

# Towards an Understanding of Residual Networks Using Neural Tangent Hierarchy (NTH)

Yuqing Li<sup>1</sup>, Tao Luo<sup>2,\*</sup> and Nung Kwan Yip<sup>3</sup>

<sup>1</sup> School of Mathematical Sciences, CMA-Shanghai, Shanghai Jiao Tong University, Shanghai, 200240, P.R. China.

<sup>2</sup> School of Mathematical Sciences, CMA-Shanghai, Institute of Natural Sciences, MOE-LSC, and Qing Yuan Research Institute, Shanghai Jiao Tong University, Shanghai, 200240, P.R. China.

<sup>3</sup> Department of Mathematics, Purdue University, IN, 47907, USA.

Received 7 December 2021; Accepted 11 April 2022

---

**Abstract.** Gradient descent yields zero training loss in polynomial time for deep neural networks despite non-convex nature of the objective function. The behavior of network in the infinite width limit trained by gradient descent can be described by the Neural Tangent Kernel (NTK) introduced in [25]. In this paper, we study dynamics of the NTK for finite width Deep Residual Network (ResNet) using the neural tangent hierarchy (NTH) proposed in [24]. For a ResNet with smooth and Lipschitz activation function, we reduce the requirement on the layer width  $m$  with respect to the number of training samples  $n$  from quartic to cubic. Our analysis suggests strongly that the particular skip-connection structure of ResNet is the main reason for its triumph over fully-connected network.

**AMS subject classifications:** 68U99, 90C26, 34A45

**Key words:** Residual networks, training process, neural tangent kernel, neural tangent hierarchy.

---

## 1 Introduction

Deep neural networks have achieved transcendent performance in a wide range of tasks such as speech recognition [9], computer vision [39], and natural language processing [8]. There are various methods to train neural networks, such as first-order gradient based methods like Gradient Descent (**GD**) and Stochastic Gradient Descent (**SGD**), which have been proven to achieve satisfactory results [20]. Experiments in [49] established

---

\*Corresponding author. *Email addresses:* liyuqing\_551@sjtu.edu.cn (Y. Li), luotao41@sjtu.edu.cn (T. Luo), yipn@purdue.edu (N. K. Yip)

that, even though with a random labeling of the training images, if one trains the state-of-the-art convolutional network for image classification using SGD, the network is still able to fit them well. There are numerous works trying to demystify such phenomenon theoretically. Du et al. [12, 15] proved that GD can obtain zero training loss for deep and shallow neural networks, and Zou et al. [52] analyzed the convergence of SGD on networks assembled with Rectified Linear Unit (**ReLU**) activation function. All these results are built upon the over-parameterized regime, and it is widely accepted that over-parameterization enables the neural network to fit all training data and bring no harm to the power of its generalization [49]. In particular, the deep neural networks that evaluated positions and selected moves for the well-known program AlphaGo are highly over-parameterized [41, 42].

Another advance is the outstanding performance of Deep Residual Network (**ResNet**), initially proposed by He et al. [22]. ResNet is arguably one of the most groundbreaking works in deep learning, in that it can train up to hundreds or even thousands of layers and still achieves compelling performance [23]. Recent works have shown that ResNet can utilize the features in transfer learning with better efficiency, and its residual link structure enables faster convergence of the training loss [45, 48]. Theoretically, Hardt and Ma [21] proved that for any residual linear networks with arbitrary depth, there are no spurious local optima. Du et al. [12] showed that in the scope of the convergence of GD via over-parameterization for different networks, training ResNet requires weaker conditions compared with fully-connected networks. Apart from that, the advantages of using residual connections remain to be discovered.

In this paper, we contribute to the further understanding of the above two aspects and make improvements in the analysis of their performance. We use the same ResNet structure as in [12]. (Details of the network structure are provided in Section 3.2.) The ResNet has  $L$  layers with width  $m$ . We will assume that the  $n$  data points are not parallel with each other. Such an assumption holds in general for a standard dataset [15]. We focus on the empirical risk minimization problem given by the quadratic loss and the activation function is 1-Lipschitz and analytic. We show that if  $m = \Omega(n^3 L^2)$ , then the empirical risk  $R_S(\boldsymbol{\theta}_t)$  under GD decays exponentially. More precisely,

$$R_S(\boldsymbol{\theta}_t) \leq R_S(\boldsymbol{\theta}_0) \exp\left(-\frac{\lambda t}{n}\right),$$

where  $\lambda$  is the least eigenvalue of  $\mathbf{K}^{[L+1]}$ , definition of which can be found in (4.2).

It is worth noticing that

- Given identical ResNet architectures, for the convergence of randomly initialized GD, our results improve upon [12] in the required number of width per layer from  $m = \Omega(n^4 L^2)$  to  $m = \Omega(n^3 L^2)$  (Corollary 4.1).
- For fully-connected network, the required amount of over-parameterization in [24] is  $m = \Omega(n^3 2^{\mathcal{O}(L)})$ . We are able to reproduce the result of Du et al. [12], showing

that the exponential dependence of  $m$  on the number of layers  $L$  can be eliminated for ResNet.

Our work is mainly motivated by the framework proposed by Huang and Yau [24], in which an infinite hierarchy of ordinary differential equations, the neural tangent hierarchy (NTH) is derived. Huang and Yau applied NTH to a fully-connected feedforward network and showed that it is possible for us to directly study the change of the neural tangent kernel (NTK) [25], and NTH outperforms kernel regressions using the corresponding limiting NTK.

Different from Huang and Yau's work in analyzing the fully-connected network, ResNet is investigated in our paper. We exploit the benefits of using ResNet architecture for training and the advantage of choosing NTH over kernel regression. In Section 5, an overview of our technique is provided.

The organization of the paper is listed as follows. In Section 2, we discuss some related works. In Section 3, we give some preliminary introductions to our problem. In Section 4, we state our main results for ResNet using NTH. In Section 5, we present an outline of our approach. We give some conclusions and future direction in Section 6. All the details of the proof are deferred to the Appendix.

## 2 Related works

In this section, we survey some previous works related to optimization aspect of neural networks.

Due to the non-convex nature of optimizing a neural network, it is challenging to locate the global optima. A popular way to analyze such optimization problems is to identify the geometric properties of each critical point. Some recent works have shown that for the set of functions satisfying: (i) all local minima are global and (ii) every saddle point possesses a negative curvature (i.e., it is non-degenerate), then GD can find a global optima [14, 17, 26, 31]. The objective functions of some shallow networks are in such set [11, 21, 38, 51]. The work [27] indicates that even for a three-layer linear network, there exists degenerate saddle points without negative curvature. So it is doubtful that global convergence of first-order methods can be guaranteed for deep neural networks.

Here we directly study the dynamics of the GD for a specific neural network architecture. This is another approach widely taken to obtain convergence results. Recently, it has been shown that if the network is over-parameterized, the SGD is able to find a global optima for two-layer networks [6, 13, 15, 16, 18, 33], deep linear networks [2, 5, 21] and ResNet [1, 12]. Jacot et al. [25] established that in the infinite width limit, the full batch GD corresponds to kernel regression predictor using the limiting NTK. Consequently, the convergence of GD for any 'infinite-width' neural network can be characterized by a fixed kernel [4, 25]. This is the cornerstone upon which rests the compelling performance of over-parameterization. In the regime of finite width, many works have suggested that the network can reduce training loss at exponential rate using GD [2, 12, 15, 16, 24]. As

the width increases, there are going to be small changes in the parameters during the whole training process [10, 52]. Such a variation of the parameters is crucial to the results we present, where the NTK of our ResNet behaves linearly in terms of its parameters throughout training (Theorem 4.2). Specifically, we use the results concerning the stability of the Gram matrices in [12] to demonstrate the benefits of choosing ResNet over fully-connected networks (Proposition C.3).

Some other works used optimal transport theory to analyze the mean field SGD dynamics of training neural networks in the large-width limit [7, 37, 40, 43]. However, their results are limited to one hidden layer networks, and their normalization factor  $1/m$  is different from our  $1/\sqrt{m}$  which is commonly employed in modern networks [19, 22].

### 3 Preliminaries

#### 3.1 Notations

We begin this section by introducing some notations that will be used in the rest of this paper. We set  $n$  for the number of input samples and  $m$  for the width of the neural network, and a special vector  $(1, 1, 1, \dots, 1)^\top \in \mathbb{R}^m$  by  $\mathbf{1} := (1, 1, 1, \dots, 1)^\top$ . We denote vector  $L^2$  norm as  $\|\cdot\|_2$ , vector or function  $L_\infty$  norm as  $\|\cdot\|_\infty$ , matrix spectral (operator) norm as  $\|\cdot\|_{2 \rightarrow 2}$ , matrix Frobenius norm as  $\|\cdot\|_F$ , matrix infinity norm as  $\|\cdot\|_{\infty \rightarrow \infty}$ , and a special matrix norm, matrix 2 to infinity norm as  $\|\cdot\|_{2 \rightarrow \infty}$ , which was shown to be useful in [15]. For a semi-positive-definite matrix  $A$ , we denote its smallest eigenvalue by  $\lambda_{\min}(A)$ . We use  $\mathcal{O}(\cdot)$  and  $\Omega(\cdot)$  for the standard Big-O and Big-Omega notations. We take  $C$  and  $c$  for some universal constants, which might vary from line to line.

Next we introduce a notion of high probability events that was also used in Huang and Yau [24, Section 1.3]. We say that an event holds with high probability if the probability of the event is at least  $1 - \exp(-m^\varepsilon)$  for some constant  $\varepsilon > 0$ . Since for a deep neural network in practice, we always have  $m \lesssim \text{poly}(n)$  and  $n \lesssim \text{poly}(m)$  [1, 28], then the intersection of a collection of many high probability events still has the same property as long as the number of events is at most polynomial in  $m$  and  $n$ .

#### 3.2 Problem setup

We shall focus on the empirical risk minimization problem given by quadratic loss:

$$\min_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{\alpha=1}^n \|f(\mathbf{x}_\alpha, \boldsymbol{\theta}) - y_\alpha\|_2^2. \quad (3.1)$$

In the above  $\{\mathbf{x}_\alpha\}_{\alpha=1}^n$  are the training inputs,  $\{y_\alpha\}_{\alpha=1}^n$  are the labels,  $f(\mathbf{x}_\alpha, \boldsymbol{\theta})$  is the prediction function, and  $\boldsymbol{\theta}$  are the parameters to be optimized, and their dependence is modeled by ResNet with  $L$  hidden layers, each of which has  $m$  neurons. Let  $\mathbf{x} \in \mathbb{R}^d$  be an input

sample, then the network has  $d$  input nodes. Let  $\mathbf{x}^{[l]}$  be the output of layer  $l$  with  $\mathbf{x}^{[0]} = \mathbf{x}$ . We consider the ResNet given below:

$$\begin{aligned} \mathbf{x}^{[1]} &= \sqrt{\frac{c_\sigma}{m}} \sigma(\mathbf{W}^{[1]} \mathbf{x}), \\ \mathbf{x}^{[l]} &= \mathbf{x}^{[l-1]} + \frac{c_{\text{res}}}{L\sqrt{m}} \sigma(\mathbf{W}^{[l]} \mathbf{x}^{[l-1]}), \quad \text{for } 2 \leq l \leq L, \end{aligned} \tag{3.2}$$

where  $\sigma(\cdot)$  is the activation function applied coordinate-wisely to its input. We assume that  $\sigma(\cdot)$  is 1-Lipschitz and smooth. The constant  $c_\sigma = (\mathbb{E}_{x \sim \mathcal{N}(0,1)} [\sigma(x)^2])^{-1}$  is a scaling factor serving for the purpose of normalization, and  $0 < c_{\text{res}} < 1$  is a small constant. Moreover, we have a series of weight matrices  $\{\mathbf{W}^{[l]}\}_{l=1}^L$ . Note that  $\mathbf{W}^{[l]} \in \mathbb{R}^{m \times d}$  for  $l = 1$ , and  $\mathbf{W}^{[l]} \in \mathbb{R}^{m \times m}$  for  $2 \leq l \leq L$ . The output function of ResNet is

$$f_{\text{res}}(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{a}^\top \mathbf{x}^{[L]}, \tag{3.3}$$

where  $\mathbf{a} \in \mathbb{R}^m$  is the weight vector of the output layer. We denote the vector containing all parameters by  $\boldsymbol{\theta} = (\text{vec}(\mathbf{W}^{[L]}), \text{vec}(\mathbf{W}^{[L-1]}), \dots, \text{vec}(\mathbf{W}^{[1]}), \mathbf{a})$ . Such a parameterization has been employed widely, see [12, 15, 30]. We shall initialize the parameter vector  $\boldsymbol{\theta}_0$  following the adopted Xavier initialization scheme [19], i.e.,  $W_{i,j}^{[l]} \sim \mathcal{N}(0,1)$ ,  $a_k \sim \mathcal{N}(0,1)$ , where  $\mathcal{N}(0,1)$  denotes the standard Gaussian distribution. Applying the continuous time GD for the loss function (3.1), we have for any time  $t \geq 0$ :

$$\partial_t \mathbf{W}_t^{[l]} = -\partial_{\mathbf{W}^{[l]}} R_S(\boldsymbol{\theta}_t), \quad l = 1, 2, \dots, L, \tag{3.4}$$

$$\partial_t \mathbf{a}_t = -\partial_{\mathbf{a}} R_S(\boldsymbol{\theta}_t). \tag{3.5}$$

We use  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  for the set of input samples,  $\sigma(\mathbf{W}^{[l]} \mathbf{x}_\alpha^{[l-1]})$  as  $\sigma_{[l]}(\mathbf{x}_\alpha)$ , and the diagonal matrix generated by the  $r$ -th derivatives of  $\sigma_{[l]}(\mathbf{x}_\alpha)$ , i.e.,  $\text{diag}(\sigma^{(r)}(\mathbf{W}^{[l]} \mathbf{x}_\alpha^{[l-1]}))$  by  $\sigma_{[l]}^{(r)}(\mathbf{x}_\alpha)$ , where  $r \geq 1$ . We also write the output function  $f_{\text{res}}(\mathbf{x}_\alpha, \boldsymbol{\theta}_t)$  as  $f_\alpha(t)$ . Moreover, we shall define a series of special matrices. Using  $\mathbf{I}_m$  to signify the identity matrix in  $\mathbb{R}^{m \times m}$ , we define for  $2 \leq l \leq L$ :

$$\mathbf{E}_{t,\alpha}^{[l]} := \left( \mathbf{I}_m + \frac{c_{\text{res}}}{L} \sigma_{[l]}^{(1)}(\mathbf{x}_\alpha) \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right). \tag{3.6}$$

The above matrices are termed skip-connection matrices. Given  $\{\mathbf{E}_{t,\alpha}^{[l]}\}_{l=2}^L$ , we let  $\mathbf{E}_{t,\alpha}^{[l:L]}$  be the direct parameterization of the end-to-end mapping realized by the group of skip-connection matrices, i.e.,  $\mathbf{E}_{t,\alpha}^{[l:L]} := \mathbf{E}_{t,\alpha}^{[L]} \mathbf{E}_{t,\alpha}^{[L-1]} \dots \mathbf{E}_{t,\alpha}^{[l]}$ , where we set  $\mathbf{E}_{t,\alpha}^{[i:j]} := \mathbf{I}_m$ ,  $i > j$  for completeness.

With the above notations, the continuous time GD dynamics (3.4) and (3.5) can be

written as

$$\partial_t \mathbf{a}_t = -\frac{1}{n} \sum_{\beta=1}^n \mathbf{x}_\beta^{[L]} (f_\beta(t) - y_\beta), \tag{3.7}$$

$$\partial_t \mathbf{W}_t^{[L]} = -\frac{1}{n} \sum_{\beta=1}^n \frac{c_{\text{res}}}{L\sqrt{m}} \text{diag} \left( \sigma_{[L]}^{(1)}(\mathbf{x}_\beta) \mathbf{a}_t \right) \mathbf{1} \otimes (\mathbf{x}_\beta^{[L-1]})^\top (f_\beta(t) - y_\beta), \tag{3.8}$$

$$\partial_t \mathbf{W}_t^{[l]} = -\frac{1}{n} \sum_{\beta=1}^n \frac{c_{\text{res}}}{L\sqrt{m}} \text{diag} \left( \sigma_{[l]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+1):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \otimes (\mathbf{x}_\beta^{[l-1]})^\top (f_\beta(t) - y_\beta), \tag{3.9}$$

for  $l = 2, 3, \dots, L-1$ ,

$$\partial_t \mathbf{W}_t^{[1]} = -\frac{1}{n} \sum_{\beta=1}^n \sqrt{\frac{c_\sigma}{m}} \text{diag} \left( \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \otimes (\mathbf{x}_\beta)^\top (f_\beta(t) - y_\beta). \tag{3.10}$$

### 3.3 Neural tangent kernel

The Neural Tangent Kernel (NTK) is introduced in Jacot et al. [25]. For any parametrized function  $f(\mathbf{x}, \boldsymbol{\theta}_t)$ , it is defined as:

$$\mathcal{K}_{\boldsymbol{\theta}_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \langle \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_\alpha, \boldsymbol{\theta}_t), \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_\beta, \boldsymbol{\theta}_t) \rangle.$$

In the situations where  $f(\mathbf{x}, \boldsymbol{\theta}_t)$  is the output of a fully-connected feedforward network with appropriate scaling factor  $1/\sqrt{m}$  for the parameters, there is an infinite width limit ( $m \rightarrow \infty$ ) of  $\mathcal{K}_{\boldsymbol{\theta}_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$ , denoted by  $\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)$ . This result allows them to capture the behavior of fully-connected feedforward network trained by GD in the infinite width limit. More precisely, the output function  $f(\mathbf{x}, \boldsymbol{\theta}_t)$  evolves as a linear differential equation:

$$\partial_t f(\mathbf{x}, \boldsymbol{\theta}_t) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{K}_\infty(\mathbf{x}, \mathbf{x}_\beta) (f(\mathbf{x}_\beta, \boldsymbol{\theta}_t) - y_\beta). \tag{3.11}$$

Note that the training dynamic is identical to the dynamics of kernel regression under gradient flow. Also we note that  $\mathcal{K}_\infty(\cdot)$  only depends on the training inputs. More importantly,  $\mathcal{K}_\infty(\cdot)$  is independent of the neural network parameters  $\boldsymbol{\theta}$  [4, 12, 15]. Similar result holds for our ResNet structure.

The finding above is groundbreaking in that it provides us an analytically tractable equation to predict the behavior of GD. However, the convergence  $\mathcal{K}_{\boldsymbol{\theta}_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$  to  $\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)$  is proved in the regime of infinite width. This is unrealistic in nature. Some concurrent works concerning various network structures [2, 3, 12, 15, 32, 44] have extended the result in [25] to the regime of finite width. For a two-layer network with ReLU, the required width  $m$  in Song and Yang [44] is  $m = \Omega(n^2 \text{poly}(\log(n)))$  under some strong assumptions on the input data. For fully-connected feedforward network, Huang and Yau requires width  $m = \Omega(n^3 \log(n) 2^{\mathcal{O}(L)})$ . Finally, for ResNet which is the main focus of

our paper, the required width for Du et al. [12] is  $m = \Omega(n^4 L^2)$  with iteration complexity  $T = \Omega(n^2 \log(\frac{1}{\epsilon}))$ . Our Corollary 4.1 only requires  $m = \Omega(n^3 L^2)$  and  $T = \Omega(n \log(\frac{1}{\epsilon}))$ .

We now write out the NTK for ResNet:

$$\begin{aligned} \partial_t(f_\alpha(t) - y_\alpha) &= -\nabla_{\theta} f_\alpha(t) \cdot \nabla_{\theta} R_S(\theta_t) \\ &= -\frac{1}{n} \nabla_{\theta} f_\alpha(t) \cdot \sum_{\beta=1}^n \nabla_{\theta} f_\beta(t) (f_\beta(t) - y_\beta) \\ &= -\frac{1}{n} \sum_{\beta=1}^n \mathcal{K}_{\theta_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta) (f_\beta(t) - y_\beta), \end{aligned} \tag{3.12}$$

using Eqs. (3.7), (3.8), (3.9) and (3.10), the NTK  $\mathcal{K}_{\theta_t}(\cdot)$  is given below

$$\mathcal{K}_{\theta_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \langle \nabla_{\theta} f_\alpha(t), \nabla_{\theta} f_\beta(t) \rangle = \sum_{l=1}^{L+1} \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta), \tag{3.13}$$

where

$$\begin{aligned} \mathcal{G}_t^{[1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \langle \partial_{\mathbf{W}^{[1]}} f_\alpha(t), \partial_{\mathbf{W}^{[1]}} f_\beta(t) \rangle \\ &= \left\langle \sqrt{\frac{c_\sigma}{m}} \sigma_{[1]}^{(1)}(\mathbf{x}_\alpha) \left(\mathbf{E}_{t,\alpha}^{[2:L]}\right)^\top \mathbf{a}_t, \sqrt{\frac{c_\sigma}{m}} \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left(\mathbf{E}_{t,\beta}^{[2:L]}\right)^\top \mathbf{a}_t \right\rangle \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle, \end{aligned}$$

for  $2 \leq l \leq L$ ,

$$\begin{aligned} \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \langle \partial_{\mathbf{W}^{[l]}} f_\alpha(t), \partial_{\mathbf{W}^{[l]}} f_\beta(t) \rangle \\ &= \left\langle \frac{c_{\text{res}}}{L\sqrt{m}} \sigma_{[l]}^{(1)}(\mathbf{x}_\alpha) \left(\mathbf{E}_{t,\alpha}^{[(l+1):L]}\right)^\top \mathbf{a}_t, \frac{c_{\text{res}}}{L\sqrt{m}} \sigma_{[l]}^{(1)}(\mathbf{x}_\beta) \left(\mathbf{E}_{t,\beta}^{[(l+1):L]}\right)^\top \mathbf{a}_t \right\rangle \langle \mathbf{x}_\alpha^{[l-1]}, \mathbf{x}_\beta^{[l-1]} \rangle, \end{aligned}$$

and finally

$$\mathcal{G}_t^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \langle \partial_{\mathbf{a}} f_\alpha(t), \partial_{\mathbf{a}} f_\beta(t) \rangle = \langle \mathbf{x}_\alpha^{[L]}, \mathbf{x}_\beta^{[L]} \rangle.$$

We remark that all the  $\mathcal{G}_t^{[l]}$  depends on  $\theta_t$  but for simplicity it is not explicitly written.

## 4 Main results

### 4.1 Activation function and input samples

In this paper, we will impose some following technical conditions on the activation function and input samples.

**Assumption 4.1.** The activation function  $\sigma(\cdot)$  is smooth, and there exists a universal constant  $0 < C_{\text{Lip}} \leq 1$  such that for any  $r \geq 1$ , its  $r$ -th derivative and the function value at 0 satisfy

$$|\sigma(0)|, \|\sigma^{(r)}(\cdot)\|_{\infty} \leq C_{\text{Lip}}. \tag{4.1}$$

We use  $C_{\text{Lip}}$  to avoid confusion with the usage of constant  $C$  in the Appendix. Here the subscript ‘Lip’ stands for Lipschitz condition. Note that Assumption 4.1 can be satisfied by using the softplus activation:

$$\sigma(x) = \ln(1 + \exp(x)).$$

Some other functions also satisfy this assumption, for instance, the sigmoid activation:

$$\sigma(x) = \frac{1}{1 + \exp(-x)}.$$

**Assumption 4.2.** The training inputs and labels satisfy  $\|\mathbf{x}_{\alpha}\|_2 = 1, |y_{\alpha}| \leq 1$ , for any  $\mathbf{x}_{\alpha} \in \mathcal{X}$ . All training inputs are non-parallel with each other, i.e.,  $\mathbf{x}_{\alpha_1} \not\parallel \mathbf{x}_{\alpha_2}$ , for any  $\alpha_1 \neq \alpha_2$ .

Assumption 4.2 guarantees that some of the Gram matrices defined in Section 4.2 are strictly positive definite.

## 4.2 Gram matrices

Recent works [15, 44, 50] have shown that the convergence of the outputs of neural networks are determined by the spectral property of Gram matrices. Here we define the key Gram matrices  $\{\mathbf{K}^{[l]}\}_{l=1}^{L+1}$  below. We more or less follow the definition of the Gram matrices partially from [12, Definition 6.1]. Also we note that the Gram matrices depend on the series of matrices  $\{\tilde{\mathbf{K}}^{[l]}\}_{l=1}^L, \{\tilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$ , and the series of vectors  $\{\tilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , which are listed out as follows, for  $2 \leq l \leq L$

$$\begin{aligned} \tilde{\mathbf{K}}_{ij}^{[0]} &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle, \\ \tilde{\mathbf{K}}_{ij}^{[1]} &= \mathbb{E}_{(u,v)^{\top} \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[0]} & \tilde{\mathbf{K}}_{ij}^{[0]} \\ \tilde{\mathbf{K}}_{ji}^{[0]} & \tilde{\mathbf{K}}_{jj}^{[0]} \end{pmatrix}\right)} c_{\sigma} \sigma(u) \sigma(v), \\ \tilde{\mathbf{b}}_i^{[1]} &= \sqrt{c_{\sigma}} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[0]})} [\sigma(u)], \end{aligned}$$



$$\begin{aligned} \tilde{\mathbf{A}}_{ij}^{[l]} &= \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[l-1]} & \tilde{\mathbf{K}}_{ij}^{[l-1]} \\ \tilde{\mathbf{K}}_{ji}^{[l-1]} & \tilde{\mathbf{K}}_{jj}^{[l-1]} \end{pmatrix}, \\ \tilde{\mathbf{K}}_{ij}^{[l]} &= \tilde{\mathbf{K}}_{ij}^{[l-1]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[l]})} \left[ \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[l-1]} \sigma(v)}{L} + \frac{c_{\text{res}} \tilde{\mathbf{b}}_j^{[l-1]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(v)}{L^2} \right], \\ \tilde{\mathbf{b}}_i^{[l]} &= \tilde{\mathbf{b}}_i^{[l-1]} + \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)], \\ \tilde{\mathbf{A}}_{ij}^{[L+1]} &= \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[L]} & \tilde{\mathbf{K}}_{ij}^{[L]} \\ \tilde{\mathbf{K}}_{ji}^{[L]} & \tilde{\mathbf{K}}_{jj}^{[L]} \end{pmatrix}, \end{aligned}$$

then we may proceed to the definitions of Gram matrices for  $l = L + 1$  and  $L$ .

**Definition 4.1.** Given the input samples  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , the Gram matrix  $\mathbf{K}^{[L+1]} \in \mathbb{R}^{n \times n}$  is defined as follows, for  $1 \leq i, j \leq n$ ,

$$\mathbf{K}_{ij}^{[L+1]} = \tilde{\mathbf{K}}_{ij}^{[L]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[L+1]})} \left[ \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[L]} \sigma(v)}{L} + \frac{c_{\text{res}} \tilde{\mathbf{b}}_j^{[L]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(v)}{L^2} \right], \quad (4.2)$$

and Gram matrix  $\mathbf{K}^{[L]} \in \mathbb{R}^{n \times n}$  is also defined as follows, for  $1 \leq i, j \leq n$ ,

$$\mathbf{K}_{ij}^{[L]} = \frac{c_{\text{res}}^2}{L^2} \tilde{\mathbf{K}}_{ij}^{[L-1]} \mathbb{E}_{(u,v)^\top \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[L]})} [\sigma^{(1)}(u) \sigma^{(1)}(v)]. \quad (4.3)$$

Note that matrix  $\mathbf{K}^{[L]}$  coincides with  $\mathbf{K}^{[H]}$  given by [12, Definition 6.1]. Now that given the definition of  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$ , we need to move forward to the definition of other Gram matrices  $\{\mathbf{K}^{[l]}\}_{l=1}^{L-1}$ . Since it is challenging to give an explicit formula for the series of matrices  $\{\mathbf{K}^{[l]}\}_{l=1}^{L-1}$ , we shall use a slightly different approach to write out the definitions for these matrices.

**Definition 4.2.** Given the input samples  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , the Gram matrices  $\mathbf{K}^{[l]} \in \mathbb{R}^{n \times n}$  are defined as follows, for  $1 \leq i, j \leq n, 2 \leq l \leq L - 1$ ,

$$\mathbf{K}_{ij}^{[l]} = \frac{c_{\text{res}}^2}{L^2} \tilde{\mathbf{K}}_{ij}^{[l-1]} \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \sigma_{[l]}^{(1)}(\mathbf{x}_i) \left( \mathbf{E}_{0,i}^{[(l+1):L]} \right)^\top \mathbf{a}_0, \sigma_{[l]}^{(1)}(\mathbf{x}_j) \left( \mathbf{E}_{0,j}^{[(l+1):L]} \right)^\top \mathbf{a}_0 \right\rangle, \quad (4.4)$$

and for  $1 \leq i, j \leq n, l = 1$ ,

$$\mathbf{K}_{ij}^{[1]} = c_\sigma \tilde{\mathbf{K}}_{ij}^{[0]} \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \sigma_{[1]}^{(1)}(\mathbf{x}_i) \left( \mathbf{E}_{0,i}^{[2:L]} \right)^\top \mathbf{a}_0, \sigma_{[1]}^{(1)}(\mathbf{x}_j) \left( \mathbf{E}_{0,j}^{[2:L]} \right)^\top \mathbf{a}_0 \right\rangle. \quad (4.5)$$

**Remark 4.1.** Thanks to the Strong Law of Large Numbers, the above limit exists [4]. Since we send  $m \rightarrow \infty$ , the Gram matrices  $\{\mathbf{K}^{[l]}\}_{l=1}^{L-1}$  only depend on the input samples and the activation patterns.

Moreover, in Section B, we show that under Assumption 4.2 and width  $m \sim n^2$ ,  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$  are strictly positive definite.

### 4.3 Convergence of gradient descent

Here we state our main theorems for the NTH of ResNet.

**Theorem 4.1.** *Under Assumptions 4.1 and 4.2, there exists an infinite family of operators  $\mathcal{K}_t^{(r)} : \mathcal{X}^r \rightarrow \mathbb{R}$ ,  $r \geq 2$  that describes the continuous time GD:*

$$\partial_t(f_\alpha(t) - y_\alpha) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{K}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta)(f_\beta(t) - y_\beta), \tag{4.6}$$

and for  $r \geq 2$ , we have

$$\partial_t \mathcal{K}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r}) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{K}_t^{(r+1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r}, \mathbf{x}_\beta)(f_\beta(t) - y_\beta). \tag{4.7}$$

Moreover, with high probability with respect to random initialization, there exist some constants  $C, C' > 0$ , such that for time  $0 \leq t \leq \sqrt{m} / (\ln m)^{C'}$ , the following holds:

$$\left\| \mathcal{K}_t^{(2)}(\cdot) \right\|_\infty \lesssim 1, \tag{4.8}$$

and for  $r \geq 3$ ,

$$\left\| \mathcal{K}_t^{(r)}(\cdot) \right\|_\infty \lesssim \frac{(\ln m)^C}{m^{r/2-1}}. \tag{4.9}$$

**Remark 4.2.** The operator  $\mathcal{K}_t^{(2)}(\cdot)$  by definition is the same as the NTK  $\mathcal{K}_{\theta_t}(\cdot)$  derived in (3.12).

We note that as  $r$  increases, the pre-factor in (4.9) explodes exponentially fast in  $r$ . However, this will not significantly affect the convergence of GD. Firstly, only some lower order kernels need to be analyzed. As is shown in the proof of Corollary 4.1, only kernels up to order  $r=4$  will be used. Secondly, we shall recall the NTK  $\mathcal{K}_t^{(2)}(\cdot)$  derived in (3.13) :

$$\mathcal{K}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \langle \nabla_{\theta} f_\alpha(t), \nabla_{\theta} f_\beta(t) \rangle = \sum_{l=1}^{L+1} \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta), \tag{4.10}$$

in the case of Huang and Yau [24], for a fully-connected feedforward network, since all those kernels  $\mathcal{G}_t^{[l]}$  are positive definite, then the sum of the least eigenvalue of all the kernels  $\mathcal{G}_t^{[l]}$  is much larger than the counterpart of a single kernel, i.e.,

$$\lambda_{\min} \left[ \mathcal{K}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} \gg \lambda_{\min} \left[ \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n}.$$

However, adding up all the kernels will not give substantial increase to the least eigenvalue for the ResNet. Since there exists a scaling factor  $\frac{1}{L^2}$  for the kernels  $\mathcal{G}_t^{[l]}$ , where

$2 \leq l \leq L$ , then heuristically, the gap of the least eigenvalues between  $\mathcal{K}_t^{(2)}(\cdot)$  and  $\mathcal{G}_t^{[L+1]}(\cdot) + \mathcal{G}_t^{[1]}(\cdot)$  is at most of order  $\mathcal{O}(\frac{L-1}{L^2}) = \mathcal{O}(\frac{1}{L})$ . Hence for ResNet, we shall see that even if the depth  $L$  gets larger, the least eigenvalue of the NTK is still concentrated on the kernels  $\mathcal{G}_t^{[L+1]}(\cdot)$  and  $\mathcal{G}_t^{[1]}(\cdot)$ . Thanks to that observation, we only need to bring the kernel  $\mathcal{G}_t^{[L+1]}(\cdot)$  to the spotlight. We omit the analysis of  $\mathcal{G}_t^{[1]}(\cdot)$  because it is not needed in our proof.

It was proven in Theorem 4.1 and other literatures [4, 12, 47] that the change of NTK during the dynamics for Deep Neural Network is bounded by  $\mathcal{O}(\frac{1}{\sqrt{m}})$ . However, it was observed by Lee et al. [30] that the time variation of the NTK is closer to  $\mathcal{O}(\frac{1}{m})$ , indicating that there exists a performance gap between the kernel regression using the limiting NTK and neural networks. Such an observation has been confirmed by Huang and Yau [24] and listed out as Corollary 2.4. in their paper. We use a different approach to obtain similar results and state them as Theorem 4.2.

**Theorem 4.2.** *Under Assumptions 4.1 and 4.2, with high probability with respect to random initialization, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^C$ , the following holds:*

$$\left\| \partial_t \mathcal{G}_t^{[L+1]}(\cdot) \right\|_\infty \lesssim \frac{(1+t)(\ln m)^C}{m}, \tag{4.11}$$

where the constant  $C$  is independent of the depth  $L$ . Moreover, the pre-factor in (4.11) is at most of order  $\mathcal{O}(L^2)$ .

As a direct consequence of Theorem 4.1, for the ResNet defined in (3.2), with width  $m \sim n^3$ , the GD converges to zero training loss at a linear rate. The precise statement is given in the following.

**Corollary 4.1.** *Under Assumptions 4.1 and 4.2, with  $\mathbf{K}^{[L+1]}$  defined in (4.2), we have that for some  $\lambda_0 > 0$ ,  $\lambda_{\min}(\mathbf{K}^{[L+1]}) > \lambda_0$ . Equipped with this, we have the following two statements.*

*There exists a small constant  $\gamma_1 > 0$ , such that for  $m = \Omega((\frac{n}{\lambda_0})^{2+\gamma_1})$ , with high probability with respect to random initialization, the following holds:*

$$\lambda_{\min} \left[ \mathcal{K}_0^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} \geq \frac{3}{4} \lambda_0. \tag{4.12}$$

Furthermore, there exists a small constant  $\gamma_2 > 0$ , such that for

$$m = \Omega \left( \left( \frac{n}{\lambda_0} \right)^{3+\gamma_2} L^2 \ln \left( \frac{1}{\varepsilon} \right)^2 \right),$$

where  $\varepsilon > 0$  is the desired accuracy for  $R_S(\boldsymbol{\theta}_t)$ , then the training loss  $R_S(\boldsymbol{\theta}_t)$  decays exponentially with respect to time  $t \geq 0$ ,

$$R_S(\boldsymbol{\theta}_t) \leq R_S(\boldsymbol{\theta}_0) \exp \left( -\frac{\lambda t}{n} \right). \tag{4.13}$$

For convenience, we summarize the above statement in the following manner. If

$$m = \max \left\{ \Omega \left( \left( \frac{n}{\lambda} \right)^{2+\gamma_1} \right), \Omega \left( \left( \frac{n}{\lambda} \right)^{3+\gamma_2} L^2 \ln \left( \frac{1}{\varepsilon} \right)^2 \right) \right\}, \tag{4.14}$$

then the continuous GD converges exponentially and reaches the training accuracy  $\varepsilon$  with time complexity

$$T = \mathcal{O} \left( \frac{n}{\lambda} \ln \left( \frac{1}{\varepsilon} \right) \right). \tag{4.15}$$

Before we end this section, we present a fair comparison of our result with others. First of all, Du et al. [12, Theorem 6.1.] required  $m = \Omega \left( \frac{n^4}{\lambda_{\min}(\mathbf{K}^{[L]})^4 L^6} \right)$ . Since there is a scaling factor  $\frac{1}{L^2}$  in  $\lambda_{\min}(\mathbf{K}^{[L]})$ , this leads to  $m = \Omega(n^4 L^2)$ . Then their GD converges with iteration complexity  $T = \Omega(n^2 L^2 \ln(\frac{1}{\varepsilon}))$ . Our Corollary 4.1 improves this result in two ways: (i) The quartic dependence on  $n$  is reduced directly to cubic dependence. (ii) A faster convergence of the training process of GD.

Second, our work serves as an extension of the NTH proposed by Huang and Yau [24], which captures the GD dynamics for a fully-connected feedforward network. We show that not only it is possible to study directly the time variation of NTK for ResNet using NTH, but that ResNet possesses more stability in many aspects than fully-connected network. In particular, we improve their results in three aspects: (i) With ResNet architecture, the dependency of the amount of over-parameterization on the depth  $L$  can be reduced from their  $2^{\mathcal{O}(L)}$  to  $L^2$ . (ii) While the time interval for the result in [24] takes the form  $0 \leq t \leq m^{\frac{p}{2(p+1)}} / (\ln m)^{C'}$  for some  $p \geq 2$ , we extend the interval to  $0 \leq t \leq \sqrt{m} / (\ln m)^{C'}$ . Moreover, we are able to show even further that the results hold true for  $t \rightarrow \infty$  using techniques from [36]. (iii) In the proof of Corollary 2.5. in [24], a further assumption on the least eigenvalue of the NTK  $\mathcal{K}_t^{(2)}(\cdot)$  has been imposed directly, we show in Appendix C that the least eigenvalue of the NTK  $\mathcal{K}_t^{(2)}(\cdot)$  can be guaranteed with high probability as long as the width  $m$  satisfies  $m = \Omega(n^2)$ .

## 5 Technique overview

In this part we first describe some technical tools and present the sketch of proofs for Theorems 4.1 and 4.2 and Corollary 4.1.

### 5.1 Replacement rules

We revisit the NTK (3.12) derived in Section 3.3,

$$\mathcal{K}_{\theta_t}(x_\alpha, x_\beta) = \sum_{l=1}^{L+1} \mathcal{G}_t^{[l]}(x_\alpha, x_\beta), \tag{5.1}$$

Notice that  $\mathcal{K}_{\theta_t}(\cdot)$  coincides with  $\mathcal{K}_t^{(2)}(\cdot)$  in (3.12), and  $\mathcal{K}_t^{(2)}(\cdot)$  is the sum of  $L+1$  terms, with each term being the inner product of vectors containing the quantities  $\mathbf{a}_t, \mathbf{x}_\alpha^{[l]}, \mathbf{E}_{t,\alpha}^{[l]}$  and  $\sigma_{[l]}^{(1)}(\mathbf{x}_\alpha)$ . We are able to write down the dynamics of  $\mathbf{a}_t, \mathbf{x}_\alpha^{[l]}, \mathbf{E}_{t,\alpha}^{[l]}$  and  $\sigma_{[l]}^{(1)}(\mathbf{x}_\alpha)$  following GD, using Eqs. (3.7), (3.8) (3.9), (3.10) and chain rules. In order to shorten the space, we perform a similar replacement rule as in Huang and Yau [24]. For instance, the dynamics of  $\mathbf{a}_t$  is written as

$$\partial_t \mathbf{a}_t = -\frac{1}{n} \sum_{\beta=1}^n \frac{1}{\sqrt{m}} \sqrt{m} \mathbf{x}_\beta^{[L]} (f_\beta(t) - y_\beta). \tag{5.2}$$

For simplicity, we symbolize the dynamics (5.2) as  $\mathbf{a}_t \rightarrow \frac{1}{\sqrt{m}} \sqrt{m} \mathbf{x}_\beta^{[L]}$ . Similarly, for the dynamics of  $\mathbf{x}_\alpha^{[l]}, 2 \leq l \leq L$ , we have

$$\begin{aligned} \sqrt{m} \mathbf{x}_\alpha^{[1]} &\rightarrow \frac{c_\sigma}{\sqrt{m}} \text{diag} \left( \sigma_{[1]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle, \\ \sqrt{m} \mathbf{x}_\alpha^{[l]} &\rightarrow \frac{c_\sigma}{\sqrt{m}} \text{diag} \left( \mathbf{E}_{t,\alpha}^{[2:l]} \sigma_{[1]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle \\ &\quad + \sum_{k=2}^l \frac{c_{\text{res}}^2}{L^2 \sqrt{m}} \text{diag} \left( \mathbf{E}_{t,\alpha}^{[(k+1):l]} \sigma_{[k]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \langle \mathbf{x}_\alpha^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle, \end{aligned}$$

and of  $\sigma_{[l]}^{(1)}(\mathbf{x}_\alpha)$ , for  $2 \leq l \leq L-1, r \geq 1$ ,

$$\begin{aligned} \sigma_{[1]}^{(r)}(\mathbf{x}_\alpha) &\rightarrow \sqrt{\frac{c_\sigma}{m}} \sigma_{[1]}^{(r+1)}(\mathbf{x}_\alpha) \text{diag} \left( \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle, \\ \sigma_{[2]}^{(r)}(\mathbf{x}_\alpha) &\rightarrow \frac{c_{\text{res}}}{L \sqrt{m}} \sigma_{[2]}^{(r+1)}(\mathbf{x}_\alpha) \text{diag} \left( \sigma_{[2]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[3:L]} \right)^\top \mathbf{a}_t \right) \langle \mathbf{x}_\alpha^{[1]}, \mathbf{x}_\beta^{[1]} \rangle \\ &\quad + \frac{c_\sigma}{\sqrt{m}} \sigma_{[2]}^{(r+1)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[2]}}{\sqrt{m}} \sigma_{[1]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle, \\ \sigma_{[l+1]}^{(r)}(\mathbf{x}_\alpha) &\rightarrow \frac{c_{\text{res}}}{L \sqrt{m}} \sigma_{[l+1]}^{(r+1)}(\mathbf{x}_\alpha) \text{diag} \left( \sigma_{[l+1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+2):L]} \right)^\top \mathbf{a}_t \right) \langle \mathbf{x}_\alpha^{[l]}, \mathbf{x}_\beta^{[l]} \rangle \\ &\quad + \sum_{k=2}^l \frac{c_{\text{res}}^2}{L^2 \sqrt{m}} \sigma_{[l+1]}^{(r+1)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \mathbf{E}_{t,\alpha}^{[(k+1):l]} \sigma_{[k]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right) \\ &\quad \quad \quad \langle \mathbf{x}_\alpha^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle \\ &\quad + \frac{c_\sigma}{\sqrt{m}} \sigma_{[l+1]}^{(r+1)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \mathbf{E}_{t,\alpha}^{[2:l]} \sigma_{[1]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle, \end{aligned}$$

and finally of  $E_{t,\alpha}^{[l]}, 2 \leq l \leq L-1$ ,

$$\begin{aligned}
 E_{t,\alpha}^{[2]} &\rightarrow \frac{c_{\text{res}}^2}{L^2\sqrt{m}} \text{diag} \left( \sigma_{[2]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[2]}^{(1)}(\mathbf{x}_\beta) \left( E_{t,\beta}^{[3:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \otimes \left( \frac{\sqrt{m} \mathbf{x}_\beta^{[1]}}{m} \right)^\top \\
 &\quad + \frac{c_{\text{res}}}{L\sqrt{m}} \sigma_{[2]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \sigma_{[2]}^{(1)}(\mathbf{x}_\beta) \left( E_{t,\beta}^{[3:L]} \right)^\top \mathbf{a}_t \right) \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[2]}}{\sqrt{m}} \langle \mathbf{x}_\alpha^{[1]}, \mathbf{x}_\beta^{[1]} \rangle \\
 &\quad + \frac{c_\sigma}{\sqrt{m}} \sigma_{[2]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[2]}}{\sqrt{m}} \sigma_{[1]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left( E_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[2]}}{\sqrt{m}} \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle, \\
 E_{t,\alpha}^{[l+1]} &\rightarrow \frac{c_{\text{res}}^2}{L^2\sqrt{m}} \text{diag} \left( \sigma_{[l+1]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[l+1]}^{(1)}(\mathbf{x}_\beta) \left( E_{t,\beta}^{[(l+2):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \otimes \left( \frac{\sqrt{m} \mathbf{x}_\beta^{[l]}}{m} \right)^\top \\
 &\quad + \frac{c_{\text{res}}^2}{L^2\sqrt{m}} \sigma_{[l+1]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \sigma_{[l+1]}^{(1)}(\mathbf{x}_\beta) \left( E_{t,\beta}^{[(l+2):L]} \right)^\top \mathbf{a}_t \right) \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \langle \mathbf{x}_\alpha^{[l]}, \mathbf{x}_\beta^{[l]} \rangle \\
 &\quad + \sum_{k=2}^l \frac{c_{\text{res}}^3}{L^3\sqrt{m}} \sigma_{[l+1]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} E_{t,\alpha}^{[(k+1):l]} \sigma_{[k]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[k]}^{(1)}(\mathbf{x}_\beta) \left( E_{t,\beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right) \\
 &\quad \quad \quad \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \langle \mathbf{x}_\alpha^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle \\
 &\quad + \frac{c_\sigma}{\sqrt{m}} \frac{c_{\text{res}}}{L} \sigma_{[l+1]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} E_{t,\alpha}^{[2:l]} \sigma_{[1]}^{(1)}(\mathbf{x}_\alpha) \sigma_{[1]}^{(1)}(\mathbf{x}_\beta) \left( E_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \\
 &\quad \quad \quad \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle.
 \end{aligned}$$

We notice that the constant  $\frac{c_{\text{res}}}{L}$  plays an important part in our proof, so that the width per layer  $m$  does not depend exponentially in depth  $L$ .

Using the above rules, the derivative for NTK  $\mathcal{K}_t^{(2)}(\cdot)$  is obtained in the following form

$$\partial_t \mathcal{K}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{K}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta) (f_\beta(t) - y_\beta),$$

where each term in  $\mathcal{K}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta)$  is the summation of all the terms generated from  $\mathcal{K}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$  by performing the replacement procedure. In order to illustrate the idea, we give out an example in the proof of Theorem 4.2 in Section 5.4.

By the same reasoning, we could obtain the higher order kernels inductively by performing all the possible replacements. For instance, for kernel  $\mathcal{K}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ , we could obtain  $\mathcal{K}_t^{(r+1)}(\mathbf{x}_{\alpha_1}, \dots, \mathbf{x}_{\alpha_r}, \mathbf{x}_\beta)$  given by the following Ordinary Differential Equation

$$\partial_t \mathcal{K}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r}) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{K}_t^{(r+1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r}, \mathbf{x}_\beta) (f_\beta(t) - y_\beta).$$

In order to describe the vectors appearing in  $\mathcal{K}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ , we need to introduce some systematic notations.

## 5.2 Hierarchical sets of kernel expressions

The hierarchy of sets are proposed originally by Huang and Yau in [24]. We denote  $\mathbb{A}_0$  the first set of expressions in the following form, which corresponds to the terms in  $\mathcal{K}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$ . We define  $\mathbb{A}_0$  as:

$$\mathbb{A}_0 \triangleq \{\mathbf{e}_s \mathbf{e}_{s-1} \cdots \mathbf{e}_1 \mathbf{e}_0 : 0 \leq s \leq 4L\}, \quad (5.3)$$

where  $\mathbf{e}_j$  is chosen following the rules:

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \{\sqrt{m}\mathbf{x}_\beta^{[1]}, \sqrt{m}\mathbf{x}_\beta^{[2]}, \dots, \sqrt{m}\mathbf{x}_\beta^{[L]}\}_{1 \leq \beta \leq n} \right\}, \quad (5.4)$$

and for  $1 \leq j \leq s$ ,

$$\mathbf{e}_j \in \left\{ \left\{ \mathbf{E}_{t,\beta'}^{[2]} \left( \mathbf{E}_{t,\beta}^{[2]} \right)^\top, \dots, \mathbf{E}_{t,\beta'}^{[L]} \left( \mathbf{E}_{t,\beta}^{[L]} \right)^\top \right\}_{1 \leq \beta \leq n}, \left\{ \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta), \dots, \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_\beta) \right\}_{1 \leq \beta \leq n} \right\}. \quad (5.5)$$

Hence, each term in  $\mathcal{K}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$  writes as

$$\frac{\langle \mathbf{v}_1(t), \mathbf{v}_2(t) \rangle}{m} \quad \text{or} \quad \frac{\langle \mathbf{v}_1(t), \mathbf{v}_2(t) \rangle}{m} \frac{\langle \mathbf{v}_3(t), \mathbf{v}_4(t) \rangle}{m},$$

where  $\mathbf{v}_1(t), \mathbf{v}_2(t), \mathbf{v}_3(t), \mathbf{v}_4(t) \in \mathbb{A}_0$ . Note that  $\mathbf{v}_i(t)$  can take the value of  $\mathbf{v}_i(t) = \sqrt{m}\mathbf{x}_\alpha$ , which are not contained in  $\mathbb{A}_0$ . However, such a singularity is not a big issue—see Appendix A.2. We remark that compared with [24],  $\mathbf{e}_j$  is chosen in a way different from ours, the counterpart in [24] is chosen from the set

$$\left\{ \left\{ \frac{\mathbf{W}_t^{[2]}}{\sqrt{m}}, \left( \frac{\mathbf{W}_t^{[2]}}{\sqrt{m}} \right)^\top, \dots, \frac{\mathbf{W}_t^{[L]}}{\sqrt{m}}, \left( \frac{\mathbf{W}_t^{[L]}}{\sqrt{m}} \right)^\top \right\}_{1 \leq \beta \leq n}, \left\{ \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta), \dots, \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_\beta) \right\}_{1 \leq \beta \leq n} \right\}.$$

Such changes arise from the change of the network structure, and it has been shown in Appendix A.2 that the group of skip-connection matrices  $\mathbf{E}_{t,\beta}^{[l]}$  possesses more stability than  $\frac{\mathbf{W}_t^{[l]}}{\sqrt{m}}$ .

Moreover, given the construction of  $\mathbb{A}_0, \mathbb{A}_1, \dots, \mathbb{A}_r$ , we denote  $\mathbb{A}_{r+1}$  to be the set of expressions in the following form:

$$\mathbb{A}_{r+1} \triangleq \{\mathbf{e}_s \mathbf{e}_{s-1} \cdots \mathbf{e}_1 \mathbf{e}_0 : 0 \leq s \leq 4L\},$$

where  $\mathbf{e}_j$  is chosen from the following sets:

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \mathbf{1}, \{\sqrt{m}\mathbf{x}_\beta^{[1]}, \sqrt{m}\mathbf{x}_\beta^{[2]}, \dots, \sqrt{m}\mathbf{x}_\beta^{[L]}\}_{1 \leq \beta \leq n} \right\}, \quad (5.6)$$

and for  $1 \leq j \leq s$ , we have that each  $e_j$  comes from one of the three following sets

$$\left\{ \left\{ \mathbf{E}_{t,\beta}^{[2]}, \left(\mathbf{E}_{t,\beta}^{[2]}\right)^\top, \dots, \mathbf{E}_{t,\beta}^{[L]}, \left(\mathbf{E}_{t,\beta}^{[L]}\right)^\top \right\}_{1 \leq \beta \leq n}, \left\{ \sigma_{[1]}^{(1)}(\mathbf{x}_\beta), \dots, \sigma_{[L]}^{(1)}(\mathbf{x}_\beta) \right\}_{1 \leq \beta \leq n} \right\},$$

$$\left\{ \text{diag}(\mathbf{g}), \mathbf{g} \in \mathbb{A}_0 \cup \mathbb{A}_1 \cup \dots \cup \mathbb{A}_r \right\},$$

$$\left\{ \sigma_{[l]}^{(u+1)}(\mathbf{x}_\beta) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \left( \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_{u+1}}, \right.$$

$$\left. \left( \frac{c_{\text{res}}}{L} \frac{\left(\mathbf{W}_t^{[l]}\right)^\top}{\sqrt{m}} \right)^{Q_{u+1}} \sigma_{[l]}^{(u+1)}(\mathbf{x}_\beta) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) : 2 \leq l \leq L, \right.$$

$$\left. 1 \leq \beta \leq n, 1 \leq u \leq r, \mathbf{g}_1, \mathbf{g}_2 \cdots \mathbf{g}_u \in \mathbb{A}_0 \cup \mathbb{A}_1 \cup \dots \cup \mathbb{A}_r \text{ and } Q_1, Q_2 \cdots Q_{u+1} \in \{0, 1\} \right\}.$$

The maximum possible total number of diag operations for any element in  $\mathbb{A}_r$  is  $r$ , i.e., if  $\mathbf{v}(t) \in \mathbb{A}_r$ , it contains at most  $r$  diag operations. We observe from the replacement rules that there will be a scaling of  $\frac{1}{\sqrt{m}}$  whenever we take derivatives. Hence inductively, each term in kernel  $\mathcal{K}_t^{(p)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_p})$  takes the form

$$\frac{1}{m^{p/2-1}} \prod_{j=1}^s \frac{\langle \mathbf{v}_{2j-1}(t), \mathbf{v}_{2j}(t) \rangle}{m}, \quad 1 \leq s \leq p, \quad \mathbf{v}_i(t) \in \mathbb{A}_0 \cup \mathbb{A}_1 \cup \dots \cup \mathbb{A}_{p-2}, \quad (5.7)$$

which is a direct consequence of Proposition A.1. Note that  $\mathbf{v}_i(t)$  can still take the value of  $\mathbf{v}_i(t) = \sqrt{m} \mathbf{x}_\alpha$ , which are not in the set  $\mathbb{A}_{r+1}$ . Huang and Yau also obtained (5.7) in equation (3.8) in [24], and they use the tensor program proposed by Yang [47] to estimate the initial value of the kernel  $\mathcal{K}_0^{(p)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_p})$ . They showed that for each vector  $\mathbf{v}_j(t)$  in (5.7) at  $t=0$ , it is a linear combination of projections of independent Gaussian vectors. Hence, if we consider such quantity

$$\eta(t) = \{ \|\mathbf{v}(t)\|_\infty : \mathbf{v}(t) \in \mathbb{A}_0 \cup \mathbb{A}_1 \cdots \cup \mathbb{A}_r \},$$

at  $t=0$ , since  $\mathbf{v}(0)$  is a linear combination of projections of independent Gaussian vectors, then with high probability,  $\eta(0) \lesssim (\ln m)^C$ . For  $t > 0$ , Huang and Yau derived a self-consistent Ordinary Differential Inequality for  $\eta(t)$  leading to the following estimates:

$$\partial_t^{(p+1)} \eta(t) \lesssim \frac{\eta(t)^{2p}}{m^{p/2}}, \quad (5.8)$$

$$\eta(0) \lesssim (\ln m)^C, \quad (5.9)$$

then it holds that  $\eta(t) \lesssim (\ln m)^C$  for time  $0 \leq t \leq m^{\frac{p}{2p+1}} / (\ln m)^C$ .



Our approach is different from theirs. Instead of using tensor programs, we use a special matrix norm, the 2 to infinity matrix norm, to show that  $\eta(0) \lesssim (\ln m)^C$ . We can then obtain a Gronwall-type inequality for  $\eta(t)$ :

$$\eta(t) \lesssim (\ln m)^C + \frac{1}{\sqrt{m}} \int_0^t \eta(s) ds,$$

for which it follows that for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ , the estimate  $\eta(t) \lesssim (\ln m)^C$  holds. Then (4.6) and (4.7) in Theorem 4.1 holds, and we are able to show that the kernels of higher order vary slowly, which brings us the proof of Theorem 4.2.

### 5.3 Least eigenvalue for randomly initialized matrix

Firstly, since  $\mathbf{K}^{[L]}$  is a recursively defined matrix, we use the results in Du et al. [12] to show that the Gram matrix  $\mathbf{K}^{[L]}$  is positive definite. Secondly, we need to analyze how the difference between  $\mathbf{G}^{[1]}$  and  $\mathbf{K}^{[1]}$ , termed the perturbation by Du et al. [12], propagates from the lower layers to the  $L$ -th layer. We quantitatively characterize how large such propagation dynamics would be and rediscover that ResNet architecture serves as a stabilizer for such propagation (Proposition C.1). Our proof is slightly different from [12], where we use the concentration inequality for Lipschitz functions. We refer readers to Appendix C for details.

### 5.4 Sketch of proof

We use Fig. 1 to illustrate the ideas of the proofs. Due to space constraints, all the proofs of the technical Lemmas and Propositions are provided in Supplementary Material. Note that for the quantities in Fig. 1,  $\lambda_0$  is the least eigenvalue of  $\tilde{\mathbf{K}}_{ij}^{[1]}$ ,  $\xi_{\infty,r}(t) = \sup_{0 \leq t' \leq t} \{\|v(t')\|_2 : v(t') \in \mathbb{A}_r\}$ , and  $\eta_{\infty,r}(t) = \sup_{0 \leq t' \leq t} \{\|v(t')\|_\infty : v(t') \in \mathbb{A}_r\}$ , where  $r \geq 0$ . Now we proceed to the Proof of Theorem 4.1.

*Proof of Theorem 4.1.* Since each term in kernel  $\mathcal{K}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ , it takes the form

$$\frac{1}{m^{r/2-1}} \prod_{j=1}^s \frac{\langle v_{2j-1}(t), v_{2j}(t) \rangle}{m}, \quad 1 \leq s \leq r, \quad v_i(t) \in \mathbb{A}_0 \cup \mathbb{A}_1 \cup \dots \cup \mathbb{A}_{r-2},$$

then for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ ,  $r \geq 3$

$$\|\mathcal{K}_t^{(2)}(\cdot)\|_\infty \lesssim \left(\frac{\xi_{\infty,0}(t)^2}{m}\right)^2 \lesssim 1,$$

$$\|\mathcal{K}_t^{(r)}(\cdot)\|_\infty \lesssim \frac{1}{m^{r/2-1}} \left(\frac{\xi_{\infty,r}(t)^2}{m}\right)^s \lesssim \frac{1}{m^{r/2-1}} \left(\frac{(c(\ln m)^C \sqrt{m})^2}{m}\right)^r \lesssim \frac{(\ln m)^{2rC}}{m^{r/2-1}}.$$

This completes the proof. □

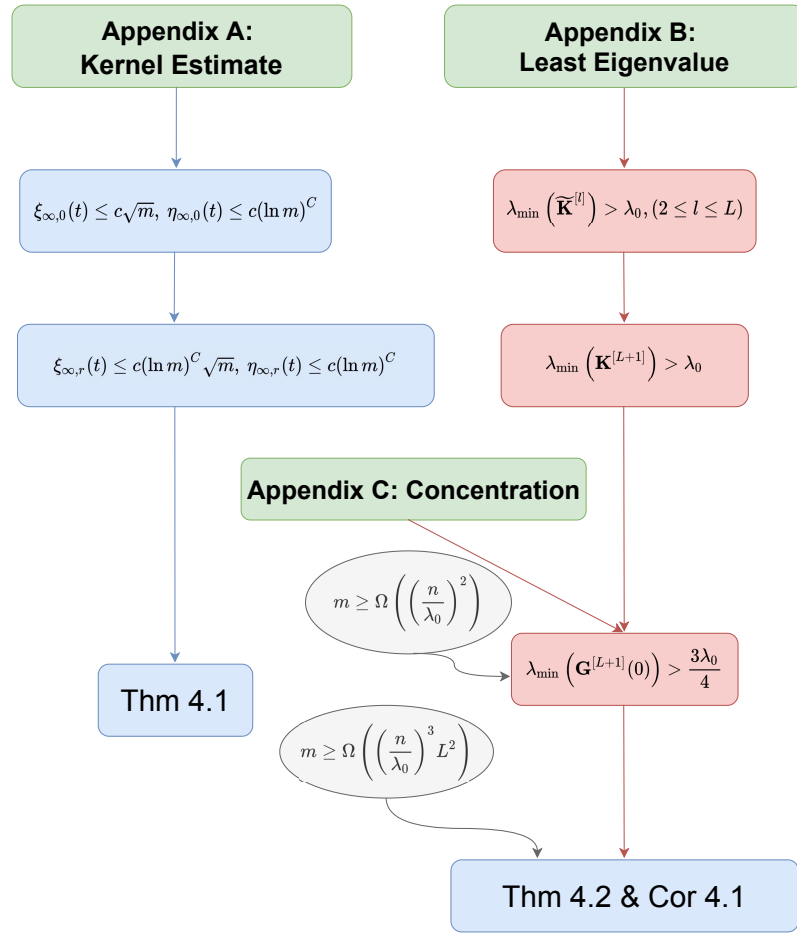


Figure 1: Flow chart for the proof of main theorems.

Now we sketch the proofs for Theorem 4.2 and Corollary 4.1. Details can be found in Appendix D.

*Sketch of the proof of Theorem 4.2.* Since there exists  $\frac{1}{L^2}$  scaling in some kernels, we use  $C(r,L)$  to denote the ‘effective terms’ in each kernel. We denote  $\mathcal{G}_t^{[L+1]}(\cdot)$  by  $\mathcal{G}_t^{(2)}(\cdot)$ , i.e.,  $\mathcal{G}_t^{(2)}(\cdot) := \mathcal{G}_t^{[L+1]}(\cdot)$ , it’s natural for us to get that  $C(2,L) = \mathcal{O}(1)$ .

Next, we apply the replacement rule, all the possible terms generated from  $\mathcal{G}_t^{(2)}(\cdot)$  are

$$\mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) = \langle \mathbf{x}_{\alpha_1}^{[L]}, \mathbf{x}_{\alpha_2}^{[L]} \rangle \rightarrow \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta),$$

$$\begin{aligned}
 & \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta) \\
 &= \frac{c_\sigma}{m} \left\langle \text{diag} \left( \mathbf{E}_{t, \alpha_1}^{[2:L]} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_{\alpha_1}) \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1}, \mathbf{x}_{\alpha_2}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_1}, \mathbf{x}_\beta \rangle \quad \text{(I)} \\
 &+ \sum_{k=2}^L \frac{c_{\text{res}}^2}{L^2 m} \left\langle \text{diag} \left( \mathbf{E}_{t, \alpha_1}^{[(k+1):L]} \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_{\alpha_1}) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1}, \mathbf{x}_{\alpha_2}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_1}^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle \quad \text{(II)} \\
 &+ \frac{c_\sigma}{m} \left\langle \text{diag} \left( \mathbf{E}_{t, \alpha_2}^{[2:L]} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_{\alpha_2}) \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1}, \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta \rangle \\
 &+ \sum_{k=2}^L \frac{c_{\text{res}}^2}{L^2 m} \left\langle \text{diag} \left( \mathbf{E}_{t, \alpha_2}^{[(k+1):L]} \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_{\alpha_2}) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1}, \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_2}^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle,
 \end{aligned}$$

and we have  $C(3, L) = \mathcal{O} \left( 2 \left( 1 + \frac{L-1}{L^2} \right) \right) = \mathcal{O} \left( 1 + \frac{1}{L} \right)$ .

Finally for  $\mathcal{G}_t^{(4)}(\cdot)$ , by symmetry, we are only going to analyze terms I and II above. Since there are at most  $(2L+2)$  symbols in term I to be replaced, each replacement will bring about up to  $(L+1)$  many terms. For term II, for each summand, there are also at most  $(2L+2)$  symbols to be replaced. Since there are  $L-1$  summands in II and each replacement will bring about up to  $(L+1)$  many terms, we have that

$$C(4, L) = \mathcal{O} \left( 2 \left( (2L+2)(L+1) + \frac{1}{L^2} (L-1)(2L+2)(L+1) \right) \right) = \mathcal{O}(L^2).$$

It holds that for time  $0 \leq t \leq \sqrt{m} / (\ln m)^{C'}$

$$\begin{aligned}
 \left| \partial_t \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| &\leq \left\| \mathcal{G}_t^{(4)}(\cdot) \right\|_\infty \sqrt{R_S(\boldsymbol{\theta}_0)} \leq C(4, L) \frac{(\ln m)^C}{m}, \\
 \left| \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| &\leq \left\| \mathcal{G}_0^{(3)}(\cdot) \right\|_\infty + t C(4, L) \frac{(\ln m)^C}{m}.
 \end{aligned}$$

Finally, we need to estimate  $\left\| \mathcal{G}_0^{(3)}(\cdot) \right\|_\infty$ . Each term in  $\mathcal{G}_0^{(3)}(\cdot)$  is of the form  $\frac{c}{m} \langle \mathbf{B} \mathbf{a}_0, \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle$  where  $\mathbf{B}$  is some specific matrix that changes from term to term. After taking conditional expectation with respect to the random variable  $\mathbf{a}_0$ , we have with high probability that

$$\frac{c}{m} \left\langle \mathbf{a}_0, \mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle \leq c \frac{(\ln m)^C}{m}. \tag{5.10}$$

Consequently, for time  $0 \leq t \leq \sqrt{m} / (\ln m)^{C'}$ ,

$$\begin{aligned}
 \left| \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| &\leq C(3, L) \frac{(\ln m)^C}{m} + t C(4, L) \frac{(\ln m)^C}{m}, \\
 \left| \partial_t \mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) \right| &\leq \left\| \mathcal{G}_t^{(3)}(\cdot) \right\|_\infty \sqrt{R_S(\boldsymbol{\theta}_0)} \leq (C(3, L) + t C(4, L)) \frac{(\ln m)^C}{m},
 \end{aligned}$$

which finishes the proof of Theorem 4.2. □

*Sketch of the proof of Corollary 4.1.* If  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\varepsilon}\right)$ , with high probability with respect to random initialization,  $\lambda_{\min}[\mathcal{K}_0^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n} > \lambda_{\min}(\mathbf{G}^{[L+1]}(0)) > \frac{3\lambda_0}{4}$ . We finish the proof of (4.12) by setting  $\lambda = \frac{3\lambda_0}{4}$ .

Concerning the change of the least eigenvalue of the NTK, from the Sketch Proof of Theorem 4.2, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ ,

$$\begin{aligned} \left\| \left( \mathcal{G}_t^{(2)} - \mathcal{G}_0^{(2)} \right) (\cdot) \right\|_{2 \rightarrow 2} &\leq \left\| \left( \mathcal{G}_t^{(2)} - \mathcal{G}_0^{(2)} \right) (\cdot) \right\|_{\mathbb{F}} \\ &\leq n \left\| \left( \mathcal{G}_t^{(2)} - \mathcal{G}_0^{(2)} \right) (\cdot) \right\|_{\infty} \leq nt(C(3,L) + tC(4,L)) \frac{(\ln m)^C}{m}. \end{aligned}$$

Let  $t^*$  satisfy

$$C(4,L)(t^*)^2 + C(3,L)t^* = \frac{\lambda m}{2(\ln m)^{C_n}}, \tag{5.11}$$

which gives

$$t^* = \frac{-C(3,L) + \sqrt{(C(3,L))^2 + 2C(4,L) \frac{\lambda m}{(\ln m)^{C_n}}}}{2C(4,L)} \geq \frac{1}{2} \sqrt{\frac{\lambda m}{C(4,L)(\ln m)^{C_n}}}.$$

Next, let  $\bar{t} := \inf \{t : \lambda_{\min}[\mathcal{K}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n} \geq \lambda/2\}$ . Naturally, we have  $t^* \leq \bar{t}$ . Using (4.6), we have that for any  $0 \leq t \leq \bar{t}$ ,  $R_S(\boldsymbol{\theta}_t) \leq \exp(-\lambda t/n) R_S(\boldsymbol{\theta}_0)$ .

Set  $R_S(\boldsymbol{\theta}_t) = \varepsilon$ . It takes time  $t \leq (n/\lambda) \ln(R_S(\boldsymbol{\theta}_0)/\varepsilon)$  for the loss  $R_S(\boldsymbol{\theta}_t)$  to reach accuracy  $\varepsilon$ . Hence, if  $t \leq (n/\lambda) \ln(R_S(\boldsymbol{\theta}_0)/\varepsilon) \leq t^* \leq \bar{t}$ , then width  $m$  is required to be

$$\frac{n}{\lambda} \ln \left( \frac{R_S(\boldsymbol{\theta}_0)}{\varepsilon} \right) \leq \frac{1}{2} \sqrt{\frac{\lambda m}{C(4,L)(\ln m)^{C_n}}}. \tag{5.12}$$

Thus we have

$$m \geq C(4,L) \left( \frac{n}{\lambda} \right)^3 (\ln m)^C \ln \left( \frac{R_S(\boldsymbol{\theta}_0)}{\varepsilon} \right)^2.$$

Since  $C(4,L) = \mathcal{O}(L^2)$ , we conclude the proof. □

## 6 Discussion

In this paper, we show that the GD on ResNet can obtain zero training loss and its training dynamic is given by an infinite hierarchy of ordinary differential equations, i.e., the NTH, which makes it possible to study the change of the NTK directly for deep neural

networks. Our proof builds on a careful analysis of the least eigenvalue of randomly initialized Gram matrix, and the uniform upper bound on kernels of higher order in the NTH.

We list some future directions for research:

- The NTH is an infinite sequence of relationship. However, Huang and Yau [24] showed that under certain conditions on the width and the data set dimension, the NTH can be truncated and the truncated version of NTH is still able to approximate the original dynamic up to any precision. We believe that for ResNet, such technical conditions can be loosened based on our result.
- In Corollary 4.1, the dependence of  $m$  on the depth  $L$  is quadratic, we believe that the dependence can be reduced even further. We conjecture that  $m$  is independent of  $L$ .
- In this paper, we focus on the GD, and we believe that it can be extended to SGD, while maintaining the linear convergence rate.
- We focus here on the training loss, but does not address the test loss. To further investigate the generalization power of ResNet, we believe that some apriori estimates for the generalization error of ResNet may be useful [34, 35].

## Acknowledgments

This work is supported by National Natural Science Foundation of China Grant No. 12101401 (T.L.), Shanghai Municipal Science and Technology Key Project No. 22JC1401500 (T.L.), and Shanghai Municipal of Science and Technology Major Project No. 2021SHZDZX0102 (T.L.).

## Appendices

### A Estimates on the kernel

#### A.1 Structure on hierarchical sets of kernel expressions

Since we have mentioned the replacement rules in Section 5.1, we haven't rigorously justified it yet. Hence we use Proposition A.1 to shed light on the structures of the elements in  $\mathbb{A}_r$ , and consequently on the structures of each term in kernel  $\mathcal{K}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ .

**Proposition A.1.** For any vector  $v(t) \in \mathbb{A}_r$ , the new vector obtained from  $v(t)$  by performing the replacement rules are the sum of terms of the following forms:

- (a).  $\frac{C}{\sqrt{m}}v'(t) : v'(t) \in \mathbb{A}_r,$
- (b).  $\frac{C}{\sqrt{m}}v'(t) \frac{\langle p, q \rangle}{m} : v'(t) \in \mathbb{A}_{r+1}, p, q \in \mathbb{A}_0,$
- (c).  $\frac{C}{\sqrt{m}}v'(t) \frac{\langle \sqrt{m}x_\alpha, \sqrt{m}x_\beta \rangle}{m} : v'(t) \in \mathbb{A}_{r+1}, p, q \in \mathbb{A}_0,$
- (d).  $\frac{C}{\sqrt{m}}v'(t) \frac{\langle p, q \rangle}{m} : v'(t) \in \mathbb{A}_{r-s+1}, p \in \mathbb{A}_s, q \in \mathbb{A}_0, \text{ for some } s \geq 1,$
- (e).  $\frac{C}{\sqrt{m}}v'(t) \frac{\langle p, q \rangle}{m} : v'(t) \in \mathbb{A}_s, p \in \mathbb{A}_{r-s+1}, q \in \mathbb{A}_0, \text{ for some } s \geq 1.$

*Proof.* The proof comes as follows. Note that the constant  $C$  listed out below might keep changing from term to term.

Since  $a_t$  appears only at the position  $e_0$ , if  $v(t) \in \mathbb{A}_r$ , based on the replacement rule

$$v(t) = e_s e_{s-1} \cdots e_1 a_t \rightarrow \tilde{v}(t) = \frac{1}{\sqrt{m}} e_s e_{s-1} \cdots e_1 \sqrt{m} x_\beta^{[L]} = \frac{1}{\sqrt{m}} v'(t),$$

then  $v'(t) = e_s e_{s-1} \cdots e_1 \sqrt{m} x_\beta^{[L]} \in \mathbb{A}_r$ .

Similarly,  $\sqrt{m} x_\alpha^{[L]}$  also appears only at  $e_0$ , then if  $v(t) \in \mathbb{A}_r$ , by the replacement rule

$$v(t) = e_s e_{s-1} \cdots e_1 \sqrt{m} x_\alpha^{[L]} \rightarrow \tilde{v}(t),$$

$$\tilde{v}(t) = \sum_k \frac{C}{\sqrt{m}} e_s e_{s-1} \cdots e_1 \text{diag}(f_k) \mathbf{1} \frac{\langle \sqrt{m} x_\alpha^{[k]}, \sqrt{m} x_\beta^{[k]} \rangle}{m} = \sum_k \frac{C}{\sqrt{m}} v'_k(t) \frac{\langle \sqrt{m} x_\alpha^{[k]}, \sqrt{m} x_\beta^{[k]} \rangle}{m},$$

given that  $f_k \in \mathbb{A}_0$ , then  $v'_k(t) \in \mathbb{A}_{r+1}$ .

Since  $\sigma_{[l]}^{(u)}(x_\alpha)$  only appears at the starting or the middle position, i.e.,  $e_j, j \geq 1$ . For  $u=1, \sigma_{[l]}^{(1)}(x_\alpha)$  has no diag operations accompanied with it, and any vector  $v(t) \in \mathbb{A}_r$  could contain  $\sigma_{[l]}^{(1)}(x_\alpha)$ , for  $r \geq 0$ ,

$$\begin{aligned} v(t) &= e_s \cdots e_{j+1} \sigma_{[l]}^{(1)}(x_\alpha) e_{j-1} \cdots e_0 \rightarrow \tilde{v}(t), \\ \tilde{v}(t) &= \frac{C}{\sqrt{m}} e_s \cdots e_{j+1} \sigma_{[l]}^{(2)}(x_\alpha) \text{diag}(f_1) e_{j-1} \cdots e_0 \frac{\langle p_1, q_1 \rangle}{m} \\ &\quad + \sum_k \frac{C}{\sqrt{m}} e_s \cdots e_{j+1} \sigma_{[l]}^{(2)}(x_\alpha) \text{diag}\left(\frac{W_t^{[l]}}{\sqrt{m}} f_k\right) e_{j-1} \cdots e_0 \frac{\langle p_k, q_k \rangle}{m} \\ &= \sum_l \frac{C}{\sqrt{m}} v'_l(t) \frac{\langle p_l, q_l \rangle}{m}, \end{aligned}$$

since  $f_k \in \mathbb{A}_0$ , then  $v'_l(t) \in \mathbb{A}_{r+1}$ , and  $p_l, q_l \in \mathbb{A}_0$ .

For  $u \neq 1$ ,  $\sigma_{[l]}^{(u)}(x_\alpha)$  has at most  $u-1$  diag operations behind it, and only vector  $v(t) \in \mathbb{A}_r$  could contain  $\sigma_{[l]}^{(u)}(x_\alpha)$ , for  $r \geq u-1$

$$v(t) = e_s \cdots e_{j+1} e_j e_{j-1} \cdots e_0 \rightarrow \tilde{v}(t),$$

with

$$e_j = \sigma_{[l]}^{(u)}(x_\alpha) \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_1} g_1 \right) \cdots \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_{u-1}} g_{u-1} \right) \left( \frac{c_{\text{res}}}{L} \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_u},$$

or  $e_j = \left( \frac{c_{\text{res}}}{L} \frac{(W_t^{[l]})^\top}{\sqrt{m}} \right)^{Q_u} \sigma_{[l]}^{(u)}(x_\alpha) \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_1} g_1 \right) \cdots \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_{u-1}} g_{u-1} \right),$

and after applying replacement rules on  $e_j \rightarrow e'_j$ ,

$$e'_j = \frac{C}{\sqrt{m}} \sigma_{[l]}^{(u+1)}(x_\alpha) \text{diag}(f_1) \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_1} g_1 \right) \cdots \left( \frac{c_{\text{res}}}{L} \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \frac{\langle p_1, q_1 \rangle}{m}$$

$$+ \sum_k \frac{C}{\sqrt{m}} \sigma_{[l]}^{(u+1)}(x_\alpha) \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_0} f_k \right) \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_1} g_1 \right)$$

$$\cdots \left( \frac{c_{\text{res}}}{L} \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \frac{\langle p_k, q_k \rangle}{m},$$

or  $e'_j = \frac{C}{\sqrt{m}} \left( \frac{c_{\text{res}}}{L} \frac{(W_t^{[l]})^\top}{\sqrt{m}} \right)^{Q_u} \sigma_{[l]}^{(u+1)}(x_\alpha) \text{diag}(f_1) \cdots \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_{u-1}} g_{u-1} \right) \frac{\langle p_1, q_1 \rangle}{m}$ 

$$+ \sum_k \frac{C}{\sqrt{m}} \left( \frac{c_{\text{res}}}{L} \frac{(W_t^{[l]})^\top}{\sqrt{m}} \right)^{Q_u} \sigma_{[l]}^{(u+1)}(x_\alpha) \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_0} f_k \right)$$

$$\cdots \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_{u-1}} g_{u-1} \right) \frac{\langle p_k, q_k \rangle}{m},$$

then

$$\tilde{v}(t) = \sum_l \frac{C}{\sqrt{m}} v'_l(t) \frac{\langle p_l, q_l \rangle}{m},$$

since  $f_k \in \mathbb{A}_0$ , then  $v'_l(t) \in \mathbb{A}_{r+1}$ , and  $p_l, q_l \in \mathbb{A}_0$ .

Since  $\frac{W_t^{[l]}}{\sqrt{m}}$  only appears at the starting or the middle position  $e_j$ , we have that if  $v(t) \in \mathbb{A}_r$ , then based on the replacement rules

$$\begin{aligned} v(t) &= e_s e_{s-1} \cdots e_{j+1} \frac{W_t^{[l]}}{\sqrt{m}} e_{j-1} \cdots e_1 e_0 \rightarrow \tilde{v}(t), \\ \tilde{v}(t) &= \frac{C}{m} e_s e_{s-1} \cdots e_{j+1} \text{diag}(g) \mathbf{1} \otimes (\mathbf{x}_\beta^{[l-1]})^\top e_{j-1} \cdots e_1 e_0 \\ &= \frac{C}{\sqrt{m}} e_s e_{s-1} \cdots e_{j+1} \text{diag}(g) \mathbf{1} \frac{\langle e_{j-1} \cdots e_1 e_0, \sqrt{m} \mathbf{x}_\beta^{[l-1]} \rangle}{m} \\ &= \frac{C}{\sqrt{m}} v'(t) \frac{\langle p, q \rangle}{m}, \end{aligned}$$

with  $v'(t) \in \mathbb{A}_{r-s+1}$ , and  $p \in \mathbb{A}_s, q \in \mathbb{A}_0$ , for some  $s \geq 1$ .

Similarly for  $\frac{(W_t^{[l]})^\top}{\sqrt{m}}$ ,

$$\begin{aligned} v(t) &= e_s e_{s-1} \cdots e_{j+1} \frac{(W_t^{[l]})^\top}{\sqrt{m}} e_{j-1} \cdots e_1 e_0 \rightarrow \tilde{v}(t), \\ \tilde{v}(t) &= \frac{C}{m} e_s e_{s-1} \cdots e_{j+1} \mathbf{x}_\beta^{[l-1]} \otimes \mathbf{1}^\top \text{diag}(g) e_{j-1} \cdots e_1 e_0 \\ &= \frac{C}{\sqrt{m}} e_s e_{s-1} \cdots e_{j+1} \sqrt{m} \mathbf{x}_\beta^{[l-1]} \frac{\langle \text{diag}(g) e_{j-1} \cdots e_1 e_0, \mathbf{1} \rangle}{m} \\ &= \frac{C}{\sqrt{m}} v'(t) \frac{\langle p, q \rangle}{m}, \end{aligned}$$

with  $v'(t) \in \mathbb{A}_{r-s}$ , and  $p \in \mathbb{A}_{r-s+1}, q \in \mathbb{A}_0$ , for some  $s \geq 1$ . Since  $E_{t,\alpha}^{[l]}$  is situations combined with  $\frac{W_t^{[l]}}{\sqrt{m}}$  and  $\sigma_{[l]}^{(1)}(\mathbf{x}_\alpha)$ , so we will skip the analysis.  $\square$

From the discussion above, if we apply Proposition A.1 to  $\mathcal{K}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$  inductively  $(r-1)$  times, for each term in kernel  $\mathcal{K}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ , it takes the form:

$$\frac{1}{m^{r/2-1}} \prod_{j=1}^s \frac{\langle v_{2j-1}(t), v_{2j}(t) \rangle}{m}, \quad 1 \leq s \leq r, \quad v_i(t) \in \mathbb{A}_0 \cup \mathbb{A}_1 \cup \dots \cup \mathbb{A}_{r-2}. \quad (\text{A.1})$$

## A.2 Apriori $L^2$ bounds for expressions in $\mathbb{A}_0$

We begin with an estimate on the empirical risk  $R_S(\theta_t)$ .

**Proposition A.2.** Under Assumptions 4.1 and 4.2, we have for  $t \geq 0$ ,

$$R_S(\theta_t) \leq R_S(\theta_0) \sim \mathcal{O}(1). \quad (\text{A.2})$$



*Proof.* We get inequality (A.2) by non-negative definiteness of kernel  $\mathcal{K}_t^{(2)}(\cdot)$ . From (4.6), we obtain that

$$\partial_t \sum_{\alpha=1}^n \|f_\alpha(t) - y_\alpha\|_2^2 = -\frac{2}{n} \sum_{\alpha,\beta=1}^n \mathcal{K}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) (f_\alpha(t) - y_\alpha)(f_\beta(t) - y_\beta) \leq 0, \quad (\text{A.3})$$

hence

$$R_S(\boldsymbol{\theta}_t) \leq R_S(\boldsymbol{\theta}_0),$$

which finish the proof of Proposition A.2.  $\square$

Our next proposition is mainly on the spectral property of the skip-connection matrices. This proposition is similar to Proposition B.1. in [24].

**Proposition A.3.** *Under Assumptions 4.1 and 4.2, we define  $\zeta(t)$  as follows*

$$\zeta(t) = \sup_{0 \leq t' \leq t} \max \left\{ 1, \frac{1}{\sqrt{m}} \left\{ \left\| \mathbf{W}_{t'}^{[2]} \right\|_{2 \rightarrow 2'}, \left\| \left( \mathbf{W}_{t'}^{[2]} \right)^\top \right\|_{2 \rightarrow 2'}, \dots \right. \right. \\ \left. \left. \dots, \left\| \mathbf{W}_{t'}^{[L]} \right\|_{2 \rightarrow 2'}, \left\| \left( \mathbf{W}_{t'}^{[L]} \right)^\top \right\|_{2 \rightarrow 2'}, \left\| \mathbf{a}_{t'} \right\|_2 \right\} \right\}, \quad (\text{A.4})$$

then with high probability with respect to the random initialization, for  $t \lesssim \sqrt{m}$

$$\zeta(t) \leq c_{w,t}, \quad (\text{A.5})$$

where  $c_{w,t} > 2$  is a constant independent of the depth of the network  $L$ .

Moreover for  $t \lesssim \sqrt{m}$ ,  $c_{w,t}$  has a uniform upper bound in  $t$ , i.e.,

$$c_{w,t} \leq \bar{c}, \quad (\text{A.6})$$

where  $\bar{c}$  is independent of depth  $L$  and time  $t$ .

*Proof.* For the purpose of proving the proposition, we shall state two lemmas, Lemmas A.1 and A.2. Lemma A.1 is given out as Lemma G.2. in Du et al. [12], also consequence of the results in [46].

**Lemma A.1.** *Given a matrix  $\mathbf{W} \in \mathbb{R}^{m \times m}$  with each entry  $W_{i,j} \sim \mathcal{N}(0,1)$ , then with probability at least  $1 - \exp\left(-\frac{(\bar{c}'_{w,0}-2)^2 m}{2}\right)$ , the following holds*

$$\|\mathbf{W}\|_{2 \rightarrow 2} \leq \bar{c}'_{w,0} \sqrt{m}, \quad (\text{A.7})$$

where  $\bar{c}'_{w,0} > 2$  is a constant.

**Remark A.1.** This event is an event that holds with high probability.

Next concerning the term  $\frac{1}{\sqrt{m}}\|\mathbf{a}_0\|_2$ , we shall state a lemma on the tail bound of the chi-square distribution, using Lemma 1 from [29]

**Lemma A.2.** *If  $Z \sim \chi^2(m)$ , then we have a tail bound*

$$\mathbb{P}(Z \geq m + 2\sqrt{mx} + 2x) \leq e^{-x}. \tag{A.8}$$

**Remark A.2.** This event is also an event that holds with high probability.

Then if we write  $2tm = m + (2t - 1)m$ , letting  $x = \frac{mt}{10}$ , we can obtain that

$$\mathbb{P}\left(\|\mathbf{a}_0\|_2^2 \geq m + 2m\left(\sqrt{t/10} + t/10\right)\right) \leq \exp(-tm/10),$$

and for  $t \geq 1$ , we have  $2t - 1 \geq 2(\sqrt{t/10} + t/10)$ . Thus, if we choose  $t$  properly, we see that such event

$$\frac{1}{\sqrt{m}}\|\mathbf{a}_0\|_2 \leq c''_{w,0}$$

holds with high probability. Hence, for  $t = 0$ ,  $\xi(0) \leq \max\{1, c'_{w,0}, c''_{w,0}\}$ . We set  $c_{w,0}$  as  $c_{w,0} = \max\{1, c'_{w,0}, c''_{w,0}\}$ , then

$$\xi(0) \leq c_{w,0}. \tag{A.9}$$

In the following we are going to show the upper bound of  $\partial_t \xi(t)$ . In order to do that, we need to estimate  $L^2$  bound on each output layer. For  $l = 1$ ,

$$\begin{aligned} \|\mathbf{x}^{[1]}\|_2 &= \sqrt{\frac{c_\sigma}{m}} \|\sigma(\mathbf{W}_t^{[1]}\mathbf{x})\|_2 \leq \sqrt{c_\sigma} \left( |\sigma(0)| + \frac{C_{\text{Lip}}}{\sqrt{m}} \|\mathbf{W}_t^{[1]}\mathbf{x}\|_2 \right) \\ &\leq \sqrt{c_\sigma} C_{\text{Lip}} (1 + \xi(t) \|\mathbf{x}\|_2) \leq C \xi(t), \end{aligned} \tag{A.10}$$

and for  $2 \leq l \leq L$ ,

$$\begin{aligned} \|\mathbf{x}^{[l]}\|_2 &\leq \|\mathbf{x}^{[l-1]}\|_2 + \frac{c_{\text{res}}}{L\sqrt{m}} \|\sigma(\mathbf{W}_t^{[l]}\mathbf{x}^{[l-1]})\|_2 \\ &\leq \|\mathbf{x}^{[l-1]}\|_2 + \frac{c_{\text{res}}}{L} \left( |\sigma(0)| + C_{\text{Lip}} \xi(t) \|\mathbf{x}^{[l-1]}\|_2 \right) \\ &\leq \|\mathbf{x}^{[l-1]}\|_2 + \frac{c_{\text{res}}}{L} \left( C_{\text{Lip}} + C_{\text{Lip}} \xi(t) \|\mathbf{x}^{[l-1]}\|_2 \right) \\ &\leq \left( 1 + \frac{2c_{\text{res}}}{L} \xi(t) \right) \|\mathbf{x}^{[l-1]}\|_2. \end{aligned} \tag{A.11}$$

Hence we can obtain an inductive relation on the 2-norm of  $\mathbf{x}^{[l]}$ .

$$\|\mathbf{x}^{[l]}\|_2 \leq C \left( 1 + \frac{2c_{\text{res}}}{L} \xi(t) \right)^{l-1} \xi(t). \tag{A.12}$$

Based on (3.7), (3.8), (3.9) and (3.10), combined with Proposition A.2

$$\begin{aligned}
 \partial_t \|\mathbf{W}_t^{[l]}\|_{2 \rightarrow 2} &\leq \frac{1}{n} \sum_{\beta=1}^n \frac{C}{\sqrt{m}} \left\| \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+1):L]} \right)^\top \mathbf{a}_t \right\|_2 \left\| \mathbf{x}_\beta^{[l-1]} \right\|_2 |f_\beta(t) - y_\beta| \\
 &\leq \frac{1}{n} \sum_{\beta=1}^n C C_{\text{Lip}} \left( 1 + \frac{c_{\text{res}} C_{\text{Lip}}}{L} \zeta(t) \right)^{L-l} \zeta(t) \left( 1 + \frac{2c_{\text{res}}}{L} \zeta(t) \right)^{l-1} \zeta(t) |f_\beta(t) - y_\beta| \\
 &\leq C \left( 1 + \frac{2c_{\text{res}}}{L} \zeta(t) \right)^{L-1} \zeta(t)^2 \sqrt{\frac{1}{n} \sum_{\beta=1}^n \|f_\beta(t) - y_\beta\|_2^2} \\
 &\leq C \left( 1 + \frac{2c_{\text{res}}}{L} \zeta(t) \right)^{L-1} \zeta(t)^2 \sqrt{R_S(\boldsymbol{\theta}_0)} \\
 &\leq C \left( 1 + \frac{2c_{\text{res}}}{L} \zeta(t) \right)^{L-1} \zeta(t)^2 \leq C \exp(2c_{\text{res}} \zeta(t)) \zeta(t)^2, \tag{A.13}
 \end{aligned}$$

$$\begin{aligned}
 \partial_t \|\mathbf{a}_t\|_2 &\leq \frac{1}{n} \sum_{\beta=1}^n \left\| \mathbf{x}_\beta^{[L]} \right\|_2 |f_\beta(t) - y_\beta| \leq C \left( 1 + \frac{2c_{\text{res}}}{L} \zeta(t) \right)^{L-1} \zeta(t) \sqrt{\frac{1}{n} \sum_{\beta=1}^n \|f_\beta(t) - y_\beta\|_2^2} \\
 &\leq C \left( 1 + \frac{2c_{\text{res}}}{L} \zeta(t) \right)^{L-1} \zeta(t) \sqrt{R_S(\boldsymbol{\theta}_0)} \\
 &\leq C \left( 1 + \frac{2c_{\text{res}}}{L} \zeta(t) \right)^{L-1} \zeta(t) \leq C \exp(2c_{\text{res}} \zeta(t)) \zeta(t). \tag{A.14}
 \end{aligned}$$

Based on (A.13) and (A.14), we have

$$\sqrt{m} \partial_t \zeta(t) \leq C \exp(2c_{\text{res}} \zeta(t)) \zeta^2(t),$$

we can obtain an integration inequality,

$$\int_{\zeta(0)}^{\zeta(t)} \frac{du}{\exp(2c_{\text{res}} u) u^2} \leq \frac{Ct}{\sqrt{m}}. \tag{A.15}$$

Hence the integration term on the LHS of (A.15) is

$$\begin{aligned}
 \int_{\zeta(0)}^{\zeta(t)} \frac{du}{\exp(2c_{\text{res}} u) u^2} &\geq \frac{1}{\exp(2c_{\text{res}} \zeta(t))} \int_{\zeta(0)}^{\zeta(t)} \frac{du}{u^2} \\
 &= \frac{1}{\exp(2c_{\text{res}} \zeta(t))} \left( \frac{1}{\zeta(0)} - \frac{1}{\zeta(t)} \right) \\
 &\geq \frac{1}{\exp(2c_{\text{res}} \zeta(t))} \left( \frac{1}{c_{w,0}} - \frac{1}{\zeta(t)} \right).
 \end{aligned}$$

We shall notice for the single variable function  $f(z)$

$$f(z) = \frac{1}{\exp(2c_{\text{res}} z)} \left( \frac{1}{c_{w,0}} - \frac{1}{z} \right),$$

maximum of  $f(z)$  can be achieved at point

$$z_0 = \frac{c_{w,0} + \sqrt{c_{w,0}^2 + 2c_{w,0}/c_{\text{res}}}}{2},$$

and  $f(z)$  is monotone increasing in the interval  $[c_{w,0}, z_0]$ . Thus, if we choose time  $t$  properly, say  $t \leq c\sqrt{m}$ ,  $c$  being small enough, the following holds

$$\zeta(t) \leq \frac{c_{w,0} + \sqrt{c_{w,0}^2 + 2c_{w,0}/c_{\text{res}}}}{2}.$$

In other words, if  $t \leq c\sqrt{m}$  for some small enough  $c > 0$ , we have

$$\zeta(t) \leq c_{w,t} \leq \frac{c_{w,0} + \sqrt{c_{w,0}^2 + 2c_{w,0}/c_{\text{res}}}}{2},$$

where the last quantity is independent of depth  $L$  and time  $t$ , and we denote this by

$$\bar{c} = \frac{c_{w,0} + \sqrt{c_{w,0}^2 + 2c_{w,0}/c_{\text{res}}}}{2},$$

which finishes the proof of Proposition A.3. □

We state the inductive relation (A.12) as a proposition.

**Proposition A.4.** *Under Assumptions 4.1 and 4.2, we have with high probability with respect to the random initialization, for time  $t \lesssim \sqrt{m}$  with  $0 \leq l \leq L$ ,*

$$\|\mathbf{x}^{[l]}\|_2 \leq C, \tag{A.16}$$

where  $C > 0$  is a constant, independent of depth  $L$ .

**Remark A.3.** We shall note that the constant  $C$  in Proposition A.4 only depends on  $c_{\text{res}}, c_{w,0}$  and  $c_{\sigma}$ . However, for a fully-connected feedforward network, (A.16) in Proposition A.4 become

$$\|\mathbf{x}^{[l]}\|_2 \leq C 2^l. \tag{A.17}$$

Note that the 2-norm for each output layer increase exponentially layer by layer for fully-connected network, showing that ResNet possesses more stability compared with fully-connected network.

Next we end this part by making an Apriori estimate on the  $L^2$ -norm for arbitrary vector  $\mathbf{v}(t) \in \mathbb{A}_0$ .

**Proposition A.5.** *Under Assumptions 4.1 and 4.2, with high probability with respect to the random initialization, uniformly for any vector  $\mathbf{v}(t) \in \mathbb{A}_0$  and time  $t \lesssim \sqrt{m}$ , the following holds*

$$\|\mathbf{v}(t)\|_2 \leq c\sqrt{m}, \quad (\text{A.18})$$

where  $c > 0$  is a constant independent of depth  $L$  and time  $t$ .

*Proof.* We shall start our analysis on the whole expressions in set  $\mathbb{A}_0$ . For any vector  $\mathbf{v}(t) \in \mathbb{A}_0$ , we can write  $\mathbf{v}(t) = \mathbf{e}_s \mathbf{e}_{s-1} \cdots \mathbf{e}_1 \mathbf{e}_0$  with  $0 \leq s \leq 4L$ .

We start with the estimate on  $\mathbf{e}_0$ , since  $\mathbf{e}_0$  is chosen following the rules:

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \left\{ \sqrt{m} \mathbf{x}_\beta^{[1]}, \sqrt{m} \mathbf{x}_\beta^{[2]}, \dots, \sqrt{m} \mathbf{x}_\beta^{[L]} \right\}_{1 \leq \beta \leq n} \right\}.$$

- (a). If  $\mathbf{e}_0 = \mathbf{a}_t$ , then by Lemma A.2, for  $t \lesssim \sqrt{m}$ ,

$$\|\mathbf{a}_t\|_2 \leq c_{w,t} \sqrt{m} \leq c\sqrt{m}.$$

- (b). If  $\mathbf{e}_0 = \sqrt{m} \mathbf{x}_\beta^{[l]}$  where  $1 \leq l \leq L$ , then based on Proposition A.4, for  $t \lesssim \sqrt{m}$

$$\left\| \sqrt{m} \mathbf{x}_\beta^{[l]} \right\|_2 = \sqrt{m} \left\| \mathbf{x}_\beta^{[l]} \right\|_2 \leq c\sqrt{m}.$$

Now we proceed to other terms in the expression  $\mathbf{e}_j$  where  $j \geq 1$ .

- (i). If  $\mathbf{e}_j = \sigma_{[l]}^{(1)}(\mathbf{x}_\beta)$ , then we have

$$\begin{aligned} \|\mathbf{v}(t)\|_2 &= \|\mathbf{e}_s \mathbf{e}_{s-1} \cdots \mathbf{e}_1 \mathbf{e}_0\|_2 \\ &= \|\mathbf{e}_s\|_{2 \rightarrow 2} \|\mathbf{e}_{s-1}\|_{2 \rightarrow 2} \cdots \|\mathbf{e}_1\|_{2 \rightarrow 2} \|\mathbf{e}_0\|_2. \end{aligned}$$

Since  $\|\sigma_{[l]}^{(1)}(\mathbf{x}_\beta)\|_{2 \rightarrow 2} \leq C_{\text{Lip}} \leq 1$ , thus for all  $j \geq 1$  with  $\mathbf{e}_j = \sigma_{[l]}^{(1)}(\mathbf{x}_\beta)$

$$\|\mathbf{v}(t)\|_2 \leq (C_{\text{Lip}})^{4L} c\sqrt{m} \leq c\sqrt{m}.$$

- (ii). If  $\mathbf{e}_j = \mathbf{E}_{t,\beta}^{[l]}$  or  $\mathbf{e}_j = (\mathbf{E}_{t,\beta}^{[l]})^\top$ , then based on Proposition A.3

$$\|\mathbf{v}(t)\|_2 = \|\mathbf{e}_s\|_{2 \rightarrow 2} \|\mathbf{e}_{s-1}\|_{2 \rightarrow 2} \cdots \|\mathbf{e}_1\|_{2 \rightarrow 2} \|\mathbf{e}_0\|_2.$$

Since

$$\left\| \mathbf{E}_{t,\beta}^{[l]} \right\|_{2 \rightarrow 2} = \left\| (\mathbf{E}_{t,\beta}^{[l]})^\top \right\|_{2 \rightarrow 2} \leq \left( 1 + \frac{c_{\text{res}} C_{\text{Lip}}}{L} \zeta(t) \right) \leq \left( 1 + \frac{c_{\text{res}} c_{w,t}}{L} \right),$$

thus for all  $j \geq 1$  with  $\mathbf{e}_j = \mathbf{E}_{t,\beta}^{[l]}$  or  $\mathbf{e}_j = (\mathbf{E}_{t,\beta}^{[l]})^\top$ ,

$$\|\mathbf{v}(t)\|_2 \leq \left( 1 + \frac{c_{\text{res}} c_{w,t}}{L} \right)^s \|\mathbf{e}_0\|_2,$$

then by taking supreme on  $0 \leq s \leq 4L$ , we have

$$\begin{aligned} \|v(t)\|_2 &\leq \left(1 + \frac{c_{\text{res}} c_{w,t}}{L}\right)^{4L} \|e_0\|_2 \\ &\leq c \exp(4c_{\text{res}} c_{w,t}) \sqrt{m} \leq c\sqrt{m}. \end{aligned}$$

Combining these two observations, we finish the proof. □

Thus, if we define the quantity  $\xi_{\infty,0}(t)$  as follows,

$$\xi_{\infty,0}(t) = \sup_{0 \leq t' \leq t} \{\|v(t')\|_2 : v(t') \in \mathbb{A}_0\}. \tag{A.19}$$

Then directly from Proposition A.5, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ , the following holds

$$\xi_{\infty,0}(t) \leq c\sqrt{m}. \tag{A.20}$$

### A.3 Apriori $L^\infty$ bounds for expressions in $\mathbb{A}_0$

In this part, we shall make estimate on the quantity  $\eta_{\infty,0}(t)$  defined below

$$\eta_{\infty,0}(t) = \sup_{0 \leq t' \leq t} \{\|v(t')\|_\infty : v(t') \in \mathbb{A}_0\}. \tag{A.21}$$

We shall begin by a lemma on the  $\|\cdot\|_\infty$  norm of a standard Gaussian vector.

**Lemma A.3.** *For any i.i.d. normal distribution  $X_1, X_2, \dots, X_m \sim \mathcal{N}(0,1)$ , it holds with high probability that the  $L^\infty$ -norm of the gaussian vector  $\mathbf{X} = (X_1, X_2, \dots, X_m)^\top$  is upper bounded by*

$$\|\mathbf{X}\|_\infty \leq (\ln m)^C,$$

for some really large constant  $C > 0$ .

*Proof.* For any  $X_i \sim \mathcal{N}(0,1)$ , we have that for some  $\varepsilon, \lambda > 0$

$$\begin{aligned} \mathbb{P}(X_i \geq \varepsilon) &= \mathbb{P}(\exp(\lambda X_i) \geq \exp(\lambda \varepsilon)) \\ &\leq \frac{\mathbb{E}(\exp(\lambda X_i))}{\exp(\lambda \varepsilon)} = \frac{\exp(\frac{1}{2}\lambda^2)}{\exp(\lambda \varepsilon)} = \exp\left(\frac{1}{2}\lambda^2 - \lambda \varepsilon\right). \end{aligned}$$

We optimize over  $\lambda$ ,

$$\mathbb{P}(X_i \geq \varepsilon) \leq \min_{\lambda > 0} \exp\left(\frac{1}{2}\lambda^2 - \lambda \varepsilon\right) = \exp\left(-\frac{\varepsilon^2}{2}\right).$$

By taking absolute value

$$\mathbb{P}(|X_i| \geq \varepsilon) \leq 2 \exp\left(-\frac{\varepsilon^2}{2}\right).$$

Hence if we take over  $m$  unions

$$\mathbb{P}(\|\mathbf{X}\|_\infty \geq \varepsilon) \leq 2m \exp\left(-\frac{\varepsilon^2}{2}\right).$$

Set  $\varepsilon = (\ln m)^C$ , we have that

$$\mathbb{P}\left(\|\mathbf{X}\|_\infty \leq (\ln m)^C\right) \geq 1 - 2m \exp\left(-\frac{(\ln m)^{2C}}{2}\right).$$

Note that when  $C > 0$  is really large,  $(\ln m)^C \approx m^\varepsilon$  for some small  $\varepsilon > 0$ . □

We now state a lemma on the matrix two to infinity norm.

**Lemma A.4.** *Given a matrix  $\mathbf{W} \in \mathbb{R}^{m \times m}$  with each entry  $W_{i,j} \sim \mathcal{N}(0,1)$ , then with high probability, the following holds*

$$\|\mathbf{W}\|_{2 \rightarrow \infty} = \sup_{\|x\|_2=1} \|\mathbf{W}x\|_\infty \leq (\ln m)^C. \tag{A.22}$$

*Proof.* Note that  $\mathbf{W}x$  shares the same distribution as the Gaussian vector  $\mathbf{X}$  in Lemma A.3, i.e.,  $\mathbf{W}x \sim \mathbf{X}$ . Then apply Lemma A.3 directly, we obtain the result. □

Finally, to evaluate  $\eta_{\infty,0}(t)$ , we need to state a lemma.

**Lemma A.5.** *Under Assumptions 4.1 and 4.2, for any vector  $v(t) \in \mathbb{A}_0$ , we can write*

$$v(t) = e_s e_{s-1} \cdots e_1 e_0, \quad 0 \leq s \leq 4L, \quad t \geq 0.$$

For some vectors in  $\mathbb{A}_0$  with **length**  $q$ , we define  $\eta_{q,0}(t)$  as

$$\eta_{q,0}(t) := \sup_{0 \leq t' \leq t} \left\{ \|\mathbf{v}_q(t')\|_\infty \mid \mathbf{v}_q(t') = \mathbf{e}_q \mathbf{e}_{q-1} \cdots \mathbf{e}_1 \mathbf{e}_0, \mathbf{v}_q(t') \in \mathbb{A}_0 \right\}. \tag{A.23}$$

Moreover, we define  $\omega(t)$  as

$$\omega(t) := \sup_{0 \leq t' \leq t} \max \left\{ \|\mathbf{W}_{t'}^{[2]}\|_{2 \rightarrow \infty}, \left\| \left( \mathbf{W}_{t'}^{[2]} \right)^\top \right\|_{2 \rightarrow \infty}, \dots, \|\mathbf{W}_{t'}^{[L]}\|_{2 \rightarrow \infty}, \left\| \left( \mathbf{W}_{t'}^{[L]} \right)^\top \right\|_{2 \rightarrow \infty} \right\},$$

then with high probability with respect to the random initialization, for  $t \lesssim \sqrt{m}$ ,

$$\eta_{q,0}(t) \leq \eta_{0,0}(t) + c \omega(t) \left( 1 + \frac{c_{\text{res}} c_{w,t}}{L} \right)^q, \tag{A.24}$$

where constant  $c > 0$  is independent of depth  $L$ , and  $c_{w,t}$  has been defined in Proposition A.3.

*Proof.* Since for any vector  $\mathbf{v}_q(t) \in \mathbb{A}_0$  of length  $q$ ,  $0 \leq q \leq 4L$ , we can write  $\mathbf{v}_q(t)$  into

$$\mathbf{v}_q(t) = \mathbf{e}_q \mathbf{e}_{q-1} \cdots \mathbf{e}_1 \mathbf{e}_0,$$

then we shall prove (A.24) by performing induction on  $q$ . Firstly, for  $q=0$ , (A.24) is trivial. While for  $q \geq 1$ , we shall investigate on the terms  $\mathbf{e}_j$  in the expression  $\mathbf{v}_q(t)$ , where  $j \geq 1$ .

- (i). If  $\mathbf{e}_j = \sigma_{[l]}^{(1)}(\mathbf{x}_\beta)$ , then we have

$$\begin{aligned} \|\mathbf{v}_q(t)\|_\infty &= \|\mathbf{e}_q \mathbf{e}_{q-1} \cdots \mathbf{e}_1 \mathbf{e}_0\|_\infty \\ &= \|\mathbf{e}_q\|_{\infty \rightarrow \infty} \|\mathbf{e}_{q-1}\|_{\infty \rightarrow \infty} \cdots \|\mathbf{e}_1\|_{\infty \rightarrow \infty} \|\mathbf{e}_0\|_\infty, \end{aligned}$$

since  $\|\sigma_{[l]}^{(1)}(\mathbf{x}_\beta)\|_{\infty \rightarrow \infty} \leq C_{\text{Lip}} \leq 1$ , we have

$$\|\mathbf{v}_q(t)\|_\infty \leq (C_{\text{Lip}})^q c(\ln m)^C \leq c(\ln m)^C.$$

- (ii). If  $\mathbf{e}_j = \mathbf{E}_{t,\beta}^{[l]}$  or  $\mathbf{e}_j = \left(\mathbf{E}_{t,\beta}^{[l]}\right)^\top$  where  $2 \leq l \leq L$ ,  $\|\mathbf{e}_j\|_{\infty \rightarrow \infty} \geq 1$ , so we need to tackle it differently

$$\begin{aligned} \|\mathbf{v}_q(t)\|_\infty &= \left\| \mathbf{E}_{t,\beta}^{[l]} \mathbf{v}_{q-1}(t) \right\|_\infty \\ &= \left\| \mathbf{v}_{q-1}(t) + \frac{c_{\text{res}}}{L} \sigma_{[l]}^{(1)}(\mathbf{x}_\beta) \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \mathbf{v}_{q-1}(t) \right\|_\infty \\ &\leq \|\mathbf{v}_{q-1}(t)\|_\infty + \frac{c_{\text{res}} C_{\text{Lip}}}{L} \left\| \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right\|_{2 \rightarrow \infty} \|\mathbf{v}_{q-1}(t)\|_2, \end{aligned}$$

or

$$\begin{aligned} \|\mathbf{v}_q(t)\|_\infty &= \left\| \left(\mathbf{E}_{t,\beta}^{[l]}\right)^\top \mathbf{v}_{q-1}(t) \right\|_\infty \\ &= \left\| \mathbf{v}_{q-1}(t) + \frac{c_{\text{res}}}{L} \left(\frac{\mathbf{W}_t^{[l]}}{\sqrt{m}}\right)^\top \sigma_{[l]}^{(1)}(\mathbf{x}_\beta) \mathbf{v}_{q-1}(t) \right\|_\infty \\ &\leq \|\mathbf{v}_{q-1}(t)\|_\infty + \frac{c_{\text{res}} C_{\text{Lip}}}{L} \left\| \left(\frac{\mathbf{W}_t^{[l]}}{\sqrt{m}}\right)^\top \right\|_{2 \rightarrow \infty} \|\mathbf{v}_{q-1}(t)\|_2, \end{aligned}$$

recall the definition of  $\omega(t)$ , we have

$$\|\mathbf{v}_q(t)\|_\infty \leq \|\mathbf{v}_{q-1}(t)\|_\infty + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \|\mathbf{v}_{q-1}(t)\|_2.$$

Based on Proposition A.3

$$\|\mathbf{v}_{q-1}(t)\|_2 \leq c \left(1 + \frac{c_{\text{res}} c_{\omega,t}}{L}\right)^{q-1} \sqrt{m},$$



then

$$\begin{aligned} \|\mathbf{v}_q(t)\|_\infty &\leq \|\mathbf{v}_{q-1}(t)\|_\infty + \frac{c_{\text{res}}}{L\sqrt{m}}\omega(t)\|\mathbf{v}_{q-1}(t)\|_2 \\ &\leq \|\mathbf{v}_{q-1}(t)\|_\infty + \frac{c_{\text{res}}}{L}\omega(t)\left(1 + \frac{c_{\text{res}}c_{w,t}}{L}\right)^{q-1}. \end{aligned}$$

Inductively we have

$$\begin{aligned} \|\mathbf{v}_q(t)\|_\infty &\leq \|\mathbf{v}_0(t)\|_\infty + \frac{c}{c_{w,t}}\omega(t)\left(1 + \frac{c_{\text{res}}c_{w,t}}{L}\right)^q \\ &\leq \|\mathbf{v}_0(t)\|_\infty + c\omega(t)\left(1 + \frac{c_{\text{res}}c_{w,t}}{L}\right)^q, \end{aligned}$$

where we use the property of a geometric sum. By taking supreme on both sides, we have

$$\eta_{q,0}(t) \leq \eta_{0,0}(t) + c\omega(t)\left(1 + \frac{c_{\text{res}}c_{w,t}}{L}\right)^q.$$

This completes the proof. □

Based on these lemmas, recall definition (A.21), we are able to make a proposition on the quantity  $\eta_{\infty,0}(t)$  at  $t=0$ .

**Proposition A.6.** *Under Assumptions 4.1 and 4.2, with high probability with respect to the random initialization*

$$\eta_{\infty,0}(0) \leq c(\ln m)^C, \tag{A.25}$$

where  $c, C > 0$  are constants independent of the depth  $L$ .

*Proof.* As always, for any vector  $\mathbf{v}(t) \in \mathbf{A}_0$ , we can write  $\mathbf{v}(t)$  as

$$\mathbf{v}(t) = \mathbf{e}_s \mathbf{e}_{s-1} \cdots \mathbf{e}_1 \mathbf{e}_0, \quad 0 \leq s \leq 4L.$$

We start with the estimate on  $\eta_{0,0}(0)$ , since  $\mathbf{e}_0$  is chosen following the rules:

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \left\{ \sqrt{m}\mathbf{x}_\beta^{[1]}, \sqrt{m}\mathbf{x}_\beta^{[2]}, \dots, \sqrt{m}\mathbf{x}_\beta^{[L]} \right\}_{1 \leq \beta \leq n} \right\}.$$

- (a). If  $\mathbf{e}_0 = \mathbf{a}_t$ , then at  $t=0$ , by Lemma A.3,

$$\|\mathbf{a}_0\|_\infty \leq (\ln m)^C.$$

- (b). If  $\mathbf{e}_0 = \sqrt{m}\mathbf{x}_\beta^{[l]}$ , starting with  $l=1$

$$\begin{aligned} \left\| \sqrt{m}\mathbf{x}_\beta^{[1]} \right\|_\infty &= \sqrt{c_\sigma} \left\| \sigma\left(\mathbf{W}_0^{[1]}\mathbf{x}_\beta\right) \right\|_\infty \\ &\leq \sqrt{c_\sigma} \left( |\sigma(0)| + C_{\text{Lip}} \left\| \mathbf{W}_0^{[1]}\mathbf{x}_\beta \right\|_\infty \right) \\ &\leq \sqrt{c_\sigma} \left( C_{\text{Lip}} + C_{\text{Lip}} \left\| \mathbf{W}_0^{[1]} \right\|_{2 \rightarrow \infty} \|\mathbf{x}_\beta\|_2 \right) \\ &\leq \sqrt{c_\sigma} C_{\text{Lip}} \left( 1 + (\ln m)^C \right) \leq c(\ln m)^C, \end{aligned}$$

moreover, for  $l \geq 1$ , based on Proposition A.4,

$$\begin{aligned} \left\| \sqrt{m} \mathbf{x}_\beta^{[l]} \right\|_\infty &\leq \left\| \sqrt{m} \mathbf{x}_\beta^{[l-1]} \right\|_\infty + \frac{c_{\text{res}}}{L} \left\| \sigma \left( \mathbf{W}_0^{[l]} \mathbf{x}_\beta^{[l-1]} \right) \right\|_\infty \\ &\leq \left\| \sqrt{m} \mathbf{x}_\beta^{[l-1]} \right\|_\infty + \frac{c_{\text{res}}}{L} \left( C_{\text{Lip}} + C_{\text{Lip}} \left\| \mathbf{W}_0^{[l]} \right\|_{2 \rightarrow \infty} \left\| \mathbf{x}_\beta^{[l-1]} \right\|_2 \right) \\ &\leq \left\| \sqrt{m} \mathbf{x}_\beta^{[l-1]} \right\|_\infty + \frac{c_{\text{res}} C_{\text{Lip}}}{L} \left( 1 + C(\ln m)^C \right) \\ &\leq \left\| \sqrt{m} \mathbf{x}_\beta^{[l-1]} \right\|_\infty + \frac{c}{L} (\ln m)^C, \end{aligned}$$

inductively for  $1 \leq l \leq L$ ,

$$\left\| \sqrt{m} \mathbf{x}_\beta^{[l]} \right\|_\infty \leq c \left( 1 + \frac{l}{L} \right) (\ln m)^C \leq c (\ln m)^C, \tag{A.26}$$

where  $c$  is independent of the depth  $L$ .

Hence we have

$$\eta_{0,0}(0) \leq c (\ln m)^C. \tag{A.27}$$

Directly from Lemma A.5

$$\begin{aligned} \eta_{q,0}(0) &\leq \eta_{0,0}(0) + c (\ln m)^C \left( 1 + \frac{c_{\text{res}} c_{w,0}}{L} \right)^q \\ &\leq c (\ln m)^C + c (\ln m)^C \exp(4c_{\text{res}} c_{w,0}) \leq c (\ln m)^C, \end{aligned}$$

by taking supreme on  $0 \leq q \leq 4L$ , we finish our proof. □

Our next proposition is on  $\eta_{\infty,0}(t)$  for time  $0 \leq t \leq \sqrt{m} / (\ln m)^{C'}$ .

**Proposition A.7.** *Under Assumptions 4.1 and 4.2, with high probability with respect to the random initialization, for time  $0 \leq t \leq \sqrt{m} / (\ln m)^{C'}$ , the following holds*

$$\eta_{\infty,0}(t) \leq c (\ln m)^C, \tag{A.28}$$

where  $c, C, C' > 0$  are constants independent of the depth  $L$ .

*Proof.* We shall start with the estimate on  $\eta_{0,0}(t)$ , since  $\mathbf{e}_0$  is chosen following the rules:

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \left\{ \sqrt{m} \mathbf{x}_\beta^{[1]}, \sqrt{m} \mathbf{x}_\beta^{[2]}, \dots, \sqrt{m} \mathbf{x}_\beta^{[L]} \right\}_{1 \leq \beta \leq n} \right\}.$$

We observe that from the replacement rules given in Section 5.1,

$$\begin{aligned} \partial_t \mathbf{a}_t &= -\frac{1}{n} \sum_{\beta=1}^n \frac{1}{\sqrt{m}} \sqrt{m} \mathbf{x}_\beta^{[L]} (f_\beta(t) - y_\beta), \\ \partial_t \sqrt{m} \mathbf{x}_\alpha^{[l]} &= -\frac{1}{n} \sum_{\beta=1}^n \frac{c_\sigma}{\sqrt{m}} \mathbf{E}_{t,\alpha}^{[2:l]} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\alpha) \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle (f_\beta(t) - y_\beta) \\ &\quad + -\frac{1}{n} \sum_{\beta=1}^n \sum_{k=2}^l \frac{c_{\text{res}}^2}{L^2 \sqrt{m}} \mathbf{E}_{t,\alpha}^{[(k+1):l]} \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\alpha) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(k+1):L]} \right)^\top \\ &\quad \quad \mathbf{a}_t \langle \mathbf{x}_\alpha^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle (f_\beta(t) - y_\beta). \end{aligned}$$

Since for  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ , then by Proposition A.4

$$\begin{aligned} \partial_t \|\mathbf{a}_t\|_\infty &\leq \frac{C}{\sqrt{m}} \left\| \sqrt{m} \mathbf{x}_\beta^{[L]} \right\|_\infty, \\ \partial_t \left\| \sqrt{m} \mathbf{x}_\alpha^{[l]} \right\|_\infty &\leq \sum_{k=1}^l \frac{C}{\sqrt{m}} \left\| \mathbf{E}_{t,\alpha}^{[(k+1):l]} \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\alpha) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right\|_\infty. \end{aligned}$$

By taking supreme on time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ , we have

$$\eta_{0,0}(t) \leq c(\ln m)^C + \frac{C}{\sqrt{m}} \int_0^t \eta_{\infty,0}(s) ds. \tag{A.29}$$

For the auxiliary term  $\omega(t)$ , from the replacement rules again, for  $2 \leq l \leq L$

$$\begin{aligned} \partial_t \mathbf{W}_t^{[l]} &= -\frac{1}{n} \sum_{\beta=1}^n \frac{c_{\text{res}}}{L \sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+1):L]} \right)^\top \mathbf{a}_t \otimes (\mathbf{x}_\beta^{[l-1]})^\top (f_\beta(t) - y_\beta), \\ \partial_t \left( \mathbf{W}_t^{[l]} \right)^\top &= -\frac{1}{n} \sum_{\beta=1}^n \frac{c_{\text{res}}}{L \sqrt{m}} \mathbf{x}_\beta^{[l-1]} \otimes \left( \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+1):L]} \right)^\top \mathbf{a}_t \right)^\top (f_\beta(t) - y_\beta), \end{aligned}$$

then by Proposition A.5

$$\begin{aligned} \partial_t \left\| \mathbf{W}_t^{[l]} \right\|_{2 \rightarrow \infty} &\leq \frac{C}{\sqrt{m}} \left\| \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+1):L]} \right)^\top \mathbf{a}_t \right\|_\infty, \\ \partial_t \left\| \left( \mathbf{W}_t^{[l]} \right)^\top \right\|_{2 \rightarrow \infty} &\leq \frac{C}{\sqrt{m}} \left\| \sqrt{m} \mathbf{x}_\beta^{[l-1]} \right\|_\infty. \end{aligned}$$

Hence by taking supreme on time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ , we have

$$\omega(t) \leq (\ln m)^C + \frac{C}{\sqrt{m}} \int_0^t \eta_{\infty,0}(s) ds. \tag{A.30}$$

Directly from Lemma A.5

$$\begin{aligned} \eta_{q,0}(t) &\leq \eta_{0,0}(t) + c \omega(t) \left(1 + \frac{c_{\text{res}} c_{w,t}}{L}\right)^q \\ &\leq \left(c(\ln m)^C + \frac{C}{\sqrt{m}} \int_0^t \eta_{\infty,0}(s) ds\right) \left(1 + \left(1 + \frac{c_{\text{res}} c_{w,t}}{L}\right)^q\right). \end{aligned}$$

Finally by taking supreme on  $0 \leq q \leq 4L$ , we have

$$\eta_{\infty,0}(t) \leq c(\ln m)^C + \frac{C}{\sqrt{m}} \int_0^t \eta_{\infty,0}(s) ds.$$

This gives us a Gronwall-type inequality, we have that

$$\eta_{\infty,0}(t) \leq c(\ln m)^C \exp\left(\frac{Ct}{\sqrt{m}}\right).$$

To sum up, for  $t \leq \sqrt{m}/(\ln m)^C$ , the following holds

$$\eta_{\infty,0}(t) \leq c(\ln m)^C, \tag{A.31}$$

which finishes the proof. □

#### A.4 Apriori $L^2$ and $L^\infty$ bounds for expression in $\mathbb{A}_r, r \geq 1$

In this part, we shall make estimates for  $\|\cdot\|_\infty$  and  $\|\cdot\|_2$  of vectors belonging to higher order sets, i.e.,  $\mathbb{A}_r, r \geq 1$ . Then it is natural for us to define several quantities for some vectors in  $\mathbb{A}_r$  with length  $q$

$$\tilde{\xi}_{q,r}(t) := \sup_{0 \leq t' \leq t} \{ \|v_q(t')\|_2 : v_q(t') = e_q e_{q-1} \cdots e_1 e_0, v_q(t') \in \mathbb{A}_r \}. \tag{A.32}$$

Note that from Propositions A.3 and A.5,

$$\tilde{\xi}_{q,0}(t) \leq c \left(1 + \frac{c_{\text{res}} c_{w,t}}{L}\right)^q \sqrt{m}. \tag{A.33}$$

Moreover, we define that

$$\tilde{\xi}_{\infty,r}(t) = \sup_{0 \leq q \leq 4L} \{ \tilde{\xi}_{q,r}(t) \}. \tag{A.34}$$

Then by taking supreme on  $0 \leq q \leq 4L$  in (A.33)

$$\tilde{\xi}_{\infty,0}(t) \leq c\sqrt{m}, \tag{A.35}$$

and recalling the definition we made in Section A.3, we similarly define

$$\eta_{q,r}(t) := \sup_{0 \leq t' \leq t} \{ \|v_q(t')\|_\infty \mid v_q(t') = e_q e_{q-1} \cdots e_1 e_0, v_q(t') \in \mathbb{A}_r \}. \tag{A.36}$$

Moreover, we also define that

$$\eta_{\infty,r}(t) = \sup_{0 \leq q \leq 4L} \{\eta_{q,r}(t)\}. \quad (\text{A.37})$$

Once again, for any vector  $v(t) \in \mathbb{A}_r$ , it can be written into

$$v(t) = e_s e_{s-1} \cdots e_1 e_0, \quad 0 \leq s \leq 4L.$$

We shall start with the estimate on  $e_0$ . Since  $e_0$  is chosen following the rules:

$$e_0 \in \left\{ a_t, \mathbf{1}, \{\sqrt{m}x_\beta^{[1]}, \sqrt{m}x_\beta^{[2]}, \dots, \sqrt{m}x_\beta^{[L]}\}_{1 \leq \beta \leq n} \right\}.$$

$\|\mathbf{1}\|_\infty = 1, \|\mathbf{1}\|_2 = \sqrt{m}$ , then for time  $0 \leq t \leq \sqrt{m}/(\ln m)^C$ , by Propositions A.3 and A.7,

$$\xi_{0,r}(t) \leq c\sqrt{m}, \quad \eta_{0,r}(t) \leq c(\ln m)^C.$$

Now we proceed to other terms in the expression  $e_j$  where  $j \geq 1$ . For each  $e_j$ , there are several cases:

- (i)  $e_j = \sigma_{[l]}^{(1)}(x_\beta)$ ,  $e_j = \mathbf{E}_{t,\beta}^{[l]}$  or  $e_j = \left(\mathbf{E}_{t,\beta}^{[l]}\right)^\top$ ,  $2 \leq l \leq L$ .
- (ii)  $e_j = \text{diag}(\mathbf{g})$ .
- (iii)

$$e_j = \sigma_{[l]}^{(u+1)}(x_\beta) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \left( \frac{c_{\text{res}} \mathbf{W}_t^{[l]}}{L \sqrt{m}} \right)^{Q_{u+1}},$$

or

$$e_j = \left( \frac{c_{\text{res}} \left( \mathbf{W}_t^{[l]} \right)^\top}{L \sqrt{m}} \right)^{Q_{u+1}} \sigma_{[l]}^{(u+1)}(x_\beta) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right).$$

By our observation, the total number of  $\text{diag}$  operations in  $v(t) \in \mathbb{A}_r$  is  $r$ , and that is how we characterize a vector belonging to different hierarchical sets. Especially if for one of those  $e_j$  belongs to case (iii), there are two scenarios:

- $Q_{u+1} = 0$ , then  $e_j$  is just multiplication of several diagonal matrices, being a special situation for case (ii).

- $Q_{u+1} = 1$ , since diagonal matrices commute,  $e_j$  writes into

$$e_j = \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \sigma_{[l]}^{(u+1)}(\mathbf{x}_\beta) \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}},$$

or

$$e_j = \left( \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^\top \sigma_{[l]}^{(u+1)}(\mathbf{x}_\beta) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right).$$

We shall take advantage of the special structure of  $e_j$ . Define a new type of skip-connection matrix,  $\tilde{\mathbf{E}}_{t,\beta}^{[l,r]}$ , for  $r \geq 2$ :

$$\tilde{\mathbf{E}}_{t,\beta}^{[l,r]} := \left( \mathbf{I}_m + \frac{c_{\text{res}}}{L} \sigma_{[l]}^{(r)}(\mathbf{x}_\beta) \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right). \tag{A.38}$$

Then we can write  $e_j$  into

$$\begin{aligned} e_j &= \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \sigma_{[l]}^{(u+1)}(\mathbf{x}_\beta) \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \\ &= \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \tilde{\mathbf{E}}_{t,\beta}^{[l,u+1]} \\ &\quad - \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right), \end{aligned}$$

or

$$\begin{aligned} e_j &= \left( \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^\top \sigma_{[l]}^{(u+1)}(\mathbf{x}_\beta) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \\ &= \left( \tilde{\mathbf{E}}_{t,\beta}^{[l,u+1]} \right)^\top \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \\ &\quad - \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right). \end{aligned}$$

To illustrate such relation, if some vector  $\bar{v}(t)$  contains  $e_j$  belonging to case (iii), we write it as

$$\begin{aligned} \bar{v}(t) &= e_s e_{s-1} \cdots e_{j+1} \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_1} g_1 \right) \\ &\quad \cdots \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_u} g_u \right) \sigma_{[l]}^{(u+1)}(x_\beta) \frac{c_{\text{res}}}{L} \frac{W_t^{[l]}}{\sqrt{m}} e_{j-1} \cdots e_0 \\ &= e_s e_{s-1} \cdots e_{j+1} \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_1} g_1 \right) \cdots \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_u} g_u \right) \tilde{E}_{t,\beta}^{[l,u+1]} e_{j-1} \cdots e_0 \\ &= e_s e_{s-1} \cdots e_{j+1} \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_1} g_1 \right) \cdots \text{diag} \left( \left( \frac{W_t^{[l]}}{\sqrt{m}} \right)^{Q_u} g_u \right) e_{j-1} \cdots e_0. \end{aligned}$$

From the above analysis, we are able to characterize an element in set  $\mathbb{A}_r$ . If  $v(t) \in \mathbb{A}_r$ , then as always, we write it as

$$v(t) = e_s e_{s-1} \cdots e_1 e_0, \quad 0 \leq s \leq 4L,$$

and there exists  $e_{j_1}, e_{j_2}, \dots, e_{j_k}$ , such that

$$\begin{aligned} e_{j_1} &= \text{diag} \left( \left( \frac{W_t^{[l_1]}}{\sqrt{m}} \right)^{Q_1} g_1 \right), \quad g_1 \in \mathbb{A}_{r_1-1}, \\ e_{j_2} &= \text{diag} \left( \left( \frac{W_t^{[l_2]}}{\sqrt{m}} \right)^{Q_2} g_2 \right), \quad g_2 \in \mathbb{A}_{r_2-1}, \\ &\vdots \\ e_{j_k} &= \text{diag} \left( \left( \frac{W_t^{[l_k]}}{\sqrt{m}} \right)^{Q_k} g_k \right), \quad g_k \in \mathbb{A}_{r_k-1}, \end{aligned}$$

with

$$r_1 + r_2 + \cdots + r_k = r, \quad r_1, r_2, \dots, r_k \in \mathbb{N}^+. \tag{A.39}$$

Eq. (A.39) serves as the counting of the number of diag operations contained in  $v(t)$ , while for other  $e_j (j \notin \{j_1, j_2, \dots, j_k, 0\})$ , chosen from the following sets

$$\left\{ E_{t,\beta}^{[l]}, \left( E_{t,\beta}^{[l]} \right)^\top : 2 \leq l \leq L \right\}_{1 \leq \beta \leq n'} \tag{A.40}$$

$$\left\{ \sigma_{[l]}^{(1)}(x_\beta) : 1 \leq l \leq L \right\}_{1 \leq \beta \leq n'} \tag{A.41}$$

$$\left\{ \tilde{E}_{t,\beta}^{[l,p]}, \left( \tilde{E}_{t,\beta}^{[l,p]} \right)^\top : 2 \leq l \leq L, p \geq 2 \right\}_{1 \leq \beta \leq n'} \tag{A.42}$$

note that the elements in set (A.40) and set (A.42) share the same matrix properties, thanks to Assumption 4.1 concerning the activation function.

Hence, in order to make estimates on  $\xi_{q,r}(t)$  and  $\eta_{q,r}(t)$ , we shall perform induction on the number of diag operations contained in each vector.

**Proposition A.8.** *Under Assumptions 4.1 and 4.2, with high probability with respect to the random initialization, for some finite  $r \geq 1$  and time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ , the following holds*

$$\xi_{\infty,r}(t) \leq c(\ln m)^C \sqrt{m}, \tag{A.43}$$

$$\eta_{\infty,r}(t) \leq c(\ln m)^C, \tag{A.44}$$

where  $c, C, C' > 0$  are constants independent of depth  $L$ .

*Proof.* Recalling the definition of  $\omega(t)$ ,  $\eta_{\infty,0}(t)$  and  $\xi_{\infty,0}(t)$ , for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ , the following holds with high probability,

$$\begin{aligned} \omega(t) &\leq c(\ln m)^C, \\ \eta_{\infty,0}(t) &\leq c(\ln m)^C, \\ \xi_{\infty,0}(t) &\leq c\sqrt{m}. \end{aligned}$$

Let's start with  $r=1$ . For any  $v(t) \in \mathbb{A}_1$ , since there is only one solution to equation (A.39), there exists one and only one index  $i$ , such that  $e_i = \text{diag}(g)$  or  $e_i = \text{diag}\left(\frac{W_t^{[l]}}{\sqrt{m}}g\right)$ , with  $g \in \mathbb{A}_0$ . Then we have

$$\begin{aligned} \xi_{i,1}(t) &\leq \sup_{g \in \mathbb{A}_0} \|\text{diag}(g)\|_{2 \rightarrow 2} \xi_{i-1,0}(t) \\ &\leq \sup_{g \in \mathbb{A}_0} \|g\|_{\infty} \xi_{i-1,0}(t) \leq \eta_{\infty,0}(t) \xi_{i-1,0}(t) \leq c(\ln m)^C \xi_{i-1,0}(t), \\ \text{or } \xi_{i,1}(t) &\leq \sup_{g \in \mathbb{A}_0} \left\| \text{diag} \left( \frac{W_t^{[l]}}{\sqrt{m}} g \right) \right\|_{2 \rightarrow 2} \xi_{i-1,0}(t) \\ &\leq \sup_{g \in \mathbb{A}_0} \left\| \frac{W_t^{[l]}}{\sqrt{m}} g \right\|_{\infty} \xi_{i-1,0}(t) \leq \frac{\omega(t)}{\sqrt{m}} \xi_{\infty,0}(t) \xi_{i-1,0}(t) \leq c(\ln m)^C \xi_{i-1,0}(t), \end{aligned}$$

for  $q > i$ ,

$$\xi_{q,1}(t) \leq \left(1 + \frac{c_{\text{res}} c_{w,t}}{L}\right) \xi_{q-1,1}(t).$$

Inductively, we have

$$\xi_{q,1}(t) \leq \left(1 + \frac{c_{\text{res}} c_{w,t}}{L}\right)^{q-i} \xi_{i,1}(t).$$



By taking supreme on  $q$  and  $i$ , this leads to

$$\zeta_{\infty,1}(t) \leq \exp(4c_{\text{res}}c_{w,t})c(\ln m)^C \zeta_{\infty,0}(t) \leq c(\ln m)^C \sqrt{m}. \tag{A.45}$$

For  $\eta_{i,1}(t)$ , we have

$$\begin{aligned} \eta_{i,1}(t) &\leq \sup_{g \in \mathbb{A}_0} \|\text{diag}(g)\|_{\infty \rightarrow \infty} \eta_{i-1,0}(t) \\ &\leq \sup_{g \in \mathbb{A}_0} \|g\|_{\infty} \eta_{i-1,0}(t) \leq \eta_{\infty,0}(t) \eta_{i-1,0}(t) \leq c(\ln m)^C \eta_{i-1,0}(t), \\ \text{or } \eta_{i,1}(t) &\leq \sup_{g \in \mathbb{A}_0} \left\| \text{diag} \left( \frac{W_t^{[l]}}{\sqrt{m}} g \right) \right\|_{\infty \rightarrow \infty} \eta_{i-1,0}(t) \\ &\leq \sup_{g \in \mathbb{A}_0} \left\| \frac{W_t^{[l]}}{\sqrt{m}} g \right\|_{\infty} \eta_{i-1,0}(t) \leq \frac{\omega(t)}{\sqrt{m}} \zeta_{\infty,0}(t) \eta_{i-1,0}(t) \leq c(\ln m)^C \eta_{i-1,0}(t), \end{aligned}$$

and for  $q > i$ , inductively

$$\begin{aligned} \eta_{q,1}(t) &\leq \eta_{q-1,1}(t) + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \zeta_{q-1,1}(t) \\ &\leq \eta_{i,1}(t) + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \zeta_{q-1,1}(t) + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \zeta_{q-2,1}(t) + \dots + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \zeta_{i,1}(t). \end{aligned}$$

Then by taking supreme on  $q$  and  $i$ , combined with (A.45), we obtain

$$\begin{aligned} \eta_{\infty,1}(t) &\leq c(\ln m)^C \eta_{\infty,0}(t) + \frac{4c_{\text{res}}}{\sqrt{m}} \omega(t) \zeta_{\infty,1}(t) \\ &\leq c(\ln m)^C + c(\ln m)^C \leq c(\ln m)^C. \end{aligned}$$

In the following we assume that (A.43) and (A.44) holds for  $1, 2, \dots, r-1$  and prove it for  $r$ .

If  $v(t) \in \mathbb{A}_r$ , then as always, we write it as

$$v(t) = e_s e_{s-1} \dots e_1 e_0, \quad 0 \leq s \leq 4L,$$

and there exists  $e_{j_1}, e_{j_2}, \dots, e_{j_k}$ , such that

$$\begin{aligned} e_{j_1} &= \text{diag} \left( \left( \frac{W_t^{[l_1]}}{\sqrt{m}} \right)^{Q_1} g_1 \right), \quad g_1 \in \mathbb{A}_{r_1-1}, \\ e_{j_2} &= \text{diag} \left( \left( \frac{W_t^{[l_2]}}{\sqrt{m}} \right)^{Q_2} g_2 \right), \quad g_2 \in \mathbb{A}_{r_2-1}, \\ &\vdots \\ e_{j_k} &= \text{diag} \left( \left( \frac{W_t^{[l_k]}}{\sqrt{m}} \right)^{Q_k} g_k \right), \quad g_k \in \mathbb{A}_{r_k-1}, \end{aligned}$$

with

$$r_1 + r_2 + \dots + r_k = r, \quad r_1, r_2, \dots, r_k \in \mathbb{N}^+.$$

Let  $i$  be the largest index among  $j_1, j_2, \dots, j_k$ , i.e.,

$$i = \max\{j_1, j_2, \dots, j_k\},$$

and without loss of generality, let  $i = j_1$ . We have  $e_i = \text{diag}(g_1)$  or  $e_i = \text{diag}\left(\frac{W_t^{[l]}}{\sqrt{m}}g_1\right)$  with  $g_1 \in \mathbb{A}_{r_1-1}$ , then

$$\begin{aligned} \zeta_{i,r}(t) &\leq \sup_{g \in \mathbb{A}_{r_1-1}} \|\text{diag}(g)\|_{2 \rightarrow 2} \zeta_{i-1,r-r_1}(t) \\ &\leq \sup_{g \in \mathbb{A}_{r_1-1}} \|g\|_{\infty} \zeta_{i-1,r-r_1}(t) \leq \eta_{\infty,r_1-1}(t) \zeta_{i-1,r-r_1}(t) \leq c(\ln m)^C \zeta_{i-1,r-r_1}(t), \\ \text{or } \zeta_{i,r}(t) &\leq \sup_{g \in \mathbb{A}_{r_1-1}} \left\| \text{diag}\left(\frac{W_t^{[l]}}{\sqrt{m}}g\right) \right\|_{2 \rightarrow 2} \zeta_{i-1,r-r_1}(t) \leq \sup_{g \in \mathbb{A}_{r_1-1}} \left\| \frac{W_t^{[l]}}{\sqrt{m}}g \right\|_{\infty} \zeta_{i-1,r-r_1}(t) \\ &\leq \frac{\omega(t)}{\sqrt{m}} \zeta_{\infty,r_1-1}(t) \zeta_{i-1,r-r_1}(t) \leq c(\ln m)^C \zeta_{i-1,r-r_1}(t), \end{aligned}$$

inductively

$$\zeta_{q,r}(t) \leq \left(1 + \frac{c_{\text{res}} c_{w,t}}{L}\right)^{q-i} \zeta_{i,r-r_1}(t).$$

Then by taking supreme on  $q$  and  $i$ , we obtain

$$\zeta_{\infty,r}(t) \leq \exp(4c_{\text{res}} c_{w,t}) c(\ln m)^C \zeta_{\infty,r-r_1}(t) \leq c(\ln m)^C \sqrt{m}. \tag{A.46}$$

For  $\eta_{i,r}(t)$ , we have

$$\begin{aligned} \eta_{i,r}(t) &\leq \sup_{g \in \mathbb{A}_{r_1-1}} \|\text{diag}(g)\|_{\infty \rightarrow \infty} \eta_{i-1,r-r_1}(t) \\ &\leq \sup_{g \in \mathbb{A}_{r_1-1}} \|g\|_{\infty} \eta_{i-1,r-r_1}(t) \leq \eta_{\infty,r_1-1}(t) \eta_{i-1,r-r_1}(t) \leq c(\ln m)^C \eta_{i-1,r-r_1}(t), \\ \text{or } \eta_{i,r}(t) &\leq \sup_{g \in \mathbb{A}_{r_1-1}} \left\| \text{diag}\left(\frac{W_t^{[l]}}{\sqrt{m}}g\right) \right\|_{\infty \rightarrow \infty} \eta_{i-1,r-r_1}(t) \leq \sup_{g \in \mathbb{A}_{r_1-1}} \left\| \frac{W_t^{[l]}}{\sqrt{m}}g \right\|_{\infty} \eta_{i-1,r-r_1}(t) \\ &\leq \frac{\omega(t)}{\sqrt{m}} \zeta_{\infty,r_1-1}(t) \eta_{i-1,r-r_1}(t) \leq c(\ln m)^C \eta_{i-1,r-r_1}(t), \end{aligned}$$

and for  $q > i$

$$\begin{aligned} \eta_{q,r}(t) &\leq \eta_{q-1,r}(t) + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \zeta_{q-1,r}(t) \\ &\leq \eta_{i,r}(t) + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \zeta_{q-1,r}(t) + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \zeta_{q-2,r}(t) + \dots + \frac{c_{\text{res}}}{L\sqrt{m}} \omega(t) \zeta_{i,r}(t), \end{aligned}$$

then by taking supreme on  $q$  and  $i$ ,

$$\begin{aligned} \eta_{\infty,r}(t) &\leq c(\ln m)^C \eta_{\infty,r-r_1}(t) + \frac{4c_{\text{res}}}{\sqrt{m}} \omega(t) \zeta_{\infty,r}(t) \\ &\leq c(\ln m)^C + c(\ln m)^C \leq c(\ln m)^C. \end{aligned}$$

Note that from the proof, for different  $r$ , the constant  $c$  grows exponentially in  $r$ , while the growth rate of  $C$  is linear. □

### B Least eigenvalue of Gram matrices

We shall recall the Gram matrices defined in Section 4.2. We first define a series of matrices  $\{\tilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\tilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$ , and a series of vectors  $\{\tilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ . Given the input samples  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $\|\mathbf{x}_i\|_2 = 1$ , for  $1 \leq i \leq n$ , and  $\mathbf{x}_i \perp \mathbf{x}_j$ , for any  $i \neq j$ ,

$$\begin{aligned} \tilde{\mathbf{K}}_{ij}^{[0]} &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle, \\ \tilde{\mathbf{K}}_{ij}^{[1]} &= \mathbb{E}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[0]} & \tilde{\mathbf{K}}_{ij}^{[0]} \\ \tilde{\mathbf{K}}_{ji}^{[0]} & \tilde{\mathbf{K}}_{jj}^{[0]} \end{pmatrix}\right)} c_\sigma \sigma(u) \sigma(v), \\ \tilde{\mathbf{b}}_i^{[1]} &= \sqrt{c_\sigma} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[0]})} [\sigma(u)], \\ \tilde{\mathbf{A}}_{ij}^{[l]} &= \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[l-1]} & \tilde{\mathbf{K}}_{ij}^{[l-1]} \\ \tilde{\mathbf{K}}_{ji}^{[l-1]} & \tilde{\mathbf{K}}_{jj}^{[l-1]} \end{pmatrix}, \\ \tilde{\mathbf{K}}_{ij}^{[l]} &= \tilde{\mathbf{K}}_{ij}^{[l-1]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[l]}\right)} \left[ \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[l-1]} \sigma(v)}{L} + \frac{c_{\text{res}} \tilde{\mathbf{b}}_j^{[l-1]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(v)}{L^2} \right], \\ \tilde{\mathbf{b}}_i^{[l]} &= \tilde{\mathbf{b}}_i^{[l-1]} + \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)], \\ \tilde{\mathbf{A}}_{ij}^{[L+1]} &= \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[L]} & \tilde{\mathbf{K}}_{ij}^{[L]} \\ \tilde{\mathbf{K}}_{ji}^{[L]} & \tilde{\mathbf{K}}_{jj}^{[L]} \end{pmatrix}. \end{aligned}$$

Given these definitions, we define that for  $2 \leq l \leq L-1$ ,

$$\mathbf{K}_{ij}^{[L+1]} = \tilde{\mathbf{K}}_{ij}^{[L]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[L+1]}\right)} \left[ \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[L]} \sigma(v)}{L} + \frac{c_{\text{res}} \tilde{\mathbf{b}}_j^{[L]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(v)}{L^2} \right], \tag{B.1}$$

$$\mathbf{K}_{ij}^{[L]} = \frac{c_{\text{res}}^2}{L^2} \tilde{\mathbf{K}}_{ij}^{[L-1]} \mathbb{E}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[L]}\right)} \left[ \sigma^{(1)}(u) \sigma^{(1)}(v) \right], \tag{B.2}$$

$$\mathbf{K}_{ij}^{[l]} = \frac{c_{\text{res}}^2}{L^2} \tilde{\mathbf{K}}_{ij}^{[l-1]} \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \sigma_{[l]}^{(1)}(\mathbf{x}_i) \left( \mathbf{E}_{0,i}^{[(l+1):L]} \right)^\top \mathbf{a}_0, \sigma_{[l]}^{(1)}(\mathbf{x}_j) \left( \mathbf{E}_{0,j}^{[(l+1):L]} \right)^\top \mathbf{a}_0 \right\rangle, \tag{B.3}$$

$$\mathbf{K}_{ij}^{[1]} = c_\sigma \tilde{\mathbf{K}}_{ij}^{[0]} \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \sigma_{[1]}^{(1)}(\mathbf{x}_i) \left( \mathbf{E}_{0,i}^{[2:L]} \right)^\top \mathbf{a}_0, \sigma_{[1]}^{(1)}(\mathbf{x}_j) \left( \mathbf{E}_{0,j}^{[2:L]} \right)^\top \mathbf{a}_0 \right\rangle. \tag{B.4}$$

We shall state two lemmas concerning full rankness of the Gram matrices, which have been stated as Lemma F.1. and Lemma F.2. in Du et al. [12].

**Lemma B.1.** Assume  $\sigma(\cdot)$  is analytic and not a polynomial function. Consider input data set as  $\mathcal{V} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ , and non-parallel with each other, i.e.,  $\mathbf{v}_j \notin \text{span}(\mathbf{v}_k)$  for any  $j \neq k$ , we define

$$\mathbf{G}(\mathcal{V})_{ij} := \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\sigma(\mathbf{w}^\top \mathbf{v}_i) \sigma(\mathbf{w}^\top \mathbf{v}_j)], \tag{B.5}$$

then  $\lambda_{\min}(\mathbf{G}(\mathcal{V})) > 0$ .

Similar to Lemma B.1, we have Lemma B.2

**Lemma B.2.** Assume  $\sigma(\cdot)$  is analytic and not a polynomial function. Consider input data set as  $\mathcal{V} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ , and non-parallel with each other, i.e.,  $\mathbf{v}_j \notin \text{span}(\mathbf{v}_k)$  for any  $j \neq k$ , we define

$$\mathbf{G}(\mathcal{V})_{ij} := \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \sigma^{(1)}(\mathbf{w}^\top \mathbf{v}_i) \sigma^{(1)}(\mathbf{w}^\top \mathbf{v}_j) (\mathbf{v}_i^\top \mathbf{v}_j) \right], \tag{B.6}$$

then  $\lambda_{\min}(\mathbf{G}(\mathcal{V})) > 0$ .

Now we proceed to quantify the least eigenvalues of these Gram matrices.

### B.1 Full rankness for $(L + 1)$ -th Gram matrix

We begin this part by a lemma on the estimate of the entry of Gram matrices,

**Lemma B.3.** Given the input samples  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $\|\mathbf{x}_i\|_2 = 1$ ,  $1 \leq i \leq n$ , and  $\mathbf{x}_i \not\parallel \mathbf{x}_j$ , for any  $i \neq j$ , then for every fixed  $l$ , where  $1 \leq l \leq L$ , each diagonal entry of  $\tilde{\mathbf{K}}^{[l]}$  is the same with each other. Also for every fixed  $l$ , where  $1 \leq l \leq L$ , each element of the vector  $\tilde{\mathbf{b}}^{[l]}$  is the same with each other, i.e.,

$$\tilde{\mathbf{K}}_{ii}^{[l_1]} = \tilde{\mathbf{K}}_{jj}^{[l_1]}, \quad \tilde{\mathbf{b}}_i^{[l_2]} = \tilde{\mathbf{b}}_j^{[l_2]}, \quad i \neq j.$$

Moreover

$$\left( 1 - \frac{l}{L} \frac{c}{\sqrt{c_\sigma}} \right)^2 \leq \tilde{\mathbf{K}}_{ii}^{[l]} \leq \left( 1 + \frac{l}{L} \frac{c}{\sqrt{c_\sigma}} \right)^2, \tag{B.7}$$

and

$$\left( \tilde{\mathbf{b}}_i^{[l]} \right)^2 < \tilde{\mathbf{K}}_{ii}^{[l]}, \tag{B.8}$$

where  $c > 0$  and only depends on  $c_{\text{res}}$  and the activation function  $\sigma(\cdot)$ .

*Proof.* We shall prove it by induction on  $l$ . Firstly, we notice that  $\tilde{\mathbf{K}}_{ii}^{[0]} = \tilde{\mathbf{K}}_{jj}^{[0]}$  for any  $i \neq j$ , this is obvious because  $\|\mathbf{x}_i\|_2 = 1$ , then  $\tilde{\mathbf{K}}_{ii}^{[0]} = \tilde{\mathbf{K}}_{jj}^{[0]} = 1$ . Next we show that it holds true for  $l=1$ .

Since based on definition, recall that  $c_\sigma = (\mathbb{E}_{x \sim \mathcal{N}(0,1)} [\sigma(x)^2])^{-1}$ ,

$$\mathbf{K}_{ii}^{[1]} = c_\sigma \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[0]})} (\sigma(u)^2) = c_\sigma \mathbb{E}_{u \sim \mathcal{N}(0,1)} (\sigma(u)^2) = 1,$$

and

$$\tilde{\mathbf{b}}_i^{[1]} = \sqrt{c_\sigma} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[0]})} [\sigma(u)] = \sqrt{c_\sigma} \mathbb{E}_{u \sim \mathcal{N}(0,1)} [\sigma(u)],$$

then

$$\left( \tilde{\mathbf{b}}_i^{[1]} \right)^2 = c_\sigma \left( \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[0]})} [\sigma(u)] \right)^2 < 1,$$

the last inequality holds because

$$\left( \mathbb{E}_{x \sim \mathcal{N}(0,1)} [\sigma(x)] \right)^2 < \mathbb{E}_{x \sim \mathcal{N}(0,1)} [\sigma(x)^2].$$

Since the quantity is independent of our choice of  $i$ ,  $\tilde{\mathbf{K}}_{ii}^{[1]} = \tilde{\mathbf{K}}_{jj}^{[1]}$ ,  $\tilde{\mathbf{b}}_i^{[1]} = \tilde{\mathbf{b}}_j^{[1]}$ , for any  $i \neq j$ .

Now we assume that it holds for superscript being  $1, 2, \dots, l-1$  and want to show that it holds for  $l$ . Hence based on definition

$$\begin{aligned} \tilde{\mathbf{K}}_{ii}^{[l]} &= \tilde{\mathbf{K}}_{ii}^{[l-1]} + \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} \left[ \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[l-1]} \sigma(u)}{L} + \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[l-1]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(u)}{L^2} \right], \\ \tilde{\mathbf{b}}_i^{[l]} &= \tilde{\mathbf{b}}_i^{[l-1]} + \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)], \end{aligned}$$

such quantities are also independent of our choice of  $i$ .

Moreover we would like to show that (B.7) and (B.8) hold for all  $l$ .

Firstly, for  $\tilde{\mathbf{b}}_i^{[l]}$ , assume (B.8) holds for  $1, 2, \dots, l-1$ , then we have

$$\begin{aligned} \left( \tilde{\mathbf{b}}_i^{[l]} \right)^2 &= \left( \tilde{\mathbf{b}}_i^{[l-1]} \right)^2 + 2\tilde{\mathbf{b}}_i^{[l-1]} \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)] + \left( \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)] \right)^2 \\ &< \left( \tilde{\mathbf{b}}_i^{[l-1]} \right)^2 + 2\tilde{\mathbf{b}}_i^{[l-1]} \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)] + \frac{c_{\text{res}}^2}{L^2} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)^2] \\ &< \tilde{\mathbf{K}}_{ii}^{[l-1]} + 2\tilde{\mathbf{b}}_i^{[l-1]} \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)] + \frac{c_{\text{res}}^2}{L^2} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)^2] = \tilde{\mathbf{K}}_{ii}^{[l]}, \end{aligned}$$

showing that (B.8) holds for  $l$ .

For  $\tilde{\mathbf{K}}_{ii}^{[l]}$ , we have

$$\begin{aligned} & \left( \sqrt{\tilde{\mathbf{K}}_{ii}^{[l-1]}} - \frac{c_{\text{res}}}{L} \sqrt{\mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)^2]} \right)^2 \\ & \leq \tilde{\mathbf{K}}_{ii}^{[l]} \leq \left( \sqrt{\tilde{\mathbf{K}}_{ii}^{[l-1]}} + \frac{c_{\text{res}}}{L} \sqrt{\mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)^2]} \right)^2. \end{aligned} \tag{B.9}$$

Since  $\sigma(\cdot)$  is  $C_{\text{Lip}}$ -Lipschitz, we have for any  $1/2 \leq \alpha \leq 2$

$$\begin{aligned} & \left| \mathbb{E}_{X \sim \mathcal{N}(0,1)} [\sigma(\alpha X)^2] - \mathbb{E}_{X \sim \mathcal{N}(0,1)} [\sigma(X)^2] \right| \\ & \leq \mathbb{E}_{X \sim \mathcal{N}(0,1)} [|\sigma(\alpha X)^2 - \sigma(X)^2|] \\ & \leq C_{\text{Lip}} |\alpha - 1| \mathbb{E}_{X \sim \mathcal{N}(0,1)} [|X(\sigma(\alpha X) + \sigma(X))|] \\ & \leq C_{\text{Lip}} |\alpha - 1| \mathbb{E}_{X \sim \mathcal{N}(0,1)} [|X| |2\sigma(0)|] + C_{\text{Lip}} |\alpha + 1| \mathbb{E}_{X \sim \mathcal{N}(0,1)} [X^2] \\ & = C_{\text{Lip}} |\alpha - 1| \left( |2\sigma(0)| \sqrt{\frac{2}{\pi}} + C_{\text{Lip}} |\alpha + 1| \right) \\ & \leq \frac{C}{c_\sigma} |\alpha - 1|. \end{aligned}$$

Then

$$\mathbb{E}_{X \sim \mathcal{N}(0,1)} [\sigma(\alpha X)^2] \leq \frac{1}{c_\sigma} + \frac{C}{c_\sigma} |\alpha - 1|.$$

By induction

$$1 - \frac{l-1}{L} \frac{c}{\sqrt{c_\sigma}} \leq \sqrt{\tilde{\mathbf{K}}_{ii}^{[l-1]}} \leq 1 + \frac{l-1}{L} \frac{c}{\sqrt{c_\sigma}},$$

setting  $\alpha = \sqrt{\tilde{\mathbf{K}}_{ii}^{[l-1]}}$ , we obtain

$$\mathbb{E}_{X \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(X)^2] \leq \frac{1}{c_\sigma} + \frac{C}{c_\sigma} \frac{l-1}{L} \frac{c}{\sqrt{c_\sigma}}.$$

Then if we choose  $c$  wisely, for example, letting

$$c = \frac{C c_{\text{res}}^2}{2\sqrt{c_\sigma}} + \sqrt{\frac{C^2 c_{\text{res}}^4}{4c_\sigma} + c_{\text{res}}^2},$$

by our choice of  $c$ , combined with (B.9)

$$\left( \sqrt{\tilde{\mathbf{K}}_{ii}^{[l-1]}} - \frac{1}{L} \frac{c}{\sqrt{c_\sigma}} \right)^2 \leq \tilde{\mathbf{K}}_{ii}^{[l]} \leq \left( \sqrt{\tilde{\mathbf{K}}_{ii}^{[l-1]}} + \frac{1}{L} \frac{c}{\sqrt{c_\sigma}} \right)^2,$$

and hence

$$\left( 1 - \frac{l}{L} \frac{c}{\sqrt{c_\sigma}} \right)^2 \leq \tilde{\mathbf{K}}_{ii}^{[l-1]} \leq \left( 1 + \frac{l}{L} \frac{c}{\sqrt{c_\sigma}} \right)^2,$$

which finishes our proof. □

Our next lemma is crucial in that it reveals a ‘covariance-type’ structure for the Gram matrices. We need to introduce a standard notation related to matrices. We denote that  $A \succeq B$  if and only if  $A - B$  is a semi-positive definite matrix, and  $A \succ B$  if and only if  $A - B$  is a strictly positive definite matrix.

**Proposition B.1.** *Given the input samples  $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$ ,  $\|x_i\|_2 = 1$ , for  $1 \leq i \leq n$ , and  $x_i \perp x_j$ ,  $i \neq j$ , then we have for every fixed  $l$ , where  $1 \leq l \leq L$ ,*

$$\tilde{\mathbf{K}}^{[l]} - \tilde{\mathbf{b}}^{[l]} \otimes (\tilde{\mathbf{b}}^{[l]})^\top \succ \tilde{\mathbf{K}}^{[l-1]} - \tilde{\mathbf{b}}^{[l-1]} \otimes (\tilde{\mathbf{b}}^{[l-1]})^\top. \tag{B.10}$$

Moreover, since

$$\tilde{\mathbf{K}}^{[1]} - \tilde{\mathbf{b}}^{[1]} \otimes (\tilde{\mathbf{b}}^{[1]})^\top \succ 0,$$

we denote that

$$\lambda_{\min} \left( \tilde{\mathbf{K}}^{[1]} - \tilde{\mathbf{b}}^{[1]} \otimes (\tilde{\mathbf{b}}^{[1]})^\top \right) = \lambda_0, \tag{B.11}$$

then we can conclude that for  $2 \leq l \leq L$ ,

$$\lambda_{\min} \left( \tilde{\mathbf{K}}^{[l]} \right) > \lambda_0, \tag{B.12}$$

where  $\lambda_0$  only depends on the activation function and input data and independent of depth  $L$ .

*Proof.* We only need to show that for  $1 \leq i, j \leq n$  and  $1 \leq l \leq L$

$$\begin{aligned} & \tilde{\mathbf{K}}_{ij}^{[l]} - \tilde{\mathbf{b}}_i^{[l]} \tilde{\mathbf{b}}_j^{[l]} \\ &= \tilde{\mathbf{K}}_{ij}^{[l-1]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N} \left( \mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[l-1]} & \tilde{\mathbf{K}}_{ij}^{[l-1]} \\ \tilde{\mathbf{K}}_{ji}^{[l-1]} & \tilde{\mathbf{K}}_{jj}^{[l-1]} \end{pmatrix} \right)} \left[ \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[l-1]} \sigma(v)}{L} + \frac{c_{\text{res}} \tilde{\mathbf{b}}_j^{[l-1]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(v)}{L^2} \right] \\ & \quad - \left( \tilde{\mathbf{b}}_i^{[l-1]} + \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)] \right) \left( \tilde{\mathbf{b}}_j^{[l-1]} + \frac{c_{\text{res}}}{L} \mathbb{E}_{v \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{jj}^{[l-1]})} [\sigma(v)] \right) \\ &= \tilde{\mathbf{K}}_{ij}^{[l-1]} - \tilde{\mathbf{b}}_i^{[l-1]} \tilde{\mathbf{b}}_j^{[l-1]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N} \left( \mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[l-1]} & \tilde{\mathbf{K}}_{ij}^{[l-1]} \\ \tilde{\mathbf{K}}_{ji}^{[l-1]} & \tilde{\mathbf{K}}_{jj}^{[l-1]} \end{pmatrix} \right)} \left[ \frac{c_{\text{res}}^2 \sigma(u) \sigma(v)}{L^2} \right] \\ & \quad - \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)] \frac{c_{\text{res}}}{L} \mathbb{E}_{v \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{jj}^{[l-1]})} [\sigma(v)] \\ &= \tilde{\mathbf{K}}_{ij}^{[l-1]} - \tilde{\mathbf{b}}_i^{[l-1]} \tilde{\mathbf{b}}_j^{[l-1]} + \frac{c_{\text{res}}^2}{L^2} \text{Cov}_{(u,v)^\top \sim \mathcal{N} \left( \mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[l-1]} & \tilde{\mathbf{K}}_{ij}^{[l-1]} \\ \tilde{\mathbf{K}}_{ji}^{[l-1]} & \tilde{\mathbf{K}}_{jj}^{[l-1]} \end{pmatrix} \right)} [\sigma(u) \sigma(v)], \end{aligned}$$

which brings us to the definition of a series of covariance matrices  $\{\mathbf{P}^{[s]}: 1 \leq s \leq L\}$ , where

$$\mathbf{P}_{ij}^{[s]} := \frac{c_{\text{res}}^2}{L^2} \text{Cov}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[s]} & \tilde{\mathbf{K}}_{ij}^{[s]} \\ \tilde{\mathbf{K}}_{ji}^{[s]} & \tilde{\mathbf{K}}_{jj}^{[s]} \end{pmatrix}\right)} [\sigma(u)\sigma(v)], \quad 1 \leq s \leq L,$$

$\mathbf{P}^{[s]}$  are covariance matrices, naturally we have  $\mathbf{P}^{[s]} \succeq 0$ , and  $\mathbf{P}^{[s]} \succ 0$  except that one sample is an exact linear function of the others. Apply Lemma B.1 directly, we can guarantee that  $\mathbf{P}^{[s]}$  is positive definite for every  $s$ . Hence, inductively we have

$$\begin{aligned} \tilde{\mathbf{K}}^{[l]} &\succeq \tilde{\mathbf{K}}^{[l]} - \tilde{\mathbf{b}}^{[l]} \otimes (\tilde{\mathbf{b}}^{[l]})^\top \\ &= \tilde{\mathbf{K}}^{[l-1]} - \tilde{\mathbf{b}}^{[l-1]} \otimes (\tilde{\mathbf{b}}^{[l-1]})^\top + \mathbf{P}^{[l-1]} \\ &\succ \tilde{\mathbf{K}}^{[l-1]} - \tilde{\mathbf{b}}^{[l-1]} \otimes (\tilde{\mathbf{b}}^{[l-1]})^\top \\ &= \tilde{\mathbf{K}}^{[l-2]} - \tilde{\mathbf{b}}^{[l-2]} \otimes (\tilde{\mathbf{b}}^{[l-2]})^\top + \mathbf{P}^{[l-2]} \\ &\quad \vdots \\ &\succ \tilde{\mathbf{K}}^{[1]} - \tilde{\mathbf{b}}^{[1]} \otimes (\tilde{\mathbf{b}}^{[1]})^\top, \end{aligned}$$

the last line brings us to the entry of  $\tilde{\mathbf{K}}^{[1]} - \tilde{\mathbf{b}}^{[1]} \otimes (\tilde{\mathbf{b}}^{[1]})^\top$ , where

$$\left(\tilde{\mathbf{K}}^{[1]} - \tilde{\mathbf{b}}^{[1]} \otimes (\tilde{\mathbf{b}}^{[1]})^\top\right)_{ij} = c_\sigma \text{Cov}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[0]} & \tilde{\mathbf{K}}_{ij}^{[0]} \\ \tilde{\mathbf{K}}_{ji}^{[0]} & \tilde{\mathbf{K}}_{jj}^{[0]} \end{pmatrix}\right)} [\sigma(u)\sigma(v)].$$

Then apply Lemma B.1 again

$$\lambda_{\min}\left(\tilde{\mathbf{K}}^{[1]} - \tilde{\mathbf{b}}^{[1]} \otimes (\tilde{\mathbf{b}}^{[1]})^\top\right) = \lambda_0 > 0,$$

and  $\lambda_0$  only depends on the input data and activation function. □

**Corollary B.1.** *Given the input samples  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $\|\mathbf{x}_i\|_2 = 1$ , for  $1 \leq i \leq n$ , and  $\mathbf{x}_i \not\parallel \mathbf{x}_j$ ,  $i \neq j$ , then we have*

$$\lambda_{\min}\left(\mathbf{K}^{[L+1]}\right) > \lambda_0, \tag{B.13}$$

where  $\lambda_0$  has been defined in (B.11).

*Proof.* The proof is quite similar to the proof of Proposition B.1. Recall that

$$\mathbf{K}_{ij}^{[L+1]} = \tilde{\mathbf{K}}_{ij}^{[L]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[L+1]}\right)} \left[ \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[L]} \sigma(v)}{L} + \frac{c_{\text{res}} \tilde{\mathbf{b}}_j^{[L]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u)\sigma(v)}{L^2} \right],$$



and we define that

$$\mathbf{b}_i^{[L+1]} := \tilde{\mathbf{b}}_i^{[L]} + \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[L]})} [\sigma(u)].$$

Then

$$\mathbf{K}_{ij}^{[L+1]} - \mathbf{b}_i^{[L+1]} \mathbf{b}_j^{[L+1]} = \tilde{\mathbf{K}}_{ij}^{[L]} - \tilde{\mathbf{b}}_i^{[L]} \tilde{\mathbf{b}}_j^{[L]} + \frac{c_{\text{res}}^2}{L^2} \text{Cov}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[L]} & \tilde{\mathbf{K}}_{ij}^{[L]} \\ \tilde{\mathbf{K}}_{ji}^{[L]} & \tilde{\mathbf{K}}_{jj}^{[L]} \end{pmatrix}\right)} [\sigma(u)\sigma(v)].$$

Hence

$$\begin{aligned} \mathbf{K}^{[L+1]} &\succeq \mathbf{K}^{[L+1]} - \mathbf{b}^{[L+1]} \otimes (\mathbf{b}^{[L+1]})^\top \\ &\succ \tilde{\mathbf{K}}^{[L]} - \tilde{\mathbf{b}}^{[L]} \otimes (\tilde{\mathbf{b}}^{[L]})^\top. \end{aligned}$$

Applying Proposition B.1 directly, we are able to finish the proof.  $\square$

By Corollary B.1, we see that  $\lambda_{\min}(\mathbf{K}^{[L+1]}) \sim \Omega(1)$ .

## B.2 Full rankness for $L$ -th Gram matrix

Our next Proposition is related to the eigenvalue of the  $L$ -th Gram matrix, whose entries concerning the derivative of the activation function. This proposition has been stated as Proposition F.2 in Du et al. [12], and we will mimic its proof.

**Proposition B.2.** *Given the input samples  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $\|\mathbf{x}_i\|_2 = 1$ , for  $1 \leq i \leq n$ , and  $\mathbf{x}_i \perp \mathbf{x}_j$ ,  $i \neq j$ , then for  $2 \leq l \leq L$*

$$\lambda_{\min}(\mathbf{K}^{[l]}) \geq \frac{c_{\text{res}}^2}{L^2} \kappa, \quad (\text{B.14})$$

where  $\kappa$  is a constant that only depends on  $\sigma(\cdot)$  and input samples, independent of depth  $L$ .

*Proof.* Based on Lemma B.3, uniformly for any  $1 \leq l \leq L$ ,

$$1/c \leq \tilde{\mathbf{K}}_{ii}^{[l]} \leq c,$$

then we can define a function  $\mathbf{G}: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$ , such that

$$\mathbf{G}(\mathbf{K})_{ij} := \mathbf{K}_{ij} \mathbb{E}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \mathbf{K}_{ii} & \mathbf{K}_{ij} \\ \mathbf{K}_{ji} & \mathbf{K}_{jj} \end{pmatrix}\right)} \sigma^{(1)}(u)\sigma^{(1)}(v).$$

Consequently, a scalar function  $g(\lambda)$  can be defined as follows:

$$g(\lambda) := \min_{\mathbf{K}: \mathbf{K} \succ 0, 1/c \leq \mathbf{K}_{ii} \leq c, \lambda_{\min}(\mathbf{K}) \geq \lambda} \lambda_{\min}(\mathbf{G}(\mathbf{K})).$$

Then Lemma B.2 guarantees that

$$g(\lambda_0) > 0.$$

Moreover, based on Proposition B.1

$$\lambda_{\min}(\tilde{\mathbf{K}}^{[L-1]}) > \lambda_0,$$

we have

$$\lambda_{\min}(\mathbf{K}^{[L]}) \geq \frac{c_{\text{res}}^2}{L^2} g(\lambda_0). \tag{B.15}$$

Let  $\kappa = g(\lambda_0)$ . Since  $\kappa$  is independent of depth  $L$ , we finish our proof. □

By Proposition B.2, we see that  $\lambda_{\min}(\mathbf{K}^{[L]}) \sim \Omega(\frac{1}{L^2})$ .

### C Random initialization of Gram matrices

In this part, we are going to show that with high probability with respect to the random initialization,

$$\lambda_{\min} \left[ \mathcal{G}_t^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} > \frac{3\lambda_0}{4},$$

where  $\lambda_0$  is defined in (B.11).

Let's get started with a lemma concerning the Gaussian concentrations.

**Lemma C.1.** *Let  $\mathbf{X} = (X_1, \dots, X_p) \in \mathbb{R}^p$ ,  $X_1, \dots, X_p$  be a vector of i.i.d. Gaussian variables from  $\mathcal{N}(0, \sigma^2)$ , and let  $f(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}$  be  $L$ -Lipschitz function, i.e.,  $|f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_2$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ , then for any  $t \geq 0$*

$$\mathbb{P}(|f(\mathbf{X}) - \mathbb{E}f(\mathbf{X})| \geq t) \leq 2\exp\left(-\frac{t^2}{2L^2\sigma^2}\right). \tag{C.1}$$

Before we proceed to the stability of the randomly initialized Gram matrix of higher order, we need to state two lemmas. The first lemma has been stated as Lemma G.3. in Du et al. [12],

**Lemma C.2.** *If  $\sigma(\cdot)$  is  $C_{\text{Lip}}$ -Lipschitz, then for  $a, b \in \mathbb{R}^+$ , with  $1/c \leq \min(a, b), \max(a, b) \leq c$  for some  $c > 0$ , then we have*

$$\left| \mathbb{E}_{z \sim \mathcal{N}(0,1)}[\sigma(az)] - \mathbb{E}_{z \sim \mathcal{N}(0,1)}[\sigma(bz)] \right| \leq C|a - b|, \tag{C.2}$$

where  $C > 0$  only depends on  $c$  and Lipschitz constant  $C_{\text{Lip}}$ .

Next lemma has been stated as Lemma G.4. in Du et al. [12]:

**Lemma C.3.** *If  $\sigma(\cdot)$  is  $C_{\text{Lip}}$ -Lipschitz, define a scalar function  $F(\mathbf{K})$  as follows:*

$$F(\mathbf{K}) = \mathbb{E}_{(u,v) \top \sim \mathcal{N}(\mathbf{0}, \mathbf{K})} [\sigma(u)\sigma(v)],$$

then for any two matrices  $\mathbf{A}, \mathbf{B}$  being

$$\mathbf{A} = \begin{pmatrix} a_1^2 & \rho_1 a_1 b_1 \\ \rho_1 a_1 b_1 & b_1^2 \end{pmatrix},$$

$$\mathbf{B} = \begin{pmatrix} a_2^2 & \rho_2 a_2 b_2 \\ \rho_2 a_2 b_2 & b_2^2 \end{pmatrix},$$

and their entries satisfying

$$1/c \leq \min(a_1, b_1), \min(a_2, b_2), \max(a_1, b_1), \max(a_2, b_2) \leq c,$$

and

$$-1 < \rho_1, \rho_2 < 1$$

for some  $c > 0$ , then we have

$$|F(\mathbf{A}) - F(\mathbf{B})| \leq C \|\mathbf{A} - \mathbf{B}\|_F \leq 2C \|\mathbf{A} - \mathbf{B}\|_\infty,$$

where the constant  $C > 0$  only relies on  $c$  and the Lipschitz constant  $C_{\text{Lip}}$ .

We shall begin with a proposition on the initial estimate of the output of each layer  $\mathbf{x}_i^{[l]}(0)$ ,

**Proposition C.1.** *Under Assumptions 4.1 and 4.2, we have that for some  $t > 0$ ,  $1 \leq i \leq n$ ,  $1 \leq l \leq L$ ,*

$$\mathbb{P} \left( \left| \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 - \sqrt{\tilde{\mathbf{K}}_{ii}^{[l]}} \right| \geq t \right) \leq \exp(-cmt^2), \tag{C.3}$$

$$\mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \tilde{\mathbf{b}}_i^{[l]} \right| \geq t \right) \leq \exp(-cmt^2), \tag{C.4}$$

where  $c > 0$  is a constant independent of depth  $L$ .

*Proof.* For  $l = 1$ , we have

$$\left\| \mathbf{x}_i^{[1]}(0) \right\|_2^2 = \frac{c_\sigma}{m} \sum_{j=1}^m \left( \sigma(\mathbf{W}^{[1]}(0) \mathbf{x}_i)_j \right)^2,$$

then

$$\mathbb{E} \left[ \left\| \mathbf{x}_i^{[1]}(0) \right\|_2^2 \right] = c_\sigma \mathbb{E}_{x \sim \mathcal{N}(0,1)} [\sigma(x)^2] = \tilde{\mathbf{K}}_{ii}^{[1]} = 1.$$

Since  $(\mathbf{W}^{[1]}(0)\mathbf{x}_i)_j$  are i.i.d standard Gaussian variables, and  $\sigma(\cdot)$  is 1-Lipschitz, then  $(\sigma(\mathbf{W}^{[1]}(0)\mathbf{x}_i)_j)$  are sub-exponential variables, and we have for  $\lambda > 0$ ,

$$\mathbb{E} \left[ \exp \lambda \left( m \|\mathbf{x}_i^{[1]}(0)\|_2^2 - m \right) \right] \leq \exp (cm\lambda^2).$$

Hence applying Markov inequality directly

$$\mathbb{P} \left( \left| \|\mathbf{x}_i^{[1]}(0)\|_2 - \sqrt{\tilde{\mathbf{K}}_i^{[1]}} \right| \geq t \right) \leq \mathbb{P} \left( \left| \|\mathbf{x}_i^{[1]}(0)\|_2^2 - 1 \right| \geq 2t \right) \leq \exp (-cmt^2),$$

and

$$\left\langle \frac{\mathbf{x}_i^{[1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle = \frac{\sqrt{c_\sigma}}{m} \sum_{j=1}^m \left( \sigma(\mathbf{W}^{[1]}(0)\mathbf{x}_i)_j \right),$$

we obtain

$$\mathbb{E} \left[ \left\langle \frac{\mathbf{x}_i^{[1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle \right] = \tilde{\mathbf{b}}_i^{[1]}.$$

We should note that  $\mathbf{x}_i^{[1]}(0)$  writes into

$$\mathbf{x}_i^{[1]}(0) = \sqrt{\frac{c_\sigma}{m}} \sigma(\mathbf{X}),$$

with  $\mathbf{X}$  being a standard normal Gaussian vector. We shall focus on the inner product function  $g^{[1]}(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$ , with

$$g^{[1]}(\mathbf{X}) = \frac{\sqrt{c_\sigma}}{m} \langle \sigma(\mathbf{X}), \mathbf{1} \rangle.$$

We have for any  $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^m$ ,

$$\begin{aligned} \left| g^{[1]}(\mathbf{X}_1) - g^{[1]}(\mathbf{X}_2) \right| &\leq \left| \frac{\sqrt{c_\sigma}}{m} \langle \sigma(\mathbf{X}_1), \mathbf{1} \rangle - \frac{\sqrt{c_\sigma}}{m} \langle \sigma(\mathbf{X}_2), \mathbf{1} \rangle \right| \\ &\leq \frac{\sqrt{c_\sigma}}{m} \langle |\mathbf{X}_1 - \mathbf{X}_2|, \mathbf{1} \rangle \leq \sqrt{\frac{c_\sigma}{m}} \|\mathbf{X}_1 - \mathbf{X}_2\|_2. \end{aligned}$$

Hence  $g^{[1]}(\cdot)$  is  $\frac{c}{\sqrt{m}}$ -Lipschitz, then apply Lemma C.1

$$\mathbb{P} \left( \left| g^{[1]}(\mathbf{X}) - \mathbb{E}g^{[1]}(\mathbf{X}) \right| \geq t \right) \leq \exp(-cmt^2).$$

Then we have

$$\mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \tilde{\mathbf{b}}_i^{[1]} \right| \geq t \right) \leq \exp(-cmt^2).$$

Our next step is to prove that (C.3) and (C.4) hold for  $l \geq 2$ , and we will prove it by induction.

Assume that (C.3) and (C.4) hold for  $1, 2, 3, \dots, l$  and want to show that they hold for  $l+1$ .

$$\mathbb{P} \left( \left| \left\| \mathbf{x}_i^{[l+1]}(0) \right\|_2 - \sqrt{\tilde{\mathbf{K}}_{ii}^{[l+1]}} \right| \geq t \right) \leq \exp(-cmt^2), \tag{C.5}$$

$$\mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l+1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \tilde{\mathbf{b}}_i^{[l+1]} \right| \geq t \right) \leq \exp(-cmt^2). \tag{C.6}$$

We recall that,

$$\mathbf{x}_i^{[l+1]}(0) = \mathbf{x}_i^{[l]}(0) + \frac{c_{\text{res}}}{L\sqrt{m}} \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right),$$

and the definition of  $\tilde{\mathbf{K}}_{ii}^{[l]}$  and  $\tilde{\mathbf{b}}_i^{[l]}$

$$\begin{aligned} \tilde{\mathbf{K}}_{ii}^{[l+1]} &= \tilde{\mathbf{K}}_{ii}^{[l]} + \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l]})} \left[ \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[l]} \sigma(u)}{L} + \frac{c_{\text{res}} \tilde{\mathbf{b}}_i^{[l]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(u)}{L^2} \right], \\ \tilde{\mathbf{b}}_i^{[l+1]} &= \tilde{\mathbf{b}}_i^{[l]} + \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[l]})} [\sigma(u)]. \end{aligned}$$

Then we have

$$\begin{aligned} \left\| \mathbf{x}_i^{[l+1]}(0) \right\|_2^2 &= \left\| \mathbf{x}_i^{[l]}(0) \right\|_2^2 + 2 \frac{c_{\text{res}}}{L} \underbrace{\left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right) \right\rangle}_{\text{I}} \\ &\quad + \frac{c_{\text{res}}^2}{L^2} \underbrace{\frac{1}{m} \left\langle \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right), \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right) \right\rangle}_{\text{II}}. \end{aligned}$$

Then we need to focus on the terms I and II. Note that for term I there is a  $\frac{1}{\sqrt{m}}$  scaling factor contained in  $\mathbf{x}_i^{[l]}(0)$ , and  $\sigma(\mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0))$  has distribution

$$\sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right) \sim \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y} \right),$$

with  $\mathbf{Y}$  being a standard normal Gaussian vector. Then we have

$$\mathbb{E} \left[ \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right) \right\rangle \right] = \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \mathbb{E} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right] \right\rangle.$$

Focusing on the inner product function  $g^{[l]}(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$ , with

$$g^{[l]}(\mathbf{Y}) = \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right\rangle,$$

we have for any  $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^m$ ,

$$\begin{aligned} \left| g^{[l]}(\mathbf{Y}_1) - g^{[l]}(\mathbf{Y}_2) \right| &\leq \frac{1}{\sqrt{m}} \left| \left\langle \mathbf{x}_i^{[l]}(0), \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y}_1 \right) \right\rangle - \left\langle \mathbf{x}_i^{[l]}(0), \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y}_2 \right) \right\rangle \right| \\ &\leq \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 |\mathbf{Y}_1 - \mathbf{Y}_2| \right\rangle \leq \frac{1}{\sqrt{m}} \left\| \mathbf{x}_i^{[l]}(0) \right\|_2^2 \|\mathbf{Y}_1 - \mathbf{Y}_2\|_2^2 \end{aligned}$$

based on our induction hypothesis,  $\left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \leq C$  with high probability. Hence  $g^{[l]}(\cdot)$  is  $\frac{C}{\sqrt{m}}$ -Lipschitz. Applying Lemma C.1 again

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right) \right\rangle - \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \mathbb{E} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right] \right\rangle \right| \geq t \right) \\ &\leq \exp(-cmt^2), \end{aligned} \tag{C.7}$$

and based on our induction hypothesis,

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \mathbb{E} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right] \right\rangle - \tilde{\mathbf{b}}_i^{[l]} \mathbb{E} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right] \right| \geq t \right) \\ &\leq \exp(-cmt^2), \end{aligned} \tag{C.8}$$

from Lemma C.2

$$\left| \mathbb{E} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right] - \mathbb{E} \left[ \sigma \left( \sqrt{\tilde{\mathbf{K}}_{ii}^{[l]}} \mathbf{Y} \right) \right] \right| \leq C \left| \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 - \sqrt{\tilde{\mathbf{K}}_{ii}^{[l]}} \right|,$$

altogether we have

$$\mathbb{P} \left( \left| \tilde{\mathbf{b}}_i^{[l]} \mathbb{E} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right] - \tilde{\mathbf{b}}_i^{[l]} \mathbb{E} \left[ \sigma \left( \sqrt{\tilde{\mathbf{K}}_{ii}^{[l]}} \mathbf{Y} \right) \right] \right| \geq t \right) \leq \exp(-cmt^2). \tag{C.9}$$

Combining (C.7), (C.8) and (C.9), we have

$$\mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right) \right\rangle - \tilde{\mathbf{b}}_i^{[l]} \mathbb{E} \left[ \sigma \left( \sqrt{\tilde{\mathbf{K}}_{ii}^{[l]}} \mathbf{Y} \right) \right] \right| \geq t \right) \leq \exp(-cmt^2). \tag{C.10}$$

Finally for term II

$$\mathbb{E} \left[ \frac{1}{m} \left\langle \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right), \sigma \left( \mathbf{W}^{[l+1]}(0) \mathbf{x}_i^{[l]}(0) \right) \right\rangle \right] = \mathbb{E}_{x \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 x \right)^2 \right],$$

since  $(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0))$  are i.i.d standard Gaussian variables, and  $\sigma(\cdot)$  is 1-Lipschitz, then  $(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0))$  are sub-exponential variables and we have

$$\begin{aligned} & \mathbb{P}\left(\left|\frac{1}{m}\langle\sigma\left(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0)\right),\sigma\left(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0)\right)\rangle-\mathbb{E}_{x\sim\mathcal{N}(0,1)}\left[\sigma\left(\left\|\mathbf{x}_i^{[l]}(0)\right\|_2x\right)^2\right]\right|\geq t\right) \\ & \leq \exp(-cmt^2). \end{aligned} \tag{C.11}$$

Applying Lemma C.3

$$\left|\mathbb{E}_{x\sim\mathcal{N}(0,1)}\left[\sigma\left(\left\|\mathbf{x}_i^{[l]}(0)\right\|_2x\right)^2\right]-\mathbb{E}_{x\sim\mathcal{N}(0,1)}\left[\sigma\left(\sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}}x\right)^2\right]\right|\leq C\left|\left\|\mathbf{x}_i^{[l]}(0)\right\|_2-\sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}}\right|,$$

then based on our induction hypothesis, we have

$$\mathbb{P}\left(\left|\mathbb{E}_{x\sim\mathcal{N}(0,1)}\left[\sigma\left(\left\|\mathbf{x}_i^{[l]}(0)\right\|_2x\right)^2\right]-\mathbb{E}_{x\sim\mathcal{N}(0,1)}\left[\sigma\left(\sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}}x\right)^2\right]\right|\geq t\right)\leq \exp(-cmt^2). \tag{C.12}$$

Combining (C.11) and (C.12)

$$\begin{aligned} & \mathbb{P}\left(\left|\frac{1}{m}\langle\sigma\left(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0)\right),\sigma\left(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0)\right)\rangle-\mathbb{E}_{x\sim\mathcal{N}(0,1)}\left[\sigma\left(\sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}}x\right)^2\right]\right|\geq t\right) \\ & \leq \exp(-cmt^2). \end{aligned} \tag{C.13}$$

Since we have

$$\begin{aligned} \left\|\mathbf{x}_i^{[l+1]}(0)\right\|_2^2 &= \left\|\mathbf{x}_i^{[l]}(0)\right\|_2^2 + \underbrace{2\frac{c_{\text{res}}}{L}\left\langle\frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}},\sigma\left(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0)\right)\right\rangle}_{\text{I}} \\ & \quad + \underbrace{\frac{c_{\text{res}}^2}{L^2}\frac{1}{m}\langle\sigma\left(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0)\right),\sigma\left(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0)\right)\rangle}_{\text{II}}, \end{aligned}$$

then

$$\begin{aligned} & \mathbb{P}\left(\left|\left\|\mathbf{x}_i^{[l+1]}(0)\right\|_2^2-\widetilde{\mathbf{K}}_{ii}^{[l+1]}\right|\geq t\left(1+\frac{c_{\text{res}}}{L}\right)^2\right) \\ & \leq \mathbb{P}\left(\left|\left\|\mathbf{x}_i^{[l]}(0)\right\|_2^2-\widetilde{\mathbf{K}}_{ii}^{[l]}\right|\geq t\right)+\mathbb{P}(\text{I}\geq t)+\mathbb{P}(\text{II}\geq t) \\ & \leq \exp(-cmt^2). \end{aligned} \tag{C.14}$$

We shall see that thanks to the  $\frac{c_{\text{res}}}{L}$  structure, with high probability the difference of  $\left| \|\mathbf{x}_i^{[l]}(0)\|_2^2 - \widetilde{\mathbf{K}}_{ii}^{[l]} \right|$  does not explode exponentially layer by layer.

For  $\widetilde{\mathbf{b}}_i^{[l+1]}$ , applying Lemma C.1, we have

$$\mathbb{P} \left( \left| \left\langle \frac{\sigma(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0))}{m}, \mathbf{1} \right\rangle - \mathbb{E} \left[ \sigma \left( \|\mathbf{x}_i^{[l]}(0)\|_2 \mathbf{Y} \right) \right] \right| \geq t \right) \leq \exp(-cmt^2). \quad (\text{C.15})$$

And applying Lemma C.2, we have

$$\mathbb{P} \left( \left| \mathbb{E} \left[ \sigma \left( \|\mathbf{x}_i^{[l]}(0)\|_2 \mathbf{Y} \right) \right] - \mathbb{E} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \mathbf{Y} \right) \right] \right| \geq t \right) \leq \exp(-cmt^2). \quad (\text{C.16})$$

Combining (C.15) and (C.16),

$$\mathbb{P} \left( \left| \left\langle \frac{\sigma(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0))}{m}, \mathbf{1} \right\rangle - \mathbb{E} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \mathbf{Y} \right) \right] \right| \geq t \right) \leq \exp(-cmt^2). \quad (\text{C.17})$$

Then

$$\begin{aligned} & \mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l+1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \widetilde{\mathbf{b}}_i^{[l+1]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right) \right) \\ & \leq \mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \widetilde{\mathbf{b}}_i^{[l]} \right| \geq t \right) + \mathbb{P} \left( \left| \left\langle \frac{\sigma(\mathbf{W}^{[l+1]}(0)\mathbf{x}_i^{[l]}(0))}{m}, \mathbf{1} \right\rangle - \mathbb{E} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \mathbf{Y} \right) \right] \right| \geq t \right) \\ & \leq \exp(-cmt^2). \end{aligned} \quad (\text{C.18})$$

We shall see again that thanks to the  $\frac{c_{\text{res}}}{L}$  structure, with high probability the difference of  $\left| \left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \widetilde{\mathbf{b}}_i^{[l]} \right|$  only has slight increment with respect to each layer  $l$ .  $\square$

Our next Proposition is on the least eigenvalue of the randomly initialized Gram matrix  $\mathbf{G}^{[1]}(0)$ .

**Proposition C.2.** *Under Assumptions 4.1 and 4.2, if  $m = \Omega\left(\left(\frac{\mu}{\lambda_0}\right)^{2+\epsilon}\right)$ , then with high probability*

$$\lambda_{\min} \left( \mathbf{G}^{[1]}(0) \right) \geq \frac{3\lambda_0}{4}, \quad (\text{C.19})$$

where  $\lambda_0$  has been defined in (B.11).



*Proof.* We have that

$$\begin{aligned} \mathbf{G}_{ij}^{[1]}(0) &= \langle \mathbf{x}_i^{[1]}(0), \mathbf{x}_j^{[1]}(0) \rangle, \\ \tilde{\mathbf{K}}_{ij}^{[0]} &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle, \\ \tilde{\mathbf{K}}_{ij}^{[1]} &= c_\sigma \mathbb{E}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[0]} & \tilde{\mathbf{K}}_{ij}^{[0]} \\ \tilde{\mathbf{K}}_{ji}^{[0]} & \tilde{\mathbf{K}}_{jj}^{[0]} \end{pmatrix}\right)} [\sigma(u)\sigma(v)]. \end{aligned}$$

Now we need to apply Lemma C.1 again, except that this time we are going to apply it to the inner product function  $h^{[1]}(\cdot) : \mathbb{R}^{2m} \rightarrow \mathbb{R}$ , with

$$h^{[1]}(\mathbf{Z}) = \frac{c_\sigma}{m} \left\langle \sigma(\mathbf{X}), \sigma(\rho\mathbf{X} + \sqrt{1-\rho^2}\mathbf{Y}) \right\rangle,$$

where  $-1 \leq \rho \leq 1$ .

Specifically with  $\mathbf{Z}^\top = (\mathbf{X}^\top, \mathbf{Y}^\top)$ , we have for any  $\mathbf{Z}_1, \mathbf{Z}_2 \in \mathbb{R}^m$ ,

$$\begin{aligned} |h^{[1]}(\mathbf{Z}_1) - h^{[1]}(\mathbf{Z}_2)| &\leq \sqrt{\frac{c_\sigma}{m}} \left\| \sigma(\rho\mathbf{X}_1 + \sqrt{1-\rho^2}\mathbf{Y}_1) \right\|_2 \sqrt{\frac{c_\sigma}{m}} \|\mathbf{X}_1 - \mathbf{X}_2\|_2 \\ &\quad + \sqrt{\frac{c_\sigma}{m}} \|\sigma(\mathbf{X}_2)\|_2 \sqrt{\frac{c_\sigma}{m}} \left( |\rho| \|\mathbf{X}_1 - \mathbf{X}_2\|_2 + \sqrt{1-\rho^2} \|\mathbf{Y}_1 - \mathbf{Y}_2\|_2 \right), \end{aligned}$$

combined with Proposition C.1, with probability  $1 - \exp(-cm)$ ,

$$\sqrt{\frac{c_\sigma}{m}} \left\| \sigma(\rho\mathbf{X}_1 + \sqrt{1-\rho^2}\mathbf{Y}_1) \right\|_2, \sqrt{\frac{c_\sigma}{m}} \|\sigma(\mathbf{X}_2)\|_2 \leq 2.$$

So we have

$$|h^{[1]}(\mathbf{Z}_1) - h^{[1]}(\mathbf{Z}_2)| \leq 8\sqrt{\frac{c_\sigma}{m}} \|\mathbf{Z}_1 - \mathbf{Z}_2\|_2,$$

hence  $h^{[1]}(\mathbf{Z})$  is  $8\sqrt{\frac{c_\sigma}{m}}$ -Lipschitz, then we shall set  $\rho = \tilde{\mathbf{K}}_{ij}^{[0]}$ ,

$$\mathbb{P}\left(\left|\mathbf{G}_{ij}^{[1]}(0) - \tilde{\mathbf{K}}_{ij}^{[1]}\right| \geq t\right) \leq \exp(-cmt^2). \tag{C.20}$$

Noting that

$$\left\| \mathbf{G}^{[1]}(0) - \tilde{\mathbf{K}}^{[1]} \right\|_{2 \rightarrow 2} \leq \left\| \mathbf{G}^{[1]}(0) - \tilde{\mathbf{K}}^{[1]} \right\|_{\text{F}} \leq n \left\| \mathbf{G}^{[1]}(0) - \tilde{\mathbf{K}}^{[1]} \right\|_{\infty},$$

based on Proposition B.1,  $\lambda_{\min}(\tilde{\mathbf{K}}^{[1]}) \geq \lambda_0$ , then if we choose  $t = \frac{\lambda_0}{4n}$  and with a union  $m^2$  such events, we have with probability  $1 - m^2 \exp(-cm\lambda_0^2/n^2)$

$$\left\| \mathbf{G}^{[1]}(0) - \tilde{\mathbf{K}}^{[1]} \right\|_{2 \rightarrow 2} \leq \frac{\lambda_0}{4}. \tag{C.21}$$

Hence if  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\varepsilon}\right)$ , we have with probability  $1 - \exp(-m^\varepsilon)$

$$\lambda_{\min}(\mathbf{G}^{[1]}(0)) \geq \lambda_{\min}(\tilde{\mathbf{K}}^{[1]}) - \left\| \mathbf{G}^{[1]}(0) - \tilde{\mathbf{K}}^{[1]} \right\|_{2 \rightarrow 2} \geq \frac{3\lambda_0}{4}. \tag{C.22}$$

This completes the proof. □

Our next proposition on the stability of the randomly initialized Gram matrix  $\mathbf{G}^{[l]}(0)$  for  $2 \leq l \leq L+1$ .

**Proposition C.3.** *Under Assumptions 4.1 and 4.2, if  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\varepsilon}\right)$ , then with high probability*

$$\lambda_{\min}(\mathbf{G}^{[l]}(0)) \geq \frac{3\lambda_0}{4}, \quad 2 \leq l \leq L+1, \tag{C.23}$$

where  $\lambda_0$  has been defined in (B.11).

*Proof.* For  $l=2$ , we shall make estimate on the norm,  $\|\mathbf{G}^{[2]}(0) - \tilde{\mathbf{K}}^{[2]}\|_\infty$ , since by definition

$$\begin{aligned} \mathbf{G}_{ij}^{[2]}(0) &= \langle \mathbf{x}_i^{[2]}(0), \mathbf{x}_j^{[2]}(0) \rangle = \mathbf{G}_{ij}^{[1]}(0) + \underbrace{\frac{c_{\text{res}}}{L} \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \sigma(\mathbf{W}^{[2]}(0)\mathbf{x}_j^{[1]}(0)) \rangle}_{\text{I}} \\ &\quad + \underbrace{\frac{c_{\text{res}}}{L} \frac{1}{\sqrt{m}} \langle \mathbf{x}_j^{[1]}(0), \sigma(\mathbf{W}^{[2]}(0)\mathbf{x}_i^{[1]}(0)) \rangle}_{\text{II}} \\ &\quad + \underbrace{\frac{c_{\text{res}}^2}{L^2} \frac{1}{m} \langle \sigma(\mathbf{W}^{[2]}(0)\mathbf{x}_i^{[1]}(0)), \sigma(\mathbf{W}^{[2]}(0)\mathbf{x}_j^{[1]}(0)) \rangle}_{\text{III}}, \\ \tilde{\mathbf{K}}_{ij}^{[1]} &= c_\sigma \mathbb{E}_{(u,v)^\top \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[0]} & \tilde{\mathbf{K}}_{ij}^{[0]} \\ \tilde{\mathbf{K}}_{ji}^{[0]} & \tilde{\mathbf{K}}_{jj}^{[0]} \end{pmatrix}\right)} [\sigma(u)\sigma(v)], \\ \tilde{\mathbf{b}}_i^{[1]} &= \sqrt{c_\sigma} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[0]})} [\sigma(u)], \\ \tilde{\mathbf{A}}_{ij}^{[2]} &= \begin{pmatrix} \tilde{\mathbf{K}}_{ii}^{[1]} & \tilde{\mathbf{K}}_{ij}^{[1]} \\ \tilde{\mathbf{K}}_{ji}^{[1]} & \tilde{\mathbf{K}}_{jj}^{[1]} \end{pmatrix}, \\ \tilde{\mathbf{K}}_{ij}^{[2]} &= \tilde{\mathbf{K}}_{ij}^{[1]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[2]})} \left[ \underbrace{\frac{c_{\text{res}}}{L} \tilde{\mathbf{b}}_i^{[1]} \sigma(v)}_{\text{I}'} + \underbrace{\frac{c_{\text{res}}}{L} \tilde{\mathbf{b}}_j^{[1]} \sigma(u)}_{\text{II}'} + \underbrace{\frac{c_{\text{res}}^2}{L^2} \sigma(u)\sigma(v)}_{\text{III}'} \right]. \end{aligned}$$

We need to tackle the difference between I and I'. In order for that, we need to write

the difference into

$$\begin{aligned} & \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \sigma(\mathbf{W}^{[2]}(0)\mathbf{x}_j^{[1]}(0)) \rangle - \tilde{\mathbf{b}}_i^{[1]} \sigma(v) \right| \\ & \leq \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \sigma(\mathbf{W}^{[2]}(0)\mathbf{x}_j^{[1]}(0)) \rangle - \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \mathbb{E} \left[ \sigma \left( \|\mathbf{x}_j^{[1]}(0)\|_2 \mathbf{Y} \right) \right] \rangle \right| \\ & \quad + \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \mathbb{E} \left[ \sigma \left( \|\mathbf{x}_j^{[1]}(0)\|_2 \mathbf{Y} \right) \right] \rangle - \tilde{\mathbf{b}}_i^{[1]} \mathbb{E} \left[ \sigma \left( \|\mathbf{x}_j^{[1]}(0)\|_2 \mathbf{Y} \right) \right] \right| \\ & \quad + \left| \tilde{\mathbf{b}}_i^{[1]} \mathbb{E} \left[ \sigma \left( \|\mathbf{x}_j^{[1]}(0)\|_2 \mathbf{Y} \right) \right] - \tilde{\mathbf{b}}_i^{[1]} \mathbb{E} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{jj}^{[1]}} \mathbf{Y} \right) \right] \right|, \end{aligned}$$

similar to the proof in Proposition C.1 with  $\mathbf{Y}$  being a standard normal Gaussian vector

$$\mathbb{P}(|\mathbf{I} - \mathbf{I}'| \geq t) \leq \exp(-cmt^2). \tag{C.24}$$

Similarly, we have

$$\mathbb{P}(|\mathbf{II} - \mathbf{II}'| \geq t) \leq \exp(-cmt^2). \tag{C.25}$$

For the difference between  $\mathbf{III}$  and  $\mathbf{III}'$ , we need to define another inner product function  $h^{[2]}(\cdot) : \mathbb{R}^{2m} \rightarrow \mathbb{R}$ , being

$$h^{[2]}(\mathbf{Z}) = \frac{1}{m} \left\langle \sigma(C_2\mathbf{X}), \sigma \left( D_2 \left( \rho\mathbf{X} + \sqrt{1-\rho^2}\mathbf{Y} \right) \right) \right\rangle,$$

with  $C_2, D_2 > 0$  being constants and  $-1 \leq \rho \leq 1$ .

Note that the form  $\mathbf{Z}^\top = (\mathbf{X}^\top, \mathbf{Y}^\top)$ , similar to  $h^{[1]}(\cdot)$  defined in the proof of Proposition C.2,  $h^{[2]}(\cdot)$  is  $\frac{C}{\sqrt{m}}$ -Lipschitz. Then we have

$$\mathbb{P} \left( \left| h^{[2]}(\mathbf{Z}) - \mathbb{E}h^{[2]}(\mathbf{Z}) \right| \geq t \right) \leq \exp(-cmt^2).$$

Hence we have

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{m} \langle \sigma(\mathbf{W}^{[2]}(0)\mathbf{x}_i^{[1]}(0)), \sigma(\mathbf{W}^{[2]}(0)\mathbf{x}_j^{[1]}(0)) \rangle - \mathbb{E}_{(u,v)^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_{ij}^{[2]})} [\sigma(u)\sigma(v)] \right| \geq t \right) \\ & \leq \exp(-cmt^2), \end{aligned} \tag{C.26}$$

with

$$\mathbf{A}_{ij}^{[2]} = \begin{pmatrix} \langle \mathbf{x}_i^{[1]}(0), \mathbf{x}_i^{[1]}(0) \rangle & \langle \mathbf{x}_i^{[1]}(0), \mathbf{x}_j^{[1]}(0) \rangle \\ \langle \mathbf{x}_j^{[1]}(0), \mathbf{x}_i^{[1]}(0) \rangle & \langle \mathbf{x}_j^{[1]}(0), \mathbf{x}_j^{[1]}(0) \rangle \end{pmatrix},$$

combined with Lemma C.3 and Proposition C.1

$$\mathbb{P} \left( \left| \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_{ij}^{[2]})} [\sigma(u)\sigma(v)] - \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[2]})} [\sigma(u)\sigma(v)] \right| \geq t \right) \leq \exp(-cmt^2). \quad (\text{C.27})$$

Combining (C.26) and (C.27)

$$\mathbb{P} (|\text{III} - \text{III}'| \geq t) \leq \exp(-cmt^2), \quad (\text{C.28})$$

we have that

$$\begin{aligned} & \mathbb{P} \left( \left| \mathbf{G}_{ij}^{[2]}(0) - \tilde{\mathbf{K}}_{ij}^{[2]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right)^2 \right) \\ & \leq \mathbb{P} \left( \left| \mathbf{G}_{ij}^{[1]}(0) - \tilde{\mathbf{K}}_{ij}^{[1]} \right| \geq t \right) + \mathbb{P} (|\text{I} - \text{I}'| \geq t) + \mathbb{P} (|\text{II} - \text{II}'| \geq t) + \mathbb{P} (|\text{III} - \text{III}'| \geq t) \\ & \leq \exp(-cmt^2). \end{aligned} \quad (\text{C.29})$$

Hence inductively, for  $2 \leq l \leq L$

$$\mathbb{P} \left( \left| \mathbf{G}_{ij}^{[l]}(0) - \tilde{\mathbf{K}}_{ij}^{[l]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right)^{2l-2} \right) \leq \exp(-cmt^2). \quad (\text{C.30})$$

Moreover,

$$\mathbb{P} \left( \left| \mathbf{G}_{ij}^{[L+1]}(0) - \mathbf{K}_{ij}^{[L+1]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right)^{2L} \right) \leq \exp(-cmt^2), \quad (\text{C.31})$$

note that we have

$$\left\| \mathbf{G}^{[L+1]}(0) - \mathbf{K}^{[L+1]} \right\|_{2 \rightarrow 2} \leq \left\| \mathbf{G}^{[L+1]}(0) - \mathbf{K}^{[L+1]} \right\|_{\text{F}} \leq n \left\| \mathbf{G}^{[L+1]}(0) - \mathbf{K}^{[L+1]} \right\|_{\infty},$$

based on Proposition B.1,  $\lambda_{\min}(\mathbf{K}^{[L+1]}) > \lambda_0$ , then if we choose  $t = \frac{\lambda_0}{4n \exp(2c_{\text{res}})}$ , for  $2 \leq l \leq L$ , with probability  $1 - \exp(-cm\lambda_0^2/n^2)$ ,

$$\left\| \mathbf{G}^{[l]}(0) - \tilde{\mathbf{K}}^{[l]} \right\|_{2 \rightarrow 2} \leq \frac{\lambda_0}{4}, \quad (\text{C.32})$$

hence if  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\epsilon}\right)$ , we have with probability  $1 - \exp(-m^\epsilon)$

$$\lambda_{\min}(\mathbf{G}^{[l]}(0)) \geq \lambda_{\min}(\tilde{\mathbf{K}}^{[l]}) - \left\| \mathbf{G}^{[l]}(0) - \tilde{\mathbf{K}}^{[l]} \right\|_{2 \rightarrow 2} > \frac{3\lambda_0}{4}. \quad (\text{C.33})$$

In particular, we have that with probability  $1 - \exp(-cm\lambda_0^2/n^2)$ ,

$$\left\| \mathbf{G}^{[L+1]}(0) - \mathbf{K}^{[L+1]} \right\|_{2 \rightarrow 2} \leq \frac{\lambda_0}{4}, \quad (\text{C.34})$$

hence if  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\epsilon}\right)$ , we have with probability  $1 - \exp(-m^\epsilon)$

$$\lambda_{\min}(\mathbf{G}^{[L+1]}(0)) \geq \lambda_{\min}(\mathbf{K}^{[L+1]}) - \left\| \mathbf{G}^{[L+1]}(0) - \mathbf{K}^{[L+1]} \right\|_{2 \rightarrow 2} > \frac{3\lambda_0}{4}. \quad (\text{C.35})$$

This completes the proof.  $\square$

### D Proof of Theorem 4.2 and Corollary 4.1

We shall begin with the detailed proof of Theorem 4.2.

*Proof of Theorem 4.2.* We are only going to use  $\mathcal{G}_t^{[L+1]}(\cdot)$  instead of the whole NTK  $\mathcal{K}_t^{(2)}(\cdot)$ , thanks to the simple structure of  $\mathcal{G}_t^{[L+1]}(\cdot)$ , we are able to bring about a more concrete proof.

Since there exists a  $\frac{1}{L^2}$  scaling in some kernels, we use  $C(r, L)$  to denote the ‘effective terms’ in each kernel and we are going to show that (4.11) holds. Firstly, we need to denote  $\mathcal{G}_t^{[L+1]}(\cdot)$  by  $\mathcal{G}_t^{[2]}(\cdot)$ , i.e.,

$$\mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) := \mathcal{G}_t^{[L+1]}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) = \langle \mathbf{x}_{\alpha_1}^{[L]}, \mathbf{x}_{\alpha_2}^{[L]} \rangle,$$

then it’s natural for us to get that  $C(2, L) = \mathcal{O}(1)$ , since there is only one term.

Secondly, by the replacement rule, all the possible terms generated from  $\mathcal{G}_t^{(2)}(\cdot)$  are

$$\begin{aligned} \mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) &= \langle \mathbf{x}_{\alpha_1}^{[L]}, \mathbf{x}_{\alpha_2}^{[L]} \rangle \rightarrow \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta), \\ \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta) &= \frac{c_\sigma}{m} \underbrace{\langle \text{diag} \left( \mathbf{E}_{t, \alpha_1}^{[2:L]} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_{\alpha_1}) \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1}, \mathbf{x}_{\alpha_2}^{[L]} \rangle}_{\text{I}} \langle \mathbf{x}_{\alpha_1}, \mathbf{x}_\beta \rangle \\ &+ \sum_{k=2}^L \frac{c_{\text{res}}^2}{L^2 m} \underbrace{\langle \text{diag} \left( \mathbf{E}_{t, \alpha_1}^{[(k+1):L]} \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_{\alpha_1}) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1}, \mathbf{x}_{\alpha_2}^{[L]} \rangle}_{\text{II}} \langle \mathbf{x}_{\alpha_1}^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle \\ &+ \frac{c_\sigma}{m} \langle \text{diag} \left( \mathbf{E}_{t, \alpha_2}^{[2:L]} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_{\alpha_2}) \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1}, \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta \rangle \\ &+ \sum_{k=2}^L \frac{c_{\text{res}}^2}{L^2 m} \langle \text{diag} \left( \mathbf{E}_{t, \alpha_2}^{[(k+1):L]} \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_{\alpha_2}) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1}, \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle. \end{aligned}$$

Thanks to the  $\frac{1}{L^2}$  scaling, we obtain that

$$C(3, L) = \mathcal{O} \left( 2 \left( 1 + \frac{L-1}{L^2} \right) \right) = \mathcal{O} \left( 1 + \frac{1}{L} \right).$$

Finally for  $\mathcal{G}_t^{(4)}(\cdot)$ , by symmetry, we are only going to analyze terms I and II. Since there are at most  $(2L+2)$  symbols in term I to be replaced, and by the replacement rules, each replacement will bring about up to  $(L+1)$  many terms. For term II, for each summand, there are also at most  $(2L+2)$  symbols to be replaced. Since there are  $L-1$  summands in II, and each replacement will bring about up to  $(L+1)$  many terms, then we have

$$C(4, L) = \mathcal{O} \left( 2 \left( (2L+2)(L+1) + \frac{1}{L^2} (L-1)(2L+2)(L+1) \right) \right) = \mathcal{O}(L^2).$$

Using (4.9) in Theorem 4.1, it holds that for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$

$$\left\| \mathcal{G}_t^{(4)}(\cdot) \right\|_\infty \leq C(4, L) \frac{(\ln m)^C}{m},$$

based on (4.7)

$$\begin{aligned} \left| \partial_t \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| &\leq \sup_{1 \leq \beta \leq n} \left| \mathcal{G}_t^{(4)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}, \mathbf{x}_\beta) \right| \sqrt{\frac{\sum_{\beta=1}^n |f_\beta(t) - y_\beta|^2}{n}} \\ &\leq \left\| \mathcal{G}_t^{(4)}(\cdot) \right\|_\infty \sqrt{R_S(\boldsymbol{\theta}_0)} \\ &\leq C(4, L) \frac{(\ln m)^C}{m}, \end{aligned}$$

then for any  $1 \leq \alpha_1, \alpha_2, \alpha_3 \leq n$ , with time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$

$$\begin{aligned} \left| \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| &\leq \left| \mathcal{G}_0^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| + t \left| \partial_t \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| \\ &\leq \left\| \mathcal{G}_0^{(3)}(\cdot) \right\|_\infty + t C(4, L) \frac{(\ln m)^C}{m}. \end{aligned}$$

Finally, we need to make estimate on  $\left\| \mathcal{G}_0^{(3)}(\cdot) \right\|_\infty$ . We shall take advantage of the  $\text{diag}(\cdot)\mathbf{1}$  structure and rewrite  $\mathcal{G}_t^{(3)}(\cdot)$  into

$$\begin{aligned} &\mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta) \\ &= \frac{c_\sigma}{m} \left\langle \mathbf{E}_{t, \alpha_1}^{[2:L]} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_{\alpha_1}) \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[2:L]} \right)^\top \mathbf{a}_t, \mathbf{x}_{\alpha_2}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_1}, \mathbf{x}_\beta \rangle \\ &\quad + \sum_{k=2}^L \frac{c_{\text{res}}^2}{L^2 m} \left\langle \mathbf{E}_{t, \alpha_1}^{[(k+1):L]} \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_{\alpha_1}) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t, \mathbf{x}_{\alpha_2}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_1}^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle \\ &\quad + \frac{c_\sigma}{m} \left\langle \mathbf{E}_{t, \alpha_2}^{[2:L]} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_{\alpha_2}) \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[2:L]} \right)^\top \mathbf{a}_t, \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta \rangle \\ &\quad + \sum_{k=2}^L \frac{c_{\text{res}}^2}{L^2 m} \left\langle \mathbf{E}_{t, \alpha_2}^{[(k+1):L]} \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_{\alpha_2}) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t, \beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t, \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \langle \mathbf{x}_{\alpha_2}^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle, \end{aligned}$$

then at time  $t=0$ , without loss of generality, each term in  $\mathcal{G}_0^{(3)}(\cdot)$  is of the form

$$\frac{c}{m} \left\langle \mathbf{B} \mathbf{a}_0, \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \left\langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \right\rangle, \quad 0 \leq l \leq L-1, \tag{D.1}$$

where  $\mathbf{B}$  is some specific matrix that changes from term to term, then we can rewrite the inner product into:

$$\frac{c}{m} \left\langle \mathbf{a}_0, \mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \left\langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \right\rangle, \tag{D.2}$$

we shall focus on the term

$$\langle \mathbf{a}_0, \mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle, \tag{D.3}$$

note that each entry of  $\mathbf{a}_0$  is i.i.d  $\mathcal{N}(0,1)$ , also based on Propositions A.3 and A.4, with high probability with respect to random initialization, for time  $0 \leq t \leq (\ln m)^{C'}$

$$\|\mathbf{B}^\top\|_{2 \rightarrow 2, \mathbf{x}_{\alpha_1}^{[L]}, \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]}} \leq c,$$

then after taking conditional expectation except for the random variable  $\mathbf{a}_0$

$$\langle \mathbf{a}_0, \mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle \sim \mathcal{N}(0, c), \tag{D.4}$$

apply Lemma A.3 directly, with high probability

$$\frac{c}{m} \langle \mathbf{a}_0, \mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle \leq c \frac{(\ln m)^C}{m}, \tag{D.5}$$

consequently

$$\|\mathcal{G}_0^{(3)}(\cdot)\|_\infty \leq C(3, L) \frac{(\ln m)^C}{m}, \tag{D.6}$$

then for any  $1 \leq \alpha_1, \alpha_2, \alpha_3 \leq n$ , with time  $0 \leq t \leq \sqrt{m} / (\ln m)^{C'}$

$$\begin{aligned} \left| \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| &\leq \|\mathcal{G}_0^{(3)}(\cdot)\|_\infty + tC(4, L) \frac{(\ln m)^C}{m} \\ &\leq C(3, L) \frac{(\ln m)^C}{m} + tC(4, L) \frac{(\ln m)^C}{m}. \end{aligned}$$

Similarly, based on (4.7), for time  $0 \leq t \leq \sqrt{m} / (\ln m)^{C'}$

$$\left| \partial_t \mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) \right| \leq \sup_{1 \leq \beta \leq n} \left| \mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta) \right| \sqrt{\frac{\sum_{\beta=1}^n |f_\beta(t) - y_\beta|^2}{n}},$$

set  $\mathbf{x}_\beta = \mathbf{x}_{\alpha_3}$

$$\begin{aligned} \left| \partial_t \mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) \right| &\leq \left( C(3, L) \frac{(\ln m)^C}{m} + tC(4, L) \frac{(\ln m)^C}{m} \right) \sqrt{\frac{\sum_{\beta=1}^n |f_\beta(t) - y_\beta|^2}{n}} \\ &\leq (C(3, L) + tC(4, L)) \frac{(\ln m)^C}{m}, \end{aligned} \tag{D.7}$$

and (D.7) finishes the proof of Theorem 4.2. □

*Proof of Corollary 4.1.* Firstly, based on Proposition C.3, if  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\epsilon}\right)$ , we have with high probability with respect to random initialization,

$$\lambda_{\min} \left[ \mathcal{K}_0^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} > \lambda_{\min} \left( \mathbf{G}^{[L+1]}(0) \right) > \frac{3\lambda_0}{4},$$

set  $\lambda = \frac{3\lambda_0}{4}$ , which finishes the proof of (4.12).

We shall move on to the change of the least eigenvalue of the NTK. Recall (D.7) in the proof of Theorem 4.2, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C'}$ ,

$$\left| \partial_t \mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) \right| \leq (C(3, L) + tC(4, L)) \frac{(\ln m)^C}{m},$$

consequently

$$\left| \mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) - \mathcal{G}_0^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) \right| \leq t(C(3, L) + tC(4, L)) \frac{(\ln m)^C}{m}.$$

The above inequality can be used to derive a bound of the change of the least eigenvalue of the  $\mathcal{G}_t^{(2)}(\cdot)$

$$\begin{aligned} \left\| \mathcal{G}_t^{(2)} - \mathcal{G}_0^{(2)} \right\|_{2 \rightarrow 2} &\leq \left\| \mathcal{G}_t^{(2)} - \mathcal{G}_0^{(2)} \right\|_{\text{F}} \leq n \left\| \mathcal{G}_t^{(2)} - \mathcal{G}_0^{(2)} \right\|_{\infty} \\ &\leq nt(C(3, L) + tC(4, L)) \frac{(\ln m)^C}{m}, \end{aligned}$$

we set  $t^*$  satisfying

$$nt^*(C(3, L) + t^*C(4, L)) \frac{(\ln m)^C}{m} = \frac{\lambda}{2},$$

rewrite the equation above, we have

$$C(4, L)(t^*)^2 + C(3, L)t^* = \frac{\lambda m}{2(\ln m)^C n}, \tag{D.8}$$

solve (D.8), we obtain that

$$t^* = \frac{-C(3, L) + \sqrt{(C(3, L))^2 + 2C(4, L) \frac{\lambda m}{(\ln m)^C n}}}{2C(4, L)}, \tag{D.9}$$

since we are in the regime of over-parametrization, for  $m$  large enough, the following holds

$$t^* \geq \frac{1}{2} \sqrt{\frac{\frac{\lambda m}{(\ln m)^C n}}{C(4, L)}} = \frac{1}{2} \sqrt{\frac{\lambda m}{C(4, L)(\ln m)^C n}}. \tag{D.10}$$



Moreover

$$\begin{aligned} \lambda_{\min} \left[ \mathcal{K}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} &\geq \lambda_{\min} \left[ \mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} \\ &\geq \lambda_{\min} \left[ \mathcal{G}_0^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} - \left\| \mathcal{G}_t^{(2)} - \mathcal{G}_0^{(2)} \right\|_{2 \rightarrow 2}, \end{aligned}$$

then let

$$\bar{t} := \inf \left\{ t : \lambda_{\min} \left[ \mathcal{K}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} \geq \lambda/2 \right\},$$

naturally

$$t^* \leq \bar{t}, \tag{D.11}$$

using (4.6), we have for any  $0 \leq t \leq \bar{t}$ ,

$$\begin{aligned} \partial_t \sum_{\alpha=1}^n \|f_\alpha(t) - y_\alpha\|_2^2 &\leq \sum_{\alpha, \beta=1}^n -\frac{2}{n} \mathcal{K}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) (f_\alpha(t) - y_\alpha)(f_\beta(t) - y_\beta) \\ &\leq -\frac{\lambda}{n} \sum_{\alpha=1}^n \|f_\alpha(t) - y_\alpha\|_2^2, \end{aligned} \tag{D.12}$$

then

$$\sum_{\alpha=1}^n \|f_\alpha(t) - y_\alpha\|_2^2 \leq \exp\left(-\frac{\lambda t}{n}\right) \sum_{\alpha=1}^n \|f_\alpha(0) - y_\alpha\|_2^2, \tag{D.13}$$

we can rewrite (D.13) into

$$R_S(\boldsymbol{\theta}_t) \leq \exp\left(-\frac{\lambda t}{n}\right) R_S(\boldsymbol{\theta}_0), \tag{D.14}$$

set  $R_S(\boldsymbol{\theta}_t) = \varepsilon$ , it takes time  $t \leq \frac{n}{\lambda} \ln\left(\frac{C'}{\varepsilon}\right)$  for loss  $R_S(\boldsymbol{\theta}_t)$  to reach accuracy  $\varepsilon$ , hence if the following holds

$$t \leq \frac{n}{\lambda} \ln\left(\frac{C'}{\varepsilon}\right) \leq t^* \leq \bar{t}, \tag{D.15}$$

then the width  $m$  is required to yield the lower bound for  $t^*$  derived in (D.10),

$$\frac{n}{\lambda} \ln\left(\frac{C'}{\varepsilon}\right) \leq \frac{1}{2} \sqrt{\frac{\lambda m}{C(4, L)(\ln m)^C n}}. \tag{D.16}$$

Then we have

$$m \geq C(4, L) \left(\frac{n}{\lambda}\right)^3 (\ln m)^C \ln\left(\frac{C'}{\varepsilon}\right)^2,$$

since  $C(4, L) = \mathcal{O}(L^2)$ , we conclude that the required width  $m$  should be

$$m = \Omega\left(\left(\frac{n}{\lambda}\right)^3 L^2 (\ln m)^C \ln\left(\frac{C'}{\varepsilon}\right)\right), \tag{D.17}$$

where  $\varepsilon$  is the desired training accuracy. □

## References

- [1] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 242–252. PMLR, 09–15 Jun 2019.
- [2] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*, 2019.
- [3] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 322–332. PMLR, 09–15 Jun 2019.
- [4] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On Exact Computation with an Infinitely Wide Neural Net. In *Advances in Neural Information Processing Systems*, pages 8139–8148, 2019.
- [5] Peter L Bartlett, David P Helmbold, and Philip M Long. Gradient Descent with Identity Initialization Efficiently Learns Positive-definite Linear Transformations by Deep Residual Networks. *Neural Comput.*, 31(3):477–502, 2019.
- [6] Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian inputs. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 605–614. JMLR, 2017.
- [7] Lenaic Chizat and Francis Bach. On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport. In *Advances in Neural Information Processing Systems*, pages 3036–3046, 2018.
- [8] Ronan Collobert and Jason Weston. A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning. In *Proceedings of the 25th International Conference on Machine Learning*, pages 160–167, 2008.
- [9] George E Dahl, Dong Yu, Li Deng, and Alex Acero. Context-dependent Pre-trained Deep Neural Networks for Large-vocabulary Speech Recognition. *IEEE Audio, Speech, Language Process.*, 20(1):30–42, 2011.
- [10] Amit Daniely, Roy Frostig, and Yoram Singer. Toward Deeper Understanding of Neural Networks: The Power of Initialization and a Dual View on Expressivity. In *Advances In Neural Information Processing Systems*, pages 2253–2261, 2016.
- [11] Simon Du and Jason Lee. On the power of over-parametrization in neural networks with quadratic activation. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1329–1338. PMLR, 10–15 Jul 2018.
- [12] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1675–1685. PMLR, 09–15 Jun 2019.
- [13] Simon Du, Jason Lee, Yuandong Tian, Aarti Singh, and Barnabas Poczos. Gradient descent learns one-hidden-layer CNN: Don’t be afraid of spurious local minima. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1339–1348. PMLR, 10–15 Jul

- 2018.
- [14] Simon S Du, Chi Jin, Jason D Lee, Michael I Jordan, Aarti Singh, and Barnabas Póczos. Gradient Descent can take Exponential Time to Escape Saddle Points. In *Advances in Neural Information Processing Systems*, pages 1067–1077, 2017.
  - [15] Simon S. Du, Xiyu Zhai, Barnabas Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019.
  - [16] Weinan E, Chao Ma, and Lei Wu. A Comparative Analysis of the Optimization and Generalization Property of Two-layer Neural Network and Random Feature Models under Gradient Descent Dynamics. *Science China Mathematics*, pages 1–24, 2020.
  - [17] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from Saddle Points Online Stochastic Gradient for Tensor Decomposition. In *Conference on Learning Theory*, pages 797–842, 2015.
  - [18] Rong Ge, Jason D. Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. In *International Conference on Learning Representations*, 2018.
  - [19] Xavier Glorot and Yoshua Bengio. Understanding the Difficulty of Training Deep Feedforward Neural Networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 249–256, 2010.
  - [20] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
  - [21] Moritz Hardt and Tengyu Ma. Identity matters in deep learning. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24–26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.
  - [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
  - [23] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely Connected Convolutional Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4700–4708, 2017.
  - [24] Jiaoyang Huang and Horng-Tzer Yau. Dynamics of deep neural networks and neural tangent hierarchy. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4542–4551. PMLR, 13–18 Jul 2020.
  - [25] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. In *Advances in Neural Information Processing Systems*, pages 8571–8580, 2018.
  - [26] Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M Kakade, and Michael I Jordan. How to Escape Saddle Points Efficiently. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1724–1732. JMLR. org, 2017.
  - [27] Kenji Kawaguchi. Deep learning without Poor Local Minima. In *Advances in Neural Information Processing Systems*, pages 586–594, 2016.
  - [28] Kenji Kawaguchi and Jiaoyang Huang. Gradient Descent Finds Global Minima for Generalizable Deep Neural Networks of Practical Sizes. In *2019 57th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 92–99. IEEE, 2019.
  - [29] Beatrice Laurent and Pascal Massart. Adaptive Estimation of a Quadratic Functional by Model Selection. *Annals of Statistics*, pages 1302–1338, 2000.
  - [30] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear mod-

- els under gradient descent. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [31] Jason D Lee, Max Simchowitz, Michael I Jordan, and Benjamin Recht. Gradient Descent only Converges to Minimizers. In *Conference on Learning Theory*, pages 1246–1257, 2016.
- [32] Yuanzhi Li and Yingyu Liang. Learning Overparameterized Neural Networks via Stochastic Gradient Descent on Structured Data. In *Advances in Neural Information Processing Systems*, pages 8157–8166, 2018.
- [33] Yuanzhi Li and Yang Yuan. Convergence Analysis of Two-layer Neural Networks with Relu Activation. In *Advances in Neural Information Processing Systems*, pages 597–607, 2017.
- [34] Yuqing Li, Tao Luo, and Chao Ma. Nonlinear weighted directed acyclic graph and a priori estimates for neural networks. *arXiv preprint arXiv:2103.16355*, 2021.
- [35] Chao Ma, Qingcan Wang, and Weinan E. A Priori Estimates of the Population Risk for Residual Networks. *arXiv preprint arXiv:1903.02154*, 2019.
- [36] Chao Ma, Qingcan Wang, Lei Wu, et al. Analysis of the Gradient Descent Algorithm for a Deep Neural Network Model with Skip-connections. *arXiv preprint arXiv:1904.05263*, 2019.
- [37] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A Mean Field View of the Landscape of Two-layer Neural Networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- [38] Quynh Nguyen and Matthias Hein. The Loss Surface of Deep and Wide Neural Networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2603–2612. JMLR, 2017.
- [39] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet Classification Using Binary Convolutional Neural Networks. In *European Conference on Computer Vision*, pages 525–542. Springer, 2016.
- [40] Grant Rotskoff and Eric Vanden-Eijnden. Parameters as Interacting Particles: Long Time Convergence and Asymptotic Error Scaling of Neural Networks. In *Advances in Neural Information Processing Systems*, pages 7146–7155, 2018.
- [41] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the Game of Go with Deep Neural Networks and Tree Search. *Nature*, 529(7587):484, 2016.
- [42] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the Game of Go without Human Knowledge. *Nature*, 550(7676):354–359, 2017.
- [43] Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: A law of large numbers. *SIAM Journal on Applied Mathematics*, 80(2):725–752, 2020.
- [44] Zhao Song and Xin Yang. Quadratic Suffices for Over-parametrization via Matrix Chernoff Bound. *arXiv preprint arXiv:1906.03593*, 2019.
- [45] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alex Alemi. Inception-v4, Inception-resnet and the Impact of Residual Connections on Learning (2016). *arXiv preprint arXiv:1602.07261*, 2016.
- [46] Roman Vershynin. Introduction to the Non-asymptotic Analysis of Random Matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [47] Greg Yang. Scaling Limits of Wide Neural Networks with Weight Sharing: Gaussian Process Behavior, Gradient Independence, and Neural Tangent Kernel Derivation. *arXiv preprint arXiv:1902.04760*, 2019.

- [48] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In Edwin R. Hancock Richard C. Wilson and William A. P. Smith, editors, *Proceedings of the British Machine Vision Conference (BMVC)*, pages 87.1–87.12. BMVA Press, September 2016.
- [49] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding Deep Learning Requires Rethinking Generalization, 2018.
- [50] Guodong Zhang, James Martens, and Roger B Grosse. Fast Convergence of Natural Gradient Descent for Over-parameterized Neural Networks. In *Advances in Neural Information Processing Systems*, pages 8080–8091, 2019.
- [51] Yi Zhou and Yingbin Liang. Critical Points of Neural Networks: Analytical Forms and Landscape Properties. *arXiv preprint arXiv:1710.11205*, 2017.
- [52] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep ReLU networks. *Machine Learning*, 109(3):467–492, 2020.