified form
 1401
 1402

$$1403 \tilde{f}_{\text{Attn}}(\mathbf{E}; \mathbf{W}^K, \mathbf{W}^Q, \mathbf{W}^V, \mathbf{W}^P) = \mathbf{E} + \begin{bmatrix} * & * \\ 0_{1 \times m} & v \end{bmatrix} \mathbf{E} \cdot \frac{\mathbf{E}^\top \begin{bmatrix} \mathbf{K} & * \\ 0_{1 \times m} & * \end{bmatrix} \mathbf{E}}{N}. \quad (\text{E.1})$$

1407 Note that we adopt the right-bottom entry $\left(\tilde{f}_{\text{Attn}}(\mathbf{E}; \mathbf{W}^K, \mathbf{W}^Q, \mathbf{W}^V, \mathbf{W}^P)\right)_{m+1, N+1}$ as prediction for a response of a
 1408 query. Then, by (E.1), we obtain
 1409
 1410

$$1411 \begin{aligned} & \left(\tilde{f}_{\text{Attn}}(\mathbf{E}; \mathbf{W}^K, \mathbf{W}^Q, \mathbf{W}^V, \mathbf{W}^P)\right)_{m+1, N+1} \\ &= \left(\begin{bmatrix} 0_{1 \times m} & v \end{bmatrix} \mathbf{E} \cdot \frac{\mathbf{E}^\top \begin{bmatrix} \mathbf{K} & * \\ 0_{1 \times m} & * \end{bmatrix} \mathbf{E}}{N} \right)_{m+1, N+1} \\ &= \begin{bmatrix} 0_{1 \times m} & v \end{bmatrix} \mathbf{E} \cdot \frac{\mathbf{E}^\top \begin{bmatrix} \mathbf{K} & * \\ 0_{1 \times m} & * \end{bmatrix} \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \\ 0 \end{bmatrix}}{N} \\ &= v [y_1 \quad \cdots \quad y_N \quad 0] \cdot \frac{\mathbf{E}^\top \begin{bmatrix} \mathbf{K} \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \\ 0 \end{bmatrix}}{N} \end{aligned}$$

1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484

$$\begin{aligned}
 & \sigma(\mathbf{W}\mathbf{X} + \mathbf{b})^\top \mathbf{K} \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \\
 &= v[y_1 \ \cdots \ y_N] \cdot \frac{\quad}{N} \\
 &= \left\langle \frac{v\mathbf{K}^\top \sigma(\mathbf{W}\mathbf{X} + \mathbf{b})\mathbf{y}}{N}, \begin{bmatrix} \sigma(\mathbf{w}_1^\top \mathbf{x} + b_1) \\ \vdots \\ \sigma(\mathbf{w}_m^\top \mathbf{x} + b_m) \end{bmatrix} \right\rangle.
 \end{aligned}$$

Letting $\mathbf{\Gamma} = v\mathbf{K}^\top$ yields equation (2.3).