

# TOWARDS AN UNDERSTANDING OF RESIDUAL NETWORKS USING NEURAL TANGENT HIERARCHY

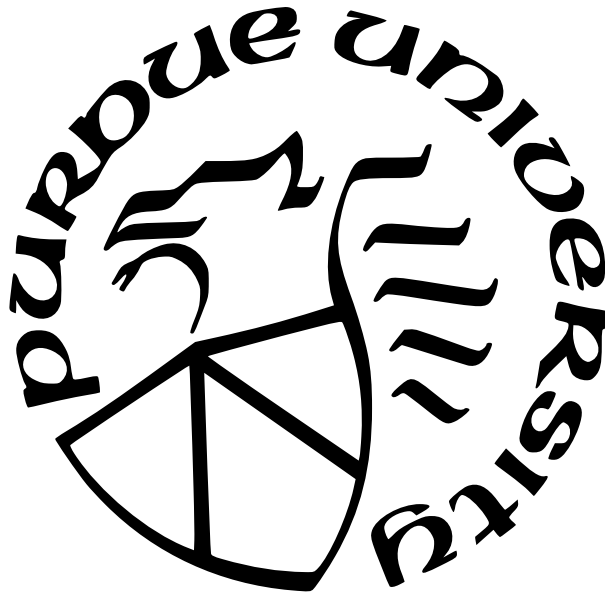
by  
Yuqing Li

A Dissertation

*Submitted to the Faculty of Purdue University*

*In Partial Fulfillment of the Requirements for the degree of*

Doctor of Philosophy



Department of Mathematics

West Lafayette, Indiana

May 2021

**THE PURDUE UNIVERSITY GRADUATE SCHOOL  
STATEMENT OF COMMITTEE APPROVAL**

**Dr. Nung Kwan Yip, Chair**

Department of Mathematics

**Dr. Vinayak Rao**

Department of Statistics Department of Statistics and Computer Science (by courtesy)

**Dr. Samy Tindel**

Department of Mathematics

**Dr. Guang Lin**

Department of Mathematics, Statistics and School of Mechanical Engineering

**Approved by:**

Dr. Plamen Stefanov

To my parents.

To my beloved meemaw and deceased papaw.

To my cutest stuffed dog, who proudly and shamelessly named himself Dr. Doge.

To my undergraduate friend, William Chengfei Wu, for his great jokes on Dr. Doge.

To my teddy bears, especially the largest one whom I found abandoned by someone else  
near a trash can.

## ACKNOWLEDGMENTS

First of all, most importantly, I shall give my sincere appreciation to my supervisor, Prof. Aaron Nung Kwan Yip, who has been giving me tremendous support on my study, research and daily life.

Next, I gratefully acknowledge Prof. Vinayak Rao, Prof. Samy Tindel, and Prof. Guang Lin for their kindness and readiness to serve on my thesis committee. I would also like to thank the faculty and staff of the Department of Mathematics for everything during my study at Purdue.

I have been fortunate to spend more than three years working with Prof. Vinayak Rao, and would like to thank him for his help and patience. Prof. Rao is an inspiration, both in how much he knows and how much he still wants to learn. I am always inspired by him via his thorough understanding on statistics and machine learning. I benefited immensely from his ideas, his weekly group discussions and his feedback, and from the general confidence that I was in very good hands.

I would like to give special thanks to my roommate, semi-academic brother, semi-academic advisor and friend, Prof. Tao Luo. When I was a second-year graduate student, he changed my perspective on mathematics and trained me to have a solid foundation in mathematics. I have been receiving numerous advice from him hereafter.

Finally and most importantly, I would like to thank my beloved mommy Limin Zhang, and my father Jun Li for their endless love and support. This family means everything to me, and I dedicate this thesis to them.

# TABLE OF CONTENTS

LIST OF TABLES . . . . .	7
LIST OF FIGURES . . . . .	8
ABSTRACT . . . . .	9
1 INTRODUCTION . . . . .	10
1.1 Introduction . . . . .	10
1.2 Fully-connected Neural Networks . . . . .	11
1.3 ResNet . . . . .	13
1.4 Main Results . . . . .	15
1.5 Assumptions and Notations . . . . .	16
2 GRADIENT DESCENT FINDS GLOBAL MINIMA . . . . .	19
2.1 Introduction . . . . .	19
2.2 Gradient Descent . . . . .	19
2.3 NTK . . . . .	21
2.4 Gram Matrices . . . . .	23
2.4.1 Gram Matrices for Fully-connected Networks . . . . .	25
2.4.2 Gram Matrices for ResNet . . . . .	26
2.5 Main Results of Du et al. . . . .	27
2.6 Main Results of Huang and Yau . . . . .	30
3 RESNET USING NTH . . . . .	33
3.1 Introduction . . . . .	33
3.2 Main Results . . . . .	33
3.3 Key Technique Number One: Kernel Structure . . . . .	37
3.3.1 Replacement Rules . . . . .	37
3.3.2 Hierarchical Sets of Kernel Expressions . . . . .	42
3.4 Key Technique Number Two: Apriori Estimates . . . . .	47

3.4.1	Apriori $L^2$ Bounds for Expressions in $\mathbb{A}_0$ . . . . .	48
3.4.2	Apriori $L^\infty$ Bounds for Expressions in $\mathbb{A}_0$ . . . . .	55
3.4.3	Apriori $L^2$ and $L^\infty$ Bounds for Expressions in $\mathbb{A}_r$ , $r \geq 1$ . . . . .	63
3.5	Proof of Theorem 3.2.1 . . . . .	71
3.6	Proof of Theorem 3.2.2 . . . . .	71
3.7	Proof of Theorem 3.2.3 . . . . .	75
4	SUMMARY AND FUTURE WORK . . . . .	78
	REFERENCES . . . . .	80
A	LEAST EIGENVALUE OF GRAM MATRICES . . . . .	86
A.1	Introduction . . . . .	86
A.2	Full Rankness for $(L + 1)$ -th Gram matrix . . . . .	87
A.3	Full Rankness for the $L$ -th Gram matrix . . . . .	93
B	RANDOM INITIALIZATION OF NTK . . . . .	95
B.1	Several Lemmas on Gaussian Concentration and Other Aspects . . . . .	95
B.2	Analysis of Random Propagation . . . . .	97
B.3	Analysis on Random Initialization . . . . .	107

## LIST OF TABLES

1.1	Notation Table . . . . .	17
-----	--------------------------	----

## LIST OF FIGURES

1.1	A fully connected layer in a deep network . . . . .	11
1.2	Multilayer deep fully-connected network . . . . .	12
1.3	Residual learning: a building block . . . . .	14
4.1	Diagram for the Proof of Main Theorems . . . . .	78



# ABSTRACT

Deep learning has become an important toolkit for data science and artificial intelligence. In contrast to its practical success across a wide range of fields, theoretical understanding of the principles behind the success of deep learning has been an issue of controversy. Optimization, as an important component of theoretical machine learning, has attracted much attention. The optimization problems induced from deep learning is often non-convex and non-smooth, which is challenging to locate the global optima. However, in practice, global convergence of first-order methods like gradient descent can be guaranteed for deep neural networks. In particular, gradient descent yields zero training loss in polynomial time for deep neural networks despite its non-convex nature. Besides that, another mysterious phenomenon is the compelling performance of Deep Residual Network (ResNet). Not only does training ResNet require weaker conditions [1], the employment of residual connections by ResNet even enables first-order methods to train the neural networks with an order of magnitude more layers [2]. Advantages arising from the usage of residual connections remain to be discovered.

In this thesis, we demystify these two phenomena accordingly. Firstly, we contribute to further understanding of gradient descent. The core of our analysis is the neural tangent hierarchy (NTH) [3] that captures the gradient descent dynamics of deep neural networks. A recent work [4] introduced the Neural Tangent Kernel (NTK) and proved that the limiting NTK describes the asymptotic behavior of neural networks trained by gradient descent in the infinite width limit. The NTH outperforms the NTK in two ways: (i) It can directly study the time variation of NTK for neural networks. (ii) It improves the result to non-asymptotic settings. Moreover, by applying NTH to 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, obtaining a state-of-the-art result. Secondly, we extend our scope of analysis to structural properties of deep neural networks. By making fair and consistent comparisons between fully-connected network and ResNet, we suggest strongly that the particular skip-connection architecture possessed by ResNet is the main reason for its triumph over fully-connected network.

# 1. INTRODUCTION

## 1.1 Introduction

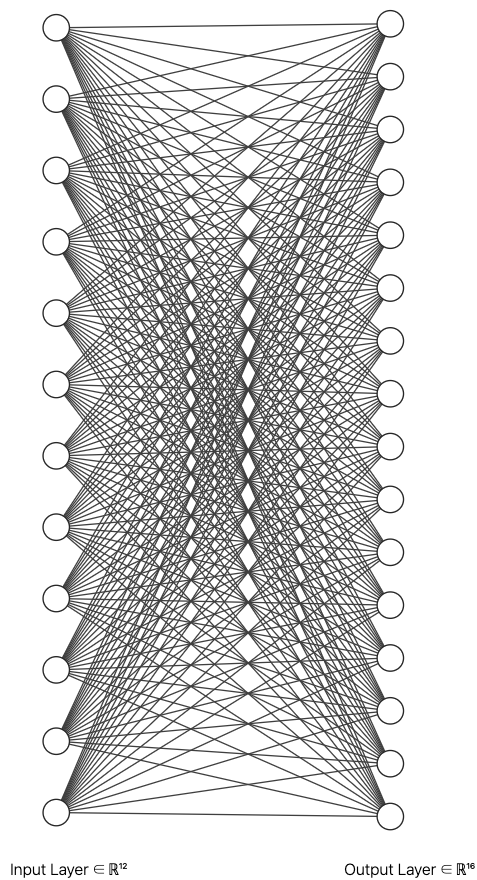
Deep neural networks have achieved transcendent performance in a wide range of tasks such as speech recognition, computer vision, and natural language processing. There are various methods to train neural networks, such as first-order gradient based methods, which have been proven to achieve satisfactory results [5]. Experiments in [6] established 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 stochastic gradient descent, the network is still able to fit them well. There are numerous works trying to demystify such phenomenon theoretically. Du et al. [7] proved that gradient descent can obtain zero training loss for two-layer networks, and Zou et al. [8] analyzed the convergence of stochastic gradient descent on networks assembled with Rectified Linear Unit (ReLU) activation function. All these neural networks are heavily overparameterized: the number of learnable parameters is much larger than the number of the training samples. It is widely accepted by the machine learning community that overparameterization enables the neural network to fit all training data, and it brings no harm to the power of its generalization, i.e., the ability to predict well on unseen data [9]. In particular, the deep neural networks that evaluated positions and selected moves for the well-known program AlphaGo are highly overparameterized [10], [11].

Another advance is the outstanding performance of Deep Residual Network (ResNet) proposed by He et al. [2]. 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. 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 [12], [13]. Theoretically, Hardt and Ma [14] proved that for any residual linear networks with arbitrary depth, there are no spurious local optima. Du et al. [1] showed that in the scope of the convergence of gradient descent via overparameterization for different networks, training ResNet requires weaker conditions compared with fully-connected networks. In this thesis, we make fair comparisons between ResNet and fully-connected networks by enforcing exactly same assumptions on input samples and acti-

vation functions, such settings would thus potentially enable us to explore more structural benefits of deep neural networks.

## 1.2 Fully-connected Neural Networks

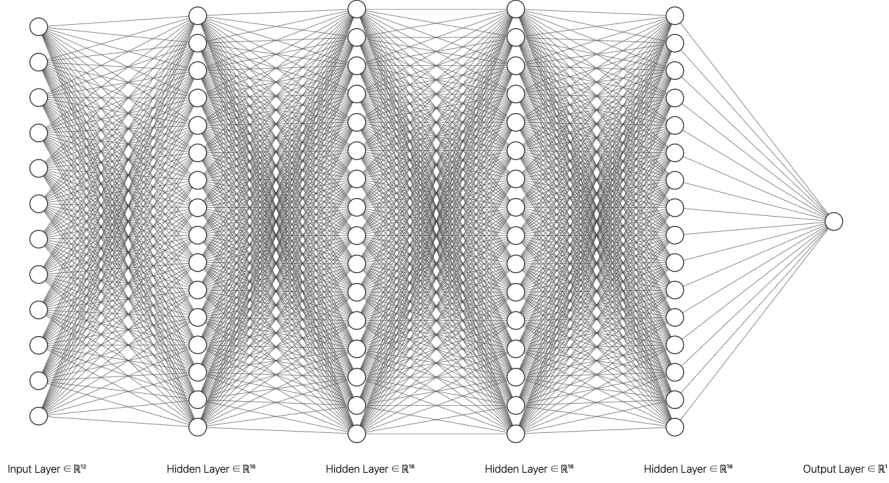
We present two types of neural network architectures respectively in Section 1.2 and Section 1.3, one is the fully-connected network, the other is ResNet. A fully-connected neural network consists of a series of fully connected layers that connect every neuron in one layer to every neuron in the other layer. Pictorially, a fully connected layer is represented in Figure 1.1.



**Figure 1.1.** A fully connected layer in a deep network

According to the classical universal approximation theorem [15], a two-layer neural network with sigmoid-like activation functions is sufficient to represent any continuous function

on the unit cube. However, its single hidden layer might be massive and the network is prone to overfitting the data. Overfitting is a serious problem in machine learning, it happens when a model learns the detail and noise in the training data to the extent that it negatively impacts the performance of the model on new data [16]. Therefore, the machine learning community has a common trend that network architecture needs to go deeper. Naturally, given fully connected layers, it is directly possible to form a network by stacking more of them, as depicted in Figure 1.2.



**Figure 1.2.** Multilayer deep fully-connected network

Let  $\mathbf{x} \in \mathbb{R}^d$  be an input sample, then the network has  $d$  input nodes. Moreover, we have a series of weight matrices  $\{\mathbf{W}^{[l]}\}_{l=1}^L$ . Note that  $\mathbf{W}^{[1]} \in \mathbb{R}^{m \times d}$  is the first weight matrix, and  $\mathbf{W}^{[l]} \in \mathbb{R}^{m \times m}$  is the weight at the  $l$ -th layer, for  $2 \leq l \leq L$ . Let  $\mathbf{x}^{[l]}$  be the output of layer  $l$ , with  $\mathbf{x}^{[0]} = \mathbf{x}$ . We consider the fully-connected network given below:

$$\begin{aligned} \mathbf{x}^{[l]} &= \frac{1}{\sqrt{m}} \sigma(\mathbf{W}^{[l]} \mathbf{x}^{[l-1]}), \\ f_{\text{nn}}(\mathbf{x}, \boldsymbol{\theta}) &= \mathbf{a}^\top \mathbf{x}^{[L]}, \end{aligned} \tag{1.1}$$

where  $\sigma(\cdot)$  is applied coordinate-wisely to its input, and  $f_{\text{nn}}(\mathbf{x}, \boldsymbol{\theta})$  is the output function. Specifically, for the case in Figure 1.2, the parameters are set to be  $d = 12, m = 16, L = 4$ .

The scaling factor  $1/\sqrt{m}$  is key to obtaining a consistent asymptotic behavior of neural networks as the width  $m$  of the hidden layers grow to infinity [4]. In the infinite-width limit,

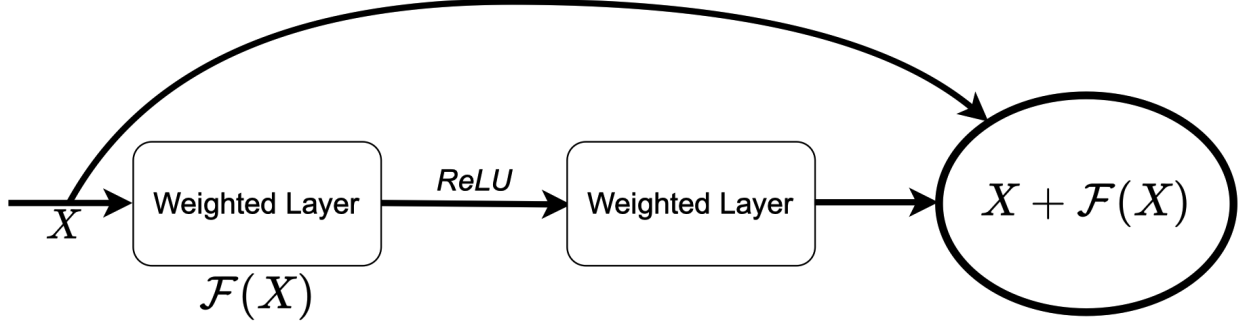
the output function at initialization converges to a Gaussian distribution, and it follows a linear differential equation during training. There are, of course, other scaling factors that are employed extensively. For instance, the mean-field scaling  $1/m$ . For two-layer networks, a line of papers [17]–[19] used optimal transport theory equipped with the mean-field scaling to establish that for infinitely wide neural networks, the empirical distribution of the neural network parameters can be described as a Wasserstein gradient flow. However, their results are limited to two-layer networks and may require an exponential amount of overparametrization. The current situation of neural network study is similar to an early era of statistical mechanics, when we observe different states of a matter at several discrete conditions without the guidance of a unified phase diagram. Therefore, inspired by that, a recently published work [20] presented a systematic and comprehensive analysis in drawing the first phase diagram for two-layer neural networks at the infinite-width limit, in pursuit of a complete characterization of its dynamical regimes and their dependence on different scaling factors.

### 1.3 ResNet

As mentioned both in Section 1.2 and [21], deeper and deeper network architectures are being developed nowadays. There is even a mathematical proof in [22] that reveals the utility of having deeper networks than that of wider networks. However, increasing network depth also introduce the issues of vanishing gradients [23] and degradation [2]. As the gradient is computed out by backward propagation, repeated multiplication with small weights renders it ineffectively small. Vanishing gradients, however, has been largely addressed by some normalizing tricks [24], [25].

When deeper networks are able to start converging, a degradation problem has been exposed: the training accuracy gets saturated and then degrades rapidly [26]. This is counter-intuitive in that by providing our model with more parameters, it shall be able to fit the training data at least as good as its predecessor. Surprisingly, such degradation is not due to overfitting. The degradation of training performances reveals that not all networks are similarly easy to optimize by brutal force. Moreover, the problem suggests that it might be

hard for the solvers to learn identity mappings with multiple nonlinear layers. Ultimately, this conjecture motivates the setup for learning small residuals and directly adding them to the input. The core idea of ResNet [2] is the employment of residual learning building block,



**Figure 1.3.** Residual learning: a building block

as shown in Figure 1.3.

We use the same notations as the fully-connected neural networks to rigorously define our version of ResNet:

$$\begin{aligned}
 \mathbf{x}^{[1]} &= \sqrt{\frac{c_\sigma}{m}} \sigma(\mathbf{W}^{[1]} \mathbf{x}), \\
 \mathbf{x}^{[l]} &= \mathbf{x}^{[l-1]} + \frac{c_{\text{res}}}{L\sqrt{m}} \sigma(\mathbf{W}^{[l]} \mathbf{x}^{[l-1]}), \quad \text{for } 2 \leq l \leq L, \\
 f_{\text{res}}(\mathbf{x}, \boldsymbol{\theta}) &= \mathbf{a}^\top \mathbf{x}^{[L]},
 \end{aligned} \tag{1.2}$$

where the constant  $c_\sigma = \left( \mathbb{E}_{z \sim \mathcal{N}(0,1)} [\sigma(z)^2] \right)^{-1}$  serves as a normalizing factor,  $\mathcal{N}(0,1)$  is the standard Gaussian distribution.  $c_{\text{res}}$  is a small constant satisfying  $0 < c_{\text{res}} < 1$ , and  $f_{\text{res}}(\mathbf{x}, \boldsymbol{\theta})$  is the output function. Note that here we use a  $\frac{c_{\text{res}}}{L}$  factor combined with the  $1/\sqrt{m}$  scaling, the  $\frac{c_{\text{res}}}{L}$  factor guarantees that the width per layer  $m$  does not blow up exponentially with respect to depth  $L$ , intuitively shown in Equation (3.70). Although Equation (1.2) differs by the standard ResNet architecture in [2], it will not be hard to generalize the analysis to architectures with skip-connections at every two or more layers. We also believe that the phase diagram in [20] can be drawn out similarly for ResNet.

## 1.4 Main Results

Owing to the non-convex nature of optimization neural networks, it is challenging to locate the global optima. A popular way to analyze such problems is to identify the geometric properties of each critical point. Some recent works have shown that for the set of functions satisfying:

- all local minima are global;
- every saddle point possesses a negative curvature (i.e. it is non-degenerate),

then gradient descent can find a global optima [27]–[29]. The objective functions of some shallow networks are in such set [14], [30]. However, even for a three-layer linear network, there exists degenerate saddle points without negative curvature [31]. So it is doubtful that global convergence of gradient descent can be ensured for deeper neural networks.

Alternatively, we directly study the dynamics of the gradient descent for specific neural network architectures. This is another approach widely taken to obtain general convergence results. Recently, it has been shown that if the network is overparameterized, gradient descent is able to find a global optima for two-layer networks [7], deep linear networks [14], [32], [33] and ResNet [34]. Jacot et al. [4] established that in the infinite width limit, the full batch gradient descent corresponds to a specific kernel regression predictor. Consequently, in the regime of infinite width, the convergence of gradient descent for neural network can be characterized by a fixed kernel [35]. This is the cornerstone upon which rests the outstanding performance of overparameterization. Inspired by the existence of such kernel, extraordinary efforts have been trying to improve it to the non-asymptotic setting, where only finite width is required.

In the regime of finite width, many works have suggested that the network can reduce training loss at exponential rate using gradient descent [8], [32]. Our thesis also belongs to this category. In this thesis, we contribute to further understanding of the gradient descent dynamics for training fully-connected networks and ResNet models. We use the same ResNet structure as in [1]. Details of the network structure are provided in Section 1.3. More importantly, we assume that the  $n$  data points are not parallel with each other. Such

assumption holds in general for standard dataset, and we focus on the empirical risk minimization problem given by the quadratic loss. We show that if  $m = \Omega(n^3 L^2)$ , then the empirical risk  $R_S(t)$  under gradient descent decays exponentially. More precisely,

$$R_S(t) \leq R_S(0) \exp\left(-\frac{\lambda t}{n}\right),$$

where  $\lambda$  is the least eigenvalue of  $\mathbf{K}^{[L+1]}$ , definition of which can be found in Definition 2.4.5 Equation (2.19). It is worth noticing that:

(1). Given identical ResNet architectures, for the convergence of randomly initialized gradient descent, our results improve upon [1] (Theorem 2.5.2) in the required number of width per layer from  $m = \Omega(n^4 L^2)$  to  $m = \Omega(n^3 L^2)$  (Theorem 3.2.3).

(2). For fully-connected network, the required amount of overparametrization in [3] is  $m = \Omega(n^3 2^{\mathcal{O}(L)})$  (Theorem 2.6.1). We are able to reproduce the result of Du et al. [1], showing that the exponential dependence of  $m$  on the number of layers  $L$  can be eliminated for ResNet, i.e.,  $m = \Omega(n^3 L^2)$  (Theorem 3.2.3).

## 1.5 Assumptions and Notations

We introduce some assumptions and notations that will be used throughout the thesis. We set  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  for the set of input samples, and we assume that:

1. all samples are of uniform length, i.e., for any  $\alpha = 1, 2, \dots, n$ ,  $\|\mathbf{x}_\alpha\|_{L_2} = 1$ ;
2. all samples are non-parallel with each other, i.e.,  $\mathbf{x}_{\alpha_1} \nparallel \mathbf{x}_{\alpha_2}$ , for any  $\alpha_1 \neq \alpha_2$ .

We use  $\sigma(\cdot)$  to denote the activation function, and we assume that:

1.  $\sigma(\cdot)$  is non-polynomial, 1-Lipschitz and (infinitely) smooth;
2. its derivative of any order is also 1-Lipschitz;
3. function value at 0 satisfy  $|\sigma(0)| \leq 1$ .

These assumptions hold for many activation functions, including the soft-plus and sigmoid activation.

Moreover, we set  $n$  for the number of input samples,  $m$  for the width of the neural network, and  $L$  for the number of hidden layers. 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  $\|\cdot\|_{2 \rightarrow \infty}$ . We set a special vector  $(1, 1, 1, \dots, 1)^\top \in \mathbb{R}^m$  by  $\mathbf{1} := (1, 1, 1, \dots, 1)^\top$ . We use  $\mathbf{I}_m$  to signify the identity matrix in  $\mathbb{R}^{m \times m}$ . For a semi-positive-definite matrix  $\mathbf{A}$ , we denote its smallest eigenvalue by  $\lambda_{\min}(\mathbf{A})$ . We use  $\mathcal{O}(\cdot)$  and  $\Omega(\cdot)$  for the standard Big-O and Big-Omega notations. Finally, we use  $\langle \cdot, \cdot \rangle$  to denote the inner product between two vectors or matrices and  $\mathcal{N}(0, 1)$  for the standard Gaussian distribution. These general notations are summarized in Table 1.1.

**Table 1.1.** Notation Table

Symbol	Representations of the symbol
$\mathcal{X}$	The set of input samples
$\sigma(\cdot)$	Nonpolynomial, 1-Lipschitz smooth function
$n$	Number of input samples
$m$	Width of a neural network
$L$	Number of hidden layers
$\ \cdot\ _2$	$L^2$ (Euclidian) norm of a vector
$\ \cdot\ _{2 \rightarrow 2}$	Operator norm of a matrix
$\ \cdot\ _F$	Frobenius norm of a matrix
$\ \cdot\ _\infty$	Infinity norm of a matrix
$\mathbf{1}$	The vector $(1, 1, 1, \dots, 1)^\top$ in $\mathbb{R}^m$
$\mathbf{I}_m$	The identity matrix in $\mathbb{R}^{m \times m}$
$\lambda_{\min}(\mathbf{A})$	Least eigenvalue of a matrix
$\mathcal{N}(0, 1)$	Standard Gaussian distribution
$\langle \cdot, \cdot \rangle$	Inner product
$\mathcal{O}(\cdot)$	Big-O notation
$\Omega(\cdot)$	Big-Omega notation

Next, since we are going to perform massive computations, some useful notations shall also be introduced but not listed out in Table 1.1. We denote  $\sigma(\mathbf{W}^{[l]} \mathbf{x}_\alpha^{[l-1]})$  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  $\boldsymbol{\sigma}_{[l]}^{(r)}(\mathbf{x}_\alpha)$ , where  $r \geq 1$ . We also write the output function  $f(\mathbf{x}_\alpha, \boldsymbol{\theta}_t)$  as  $f_\alpha(t)$ . Moreover, we define a series of special matrices,

$$\mathbf{E}_{t,\alpha}^{[l]} := \left( \mathbf{I}_m + \frac{c_{\text{res}}}{L} \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\alpha) \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right), \quad 2 \leq l \leq L. \quad (1.3)$$

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.,  $\mathbf{E}_{t,\alpha}^{[l:L]} := \mathbf{E}_{t,\alpha}^{[L]}\mathbf{E}_{t,\alpha}^{[L-1]}\cdots\mathbf{E}_{t,\alpha}^{[l]}$ , where we set  $\mathbf{E}_{t,\alpha}^{[i:j]} := \mathbf{I}_m$ ,  $i > j$  for the purpose of completeness.

Finally, we introduce a notion of high probability events that has been commonly used, for instance, see [3, 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  $n \lesssim \text{poly}(m)$ , 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  $n$ .

## 2. GRADIENT DESCENT FINDS GLOBAL MINIMA

### 2.1 Introduction

In this chapter, we begin with an introduction to empirical risk minimization problems. We focus particularly on the empirical risk minimization problem with a quadratic loss:

$$\min_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{\alpha=1}^n \|f(\mathbf{x}_\alpha, \boldsymbol{\theta}) - y_\alpha\|_2^2. \quad (2.1)$$

In the above equation,  $\{\mathbf{x}_\alpha\}_{\alpha=1}^n$  are the training inputs,  $\{y_\alpha\}_{\alpha=1}^n$  are the labels.  $f(\mathbf{x}_\alpha, \boldsymbol{\theta})$  is the prediction function, which in our case is a neural network, and  $\boldsymbol{\theta}$  are the parameters to be optimized.

### 2.2 Gradient Descent

To learn the deep neural network, we introduce the randomly initialized gradient descent algorithm to find the global minimizer of the empirical loss Equation (2.1). The core of the algorithm consists of two steps. The first step incorporates a random initialization of parameters. As the vector containing all parameters is denoted by

$$\boldsymbol{\theta} = \left( \text{vec}(\mathbf{W}^{[L]}), \text{vec}(\mathbf{W}^{[L-1]}), \dots, \text{vec}(\mathbf{W}^{[1]}), \mathbf{a} \right),$$

where  $\text{vec}$  is the standard vectorization operation, we initialize the parameters following the adopted Xavier initialization scheme [25], i.e.,

$$W_{i,j}^{[l]} \sim \mathcal{N}(0, 1), \quad a_k \sim \mathcal{N}(0, 1), \quad 1 \leq l \leq L, \quad 1 \leq i, j, k \leq m. \quad (2.2)$$

The second step is to train all layers of the neural network with continuous time gradient descent (gradient flow): for  $l = 1, 2, \dots, L$ , and time  $t \geq 0$ ,

$$\begin{aligned} \partial_t \mathbf{W}_t^{[l]} &= -\partial_{\mathbf{W}^{[l]}} R_S(\boldsymbol{\theta}_t), \\ \partial_t \mathbf{a}_t &= -\partial_{\mathbf{a}} R_S(\boldsymbol{\theta}_t). \end{aligned} \quad (2.3)$$

A discrete version of gradient descent can be found in other literature [7], where the parameters are updated via

$$\begin{aligned}\mathbf{W}^{[l]}(k) &= \mathbf{W}^{[l]}(k-1) - \eta \frac{\partial R_S(\boldsymbol{\theta}(k-1))}{\partial \mathbf{W}^{[l]}(k-1)}, \\ \mathbf{a}(k) &= \mathbf{a}(k-1) - \eta \frac{\partial R_S(\boldsymbol{\theta}(k-1))}{\partial \mathbf{a}(k-1)},\end{aligned}\tag{2.4}$$

with  $k$  being the index of the step to be taken, and  $\eta > 0$  being the step size.

We shall write out the dynamics (2.3) respectively for fully-connected network and ResNet. Recall that the output function  $f(\mathbf{x}_\beta, \boldsymbol{\theta}_t)$  is denoted by  $f_\beta(t)$ , then for fully-connected network, dynamics (2.3) reads

$$\begin{aligned}\partial_t \mathbf{a}_t &= -\frac{1}{n} \sum_{\beta=1}^n \mathbf{x}_\beta^{[L]} (f_\beta(t) - y_\beta), \\ \partial_t \mathbf{W}_t^{[l]} &= -\frac{1}{n} \sum_{\beta=1}^n \left( \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) \frac{(\mathbf{W}_t^{[l+1]})^\top}{\sqrt{m}} \cdots \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_\beta) \frac{\mathbf{a}_t}{\sqrt{m}} \right) \otimes (\mathbf{x}_\beta^{[l-1]})^\top (f_\beta(t) - y_\beta),\end{aligned}\tag{2.5}$$

and for ResNet, dynamics (2.3) reads

$$\begin{aligned}\partial_t \mathbf{a}_t &= -\frac{1}{n} \sum_{\beta=1}^n \mathbf{x}_\beta^{[L]} (f_\beta(t) - y_\beta), \\ \partial_t \mathbf{W}_t^{[L]} &= -\frac{1}{n} \sum_{\beta=1}^n \frac{c_{\text{res}}}{L\sqrt{m}} \text{diag} \left( \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_\beta) \mathbf{a}_t \right) \mathbf{1} \otimes (\mathbf{x}_\beta^{[L-1]})^\top (f_\beta(t) - y_\beta), \\ &\text{for } l = 2, 3, \dots, L-1, \\ \partial_t \mathbf{W}_t^{[l]} &= -\frac{1}{n} \sum_{\beta=1}^n \frac{c_{\text{res}}}{L\sqrt{m}} \text{diag} \left( \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+1):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \otimes (\mathbf{x}_\beta^{[l-1]})^\top (f_\beta(t) - y_\beta), \\ \partial_t \mathbf{W}_t^{[1]} &= -\frac{1}{n} \sum_{\beta=1}^n \sqrt{\frac{c_\sigma}{m}} \text{diag} \left( \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \otimes (\mathbf{x}_\beta)^\top (f_\beta(t) - y_\beta).\end{aligned}\tag{2.6}$$

A recent line of work tries to understand the optimization process of training deep neural networks from the perspective of over-parameterization and random weight initialization. It has been observed that over-parameterization and proper random initialization can help the optimization in training neural networks, and various theoretical results have been established [32], [36]. Our results mainly build on two ideas from previous works on gradient

descent. First, we use the observation by Du et al. [1] that the required width  $m$  relies heavily on the structure of networks. Second, we applied the neural tangent hierarchy (NTH) to ResNet, a framework initially proposed by Huang and Yau [3], to directly study the change of its neural tangent kernel (NTK) [4].

With these in mind, we start with a review of NTK in Section 2.3. In Section 2.4 we introduce the concept of Gram matrices. In Section 2.5 we give out the outline of analysis and main results obtained by Du et al. [1] without proof, and we treat the counterparts in Huang and Yau [3] similarly in Section 2.6.

### 2.3 NTK

A flurry of recent papers in theoretical deep learning endeavor to tackle the common theme of analyzing neural networks in the infinite-width limit. At first glance, this limit may seem impractical and even pointless to study. But for mathematicians, there is a tradition of deriving insights into questions by studying them in the infinite limit, which usually tends to be easier in theory. As it turns out, neural networks in this regime simplify to linear models with a regression kernel called the NTK [4].

We shall refer the readers to the connection between infinitely wide neural networks and kernel methods [37]. Specifically, for any parametrized function  $f(\mathbf{x}, \boldsymbol{\theta}_t)$  equipped with two inputs  $\mathbf{x}_\alpha, \mathbf{x}_\beta$  ( $\mathbf{x}_\alpha, \mathbf{x}_\beta$  could be identical), the corresponding kernel is

$$\mathcal{G}_{\boldsymbol{\theta}_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \langle \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_\alpha, \boldsymbol{\theta}_t), \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_\beta, \boldsymbol{\theta}_t) \rangle. \quad (2.7)$$

The key difference between the kernel above and the one in [37] is that our kernel is defined through the inner product between the gradients of the function with respect to its parameters, while the counterpart in [37] comprises the product of the output function. Emergence of the gradient arises from the usage of gradient descent.

In the situations where  $f(\mathbf{x}, \boldsymbol{\theta}_t)$  is the output of a fully-connected network or ResNet introduced in Chapter 1, it consists of a series of kernels  $\{\mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)\}_{l=1}^{L+1}$ , i.e.,

$$\begin{aligned}
\mathcal{G}_{\boldsymbol{\theta}_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \langle \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_\alpha, \boldsymbol{\theta}_t), \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_\beta, \boldsymbol{\theta}_t) \rangle \\
&= \langle \partial_{\mathbf{a}} f(\mathbf{x}_\alpha, \boldsymbol{\theta}_t), \partial_{\mathbf{a}} f(\mathbf{x}_\beta, \boldsymbol{\theta}_t) \rangle + \sum_{l=1}^L \langle \partial_{\mathbf{W}^{[l]}} f(\mathbf{x}_\alpha, \boldsymbol{\theta}_t), \partial_{\mathbf{W}^{[l]}} f(\mathbf{x}_\beta, \boldsymbol{\theta}_t) \rangle \\
&:= \mathcal{G}_t^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) + \sum_{l=1}^L \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \\
&= \sum_{l=1}^{L+1} \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta).
\end{aligned} \tag{2.8}$$

For fully-connected network,  $\mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$  individually reads

$$\begin{aligned}
\mathcal{G}_t^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \langle \mathbf{x}_\alpha^{[L]}, \mathbf{x}_\beta^{[L]} \rangle, \\
&\text{for } 1 \leq l \leq L \\
\mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \left\langle \frac{1}{\sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\alpha) \frac{(\mathbf{W}_t^{[l+1]})^\top}{\sqrt{m}} \cdots \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_\alpha) \mathbf{a}_t, \right. \\
&\quad \left. \frac{1}{\sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) \frac{(\mathbf{W}_t^{[l+1]})^\top}{\sqrt{m}} \cdots \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_\beta) \mathbf{a}_t \right\rangle \langle \mathbf{x}_\alpha^{[l-1]}, \mathbf{x}_\beta^{[l-1]} \rangle,
\end{aligned} \tag{2.9}$$

while for ResNet

$$\begin{aligned}
\mathcal{G}_t^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \langle \mathbf{x}_\alpha^{[L]}, \mathbf{x}_\beta^{[L]} \rangle, \\
&\text{for } 2 \leq l \leq L \\
\mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \left\langle \frac{c_{\text{res}}}{L\sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\alpha) (\mathbf{E}_{t,\alpha}^{[(l+1):L]})^\top \mathbf{a}_t, \right. \\
&\quad \left. \frac{c_{\text{res}}}{L\sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) (\mathbf{E}_{t,\beta}^{[(l+1):L]})^\top \mathbf{a}_t \right\rangle \langle \mathbf{x}_\alpha^{[l-1]}, \mathbf{x}_\beta^{[l-1]} \rangle, \\
\mathcal{G}_t^{[1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \left\langle \sqrt{\frac{c_\sigma}{m}} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\alpha) (\mathbf{E}_{t,\alpha}^{[2:L]})^\top \mathbf{a}_t, \right. \\
&\quad \left. \sqrt{\frac{c_\sigma}{m}} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) (\mathbf{E}_{t,\beta}^{[2:L]})^\top \mathbf{a}_t \right\rangle \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle
\end{aligned} \tag{2.10}$$

Up to this point, we have only given out the exact computations of the regression kernel, where the property of infinite width have not been used. In the large width limit, it turns

out that the time-varying kernel  $\mathcal{G}_{\theta_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$  is close to a deterministic kernel  $\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)$ , the limiting NTK. This property is proved in two steps. Firstly, at initial stage ( $t = 0$ ), with appropriate scaling factor  $1/\sqrt{m}$  for the parameters  $\boldsymbol{\theta}_0$ , there exists an infinite width limit ( $m \rightarrow \infty$ ) of  $\mathcal{G}_{\theta_0}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$ , denoted by  $\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)$ . Secondly, the kernel  $\mathcal{G}_{\theta_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$  itself barely changes during the entire training process, i.e.  $\mathcal{G}_{\theta_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \approx \mathcal{G}_{\theta_0}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$ . Hence, as  $m \rightarrow \infty$ ,  $\mathcal{G}_{\theta_t}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \approx \mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)$  for all  $t > 0$ .

The results above allow us to capture the behavior of networks trained by gradient descent. In large width limit, a single output function  $f(\mathbf{x}, \boldsymbol{\theta}_t)$  evolves as a linear differential equation

$$\partial_t (f(\mathbf{x}, \boldsymbol{\theta}_t) - y) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{K}_\infty(\mathbf{x}, \mathbf{x}_\beta) (f(\mathbf{x}_\beta, \boldsymbol{\theta}_t) - y_\beta), \quad (2.11)$$

where  $\mathcal{K}_\infty(\cdot)$  only depends on the training inputs. More importantly, it is independent of the neural network parameters. As a direct consequence of Equation (2.11), we have

$$\begin{aligned} \partial_t \sum_{\alpha=1}^n \|f(\mathbf{x}_\alpha, \boldsymbol{\theta}_t) - y_\alpha\|_2^2 &= -\frac{2}{n} \sum_{\alpha, \beta=1}^n \mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta) (f(\mathbf{x}_\alpha, \boldsymbol{\theta}_t) - y_\alpha) (f(\mathbf{x}_\beta, \boldsymbol{\theta}_t) - y_\beta) \\ &\leq -\frac{2}{n} \lambda_{\min} \left( [\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n} \right) \sum_{\gamma=1}^n \|f(\mathbf{x}_\gamma, \boldsymbol{\theta}_t) - y_\gamma\|_2^2. \end{aligned} \quad (2.12)$$

From Equation (2.12), we observe that the empirical loss (2.1) converges at a linear rate determined by the least eigenvalue of matrix  $[\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n}$ . To sum up, in the regime of infinite width, the empirical loss converges exponentially to zero regardless of the fully-connected or ResNet structure, and its convergence rate relies heavily on the limiting NTK.

## 2.4 Gram Matrices

In linear algebra, the Gram matrix  $\mathbf{K}$  of a set of real vectors  $\{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_n\}$  in an inner product space is the Hermitian matrix of inner products, whose entries are given by  $\mathbf{K}_{ij} := \langle \mathbf{k}_i, \mathbf{k}_j \rangle$ . Consequently, if we denote the column matrix of  $\{\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_n\}$  by  $\mathbf{D}$ , then  $\mathbf{K} = \mathbf{D}^\top \mathbf{D}$ .

Jacot et al. [4] pointed out that convergence of the empirical loss (2.1) is related to the positive-definiteness of the limiting NTK. For any parametrized function  $f(\mathbf{x}, \boldsymbol{\theta}_t)$ , its limiting NTK reads

$$\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \mathbb{E}_{\boldsymbol{\theta}_0 \sim \mathcal{W}} \langle \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_\alpha, \boldsymbol{\theta}_0), \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_\beta, \boldsymbol{\theta}_0) \rangle,$$

where  $\mathbf{x}_\alpha, \mathbf{x}_\beta$  are two inputs, and  $\mathcal{W}$  is the initial distribution over  $\boldsymbol{\theta}$ .

In our cases, where  $f(\mathbf{x}, \boldsymbol{\theta}_t)$  is the output of a fully-connected network or ResNet, and  $\boldsymbol{\theta} = (\text{vec}(\mathbf{W}^{[L]}), \text{vec}(\mathbf{W}^{[L-1]}), \dots, \text{vec}(\mathbf{W}^{[1]}), \mathbf{a})$ , its limiting NTK is the sum of a series of kernels  $\{\mathcal{K}^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)\}_{l=1}^{L+1}$ , i.e.,

$$\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \sum_{l=1}^{L+1} \mathcal{K}^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta),$$

where

$$\begin{aligned} \mathcal{K}^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \mathbb{E}_{\boldsymbol{\theta}_0 \sim \mathcal{W}} \langle \partial_{\mathbf{a}} f(\mathbf{x}_\alpha, \boldsymbol{\theta}_0), \partial_{\mathbf{a}} f(\mathbf{x}_\beta, \boldsymbol{\theta}_0) \rangle, \\ \mathcal{K}^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \mathbb{E}_{\boldsymbol{\theta}_0 \sim \mathcal{W}} \langle \partial_{\mathbf{W}^{[l]}} f(\mathbf{x}_\alpha, \boldsymbol{\theta}_0), \partial_{\mathbf{W}^{[l]}} f(\mathbf{x}_\beta, \boldsymbol{\theta}_0) \rangle, \quad 1 \leq l \leq L. \end{aligned}$$

We remark that at  $t = 0$ , respectively for all  $l$ , the limiting values of  $\mathcal{G}_0^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$ , whose definition can be found in Equation (2.9) and Equation (2.10), equal to  $\mathcal{K}^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)$ . This is a crucial finding in [38, Corollary 2.4.]. As is shown in Equation (2.12), dynamics of the empirical loss is governed by spectral property of matrix  $[\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n}$ . With some abuse of notations, we name matrix  $\mathbf{K}_\infty := [\mathcal{K}_\infty(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n}$  by *Gram matrix*. Also with additional abuse of terminologies, for  $1 \leq l \leq L+1$ , the matrices  $\mathbf{K}^{[l]} := [\mathcal{K}^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n}$  are entitled *Gram matrices*, whose definitions can be traced back to [1, Definition 5.1, Definition 6.1].

As stated in [38, Corollary 2.4], to assure the convergence of NTK,  $\boldsymbol{\theta}$  is required to be suitably randomized. As mentioned in Section 2.2, we initialize  $\boldsymbol{\theta}$  using the adopted Xavier initialization scheme (2.2). Furthermore, it shall be noted that derivation of the Gram matrices depends on the series of matrices  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$  and vectors  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$  to be defined immediately in the next two sections.



### 2.4.1 Gram Matrices for Fully-connected Networks

In this section, as a warm up, we give out definitions of Gram matrices generated by fully-connected networks. We are able to write out the explicit formulas for  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$ . However, it is not the case for  $\{\mathbf{K}^{[l]}\}_{l=1}^{L-1}$ . A slightly different approach shall be taken to write out the expressions for them.

**Definition 2.4.1.** *Given input samples  $\mathcal{X}$  and activation function  $\sigma(\cdot)$ , the matrices  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$  are defined recursively: for  $1 \leq i, j \leq n, 1 \leq l \leq L$*

$$\begin{aligned}\widetilde{\mathbf{K}}_{ij}^{[0]} &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle, \\ \widetilde{\mathbf{A}}_{ij}^{[l]} &= \begin{pmatrix} \widetilde{\mathbf{K}}_{ii}^{[l-1]} & \widetilde{\mathbf{K}}_{ij}^{[l-1]} \\ \widetilde{\mathbf{K}}_{ji}^{[l-1]} & \widetilde{\mathbf{K}}_{jj}^{[l-1]} \end{pmatrix}, \\ \widetilde{\mathbf{K}}_{ij}^{[l]} &= \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{A}}_{ij}^{[l]})} \sigma(u) \sigma(v), \\ \widetilde{\mathbf{A}}_{ij}^{[L+1]} &= \begin{pmatrix} \widetilde{\mathbf{K}}_{ii}^{[L]} & \widetilde{\mathbf{K}}_{ij}^{[L]} \\ \widetilde{\mathbf{K}}_{ji}^{[L]} & \widetilde{\mathbf{K}}_{jj}^{[L]} \end{pmatrix}.\end{aligned}\tag{2.13}$$

**Definition 2.4.2** (Gram Matrices  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$  for Fully-connected Networks). *Given  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$ , Gram matrices  $\mathbf{K}^{[L+1]}, \mathbf{K}^{[L]} \in \mathbb{R}^{n \times n}$  are defined as follows: for  $1 \leq i, j \leq n$*

$$\mathbf{K}_{ij}^{[L+1]} = \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{A}}_{ij}^{[L+1]})} \sigma(u) \sigma(v),\tag{2.14}$$

$$\mathbf{K}_{ij}^{[L]} = \widetilde{\mathbf{K}}_{ij}^{[L-1]} \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{A}}_{ij}^{[L]})} [\sigma^{(1)}(u) \sigma^{(1)}(v)].\tag{2.15}$$

**Definition 2.4.3.** *Given  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$ , Gram matrices  $\mathbf{K}^{[l]} \in \mathbb{R}^{n \times n}$  are defined as follows: for  $1 \leq i, j \leq n, 1 \leq l \leq L-1$*

$$\begin{aligned}\mathbf{K}_{ij}^{[l]} &= \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_i) \frac{(\mathbf{W}_0^{[l+1]})^\top}{\sqrt{m}} \cdots \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_i) \mathbf{a}_0, \right. \\ &\quad \left. \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_j) \frac{(\mathbf{W}_0^{[l+1]})^\top}{\sqrt{m}} \cdots \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_j) \mathbf{a}_0 \right\rangle \langle \mathbf{x}_i^{[l-1]}, \mathbf{x}_j^{[l-1]} \rangle.\end{aligned}\tag{2.16}$$

## 2.4.2 Gram Matrices for ResNet

We give out the derivation of Gram matrices generated by ResNet with slight abuse of notations. Similarly, only  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$  can be written into closed forms.

**Definition 2.4.4.** *Given input samples  $\mathcal{X}$  and activation function  $\sigma(\cdot)$ , the matrices  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$  and vectors  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$  are defined recursively: for  $1 \leq i, j \leq n, 2 \leq l \leq L$*

$$\begin{aligned}
\widetilde{\mathbf{K}}_{ij}^{[0]} &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle, \\
\widetilde{\mathbf{A}}_{ij}^{[1]} &= \begin{pmatrix} \widetilde{\mathbf{K}}_{ii}^{[0]} & \widetilde{\mathbf{K}}_{ij}^{[0]} \\ \widetilde{\mathbf{K}}_{ji}^{[0]} & \widetilde{\mathbf{K}}_{jj}^{[0]} \end{pmatrix}, \\
\widetilde{\mathbf{K}}_{ij}^{[1]} &= \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{A}}_{ij}^{[1]})} c_\sigma \sigma(u) \sigma(v), \\
\widetilde{\mathbf{b}}_i^{[1]} &= \sqrt{c_\sigma} \mathbb{E}_{u \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{K}}_{ii}^{[0]})} [\sigma(u)], \\
\widetilde{\mathbf{A}}_{ij}^{[l]} &= \begin{pmatrix} \widetilde{\mathbf{K}}_{ii}^{[l-1]} & \widetilde{\mathbf{K}}_{ij}^{[l-1]} \\ \widetilde{\mathbf{K}}_{ji}^{[l-1]} & \widetilde{\mathbf{K}}_{jj}^{[l-1]} \end{pmatrix}, \\
\widetilde{\mathbf{K}}_{ij}^{[l]} &= \widetilde{\mathbf{K}}_{ij}^{[l-1]} + \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{A}}_{ij}^{[l]})} \left[ \frac{c_{\text{res}} \widetilde{\mathbf{b}}_i^{[l-1]} \sigma(v)}{L} + \frac{c_{\text{res}} \widetilde{\mathbf{b}}_j^{[l-1]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(v)}{L^2} \right], \\
\widetilde{\mathbf{b}}_i^{[l]} &= \widetilde{\mathbf{b}}_i^{[l-1]} + \frac{c_{\text{res}}}{L} \mathbb{E}_{u \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{K}}_{ii}^{[l-1]})} [\sigma(u)], \\
\widetilde{\mathbf{A}}_{ij}^{[L+1]} &= \begin{pmatrix} \widetilde{\mathbf{K}}_{ii}^{[L]} & \widetilde{\mathbf{K}}_{ij}^{[L]} \\ \widetilde{\mathbf{K}}_{ji}^{[L]} & \widetilde{\mathbf{K}}_{jj}^{[L]} \end{pmatrix}.
\end{aligned} \tag{2.17}$$

**Definition 2.4.5** (Gram Matrices  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$  for ResNet). *Given  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$ ,  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , Gram matrices  $\mathbf{K}^{[L+1]}, \mathbf{K}^{[L]} \in \mathbb{R}^{n \times n}$  are defined as follows: for  $1 \leq i, j \leq n$*

$$\mathbf{K}_{ij}^{[L+1]} = \widetilde{\mathbf{K}}_{ij}^{[L]} + \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{A}}_{ij}^{[L+1]})} \left[ \frac{c_{\text{res}} \widetilde{\mathbf{b}}_i^{[L]} \sigma(v)}{L} + \frac{c_{\text{res}} \widetilde{\mathbf{b}}_j^{[L]} \sigma(u)}{L} + \frac{c_{\text{res}}^2 \sigma(u) \sigma(v)}{L^2} \right], \tag{2.18}$$

$$\mathbf{K}_{ij}^{[L]} = \frac{c_{\text{res}}^2}{L^2} \widetilde{\mathbf{K}}_{ij}^{[L-1]} \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{A}}_{ij}^{[L]})} [\sigma^{(1)}(u) \sigma^{(1)}(v)]. \tag{2.19}$$

**Definition 2.4.6.** Given  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$ ,  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , Gram matrices  $\mathbf{K}^{[l]} \in \mathbb{R}^{n \times n}$  are defined as follows, for  $1 \leq i, j \leq n, 2 \leq l \leq L-1$ ,

$$\mathbf{K}_{ij}^{[l]} = \frac{c_{\text{res}}^2}{L^2} \widetilde{\mathbf{K}}_{ij}^{[l-1]} \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_i) \left( \mathbf{E}_{0,i}^{[(l+1):L]} \right)^\top \mathbf{a}_0, \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_j) \left( \mathbf{E}_{0,j}^{[(l+1):L]} \right)^\top \mathbf{a}_0 \right\rangle, \quad (2.20)$$

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

$$\mathbf{K}_{ij}^{[1]} = c_\sigma \widetilde{\mathbf{K}}_{ij}^{[0]} \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_i) \left( \mathbf{E}_{0,i}^{[2:L]} \right)^\top \mathbf{a}_0, \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_j) \left( \mathbf{E}_{0,j}^{[2:L]} \right)^\top \mathbf{a}_0 \right\rangle. \quad (2.21)$$

We remark that the  $\frac{1}{m}$  scalings in (2.16), (2.20) and (2.21) originate from the inner product between the gradients, i.e.,  $\frac{1}{m} = \frac{1}{\sqrt{m}} \times \frac{1}{\sqrt{m}}$ . Thanks to the Strong Law of Large Numbers, the above limits (2.16), (2.20) and (2.21) exist [38, Corollary 2.4]. As we send  $m \rightarrow \infty$ , Gram matrices  $\{\mathbf{K}^{[l]}\}_{l=1}^{L-1}$  only depend on the input samples and the activation patterns. Hence we conclude that all Gram matrices  $\{\mathbf{K}^{[l]}\}_{l=1}^{L+1}$  only depend on the input samples and independent of  $\boldsymbol{\theta}$ .

## 2.5 Main Results of Du et al.

Jacot et al. [4, Proposition 2] proved the positive-definiteness of the limiting NTK when the data is supported on the sphere and the non-linearity is non-polynomial. Du et al. [1, Proposition F.1, Proposition F.2] extended their results and showed that as long as the input training data is not degenerate and supported on the unit sphere,  $\lambda_{\min}(\mathbf{K}^{[L]})$  is strictly positive. Since gram matrices of all orders are positive semi-definite,  $\lambda_{\min}(\mathbf{K}^{[L]})$  is an explicit lower bound of the least eigenvalue of the limiting NTK matrix  $\mathbf{K}_\infty$ , i.e.,  $\lambda_{\min}(\mathbf{K}_\infty) \geq \lambda_{\min}(\mathbf{K}^{[L]})$ . Moreover, using the contribution of all the gram matrices to the minimum eigenvalue can potentially improve the convergence rate, as is shown in [3, Corollary 2.5].

The high-level analysis framework of Du et al. consists of mainly two components. At first stage, they showed that at  $t = 0$ ,  $\left[ \mathcal{G}_0^{[L]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n}$  is close to  $\mathbf{K}^{[L]}$  via repeated application of concentration inequality. Instead of sending the width  $m$  of every layer to  $\infty$ , as is the setting for Jacot et al., one only needs it to be greater than a finite threshold of

order  $\Omega(n^2)$  in order for  $\lambda_{\min}\left(\left[\mathcal{G}_0^{[L]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)\right]_{1 \leq \alpha, \beta \leq n}\right)$  to maintain a lower bound with high probability. We observe that  $\mathbf{K}^{[L]}$  is recursively defined, and so is  $\left[\mathcal{G}_0^{[L]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)\right]_{1 \leq \alpha, \beta \leq n}$ . Due to the randomness inherited from the initialization scheme and the introduction of finite threshold of  $m$ , inevitably we have some perturbation in the first layer, and how this perturbation propagates to the  $L$ -th layer shall be analyzed carefully. Du et al. derived a general formulation that allows readers to analyze the initialization behavior for the fully-connected network, ResNet, and even convolutional ResNet in a unified way.

One important finding in the perturbation analysis is that ResNet architecture makes the propagation more stable. For fully-connected network, such perturbation propagates to the  $L$ -th layer exponentially, hence forcing the threshold of  $m$  to maintain exponential dependency on the depth  $L$ . Heuristically speaking, let  $\mathcal{E}_1$  be the perturbation in the first layer, then  $\mathcal{E}_L$ , perturbation in the  $L$ -th layer, admits the form

$$\mathcal{E}_L \leq 2^{\mathcal{O}(L)} \mathcal{E}_1. \quad (2.22)$$

However, for ResNet, thanks to the skip connection structure, its counterpart reads

$$\mathcal{E}_L \leq \left(1 + \mathcal{O}\left(\frac{1}{L}\right)\right)^L \mathcal{E}_1. \quad (2.23)$$

Therefore, the issue of exponential explosion can be avoided, and the above analysis sheds light on the benefit of using ResNet architecture for training.

At the second stage, for  $t > 0$ , the averaged Frobenius norm  $\frac{1}{\sqrt{m}} \left\| \mathbf{W}_t^{[l]} - \mathbf{W}_0^{[l]} \right\|_{\text{F}}$  is used to control the absolute change of eigenvalues, so that the lower bound of the eigenvalue of matrix  $\left[\mathcal{G}_t^{[L]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)\right]_{1 \leq \alpha, \beta \leq n}$  can be guaranteed during the whole training process. In purpose of bounding the averaged Frobenius norm, another threshold of order  $\Omega(n^4)$  is required for width  $m$ . Such analysis, whose high-level intuition is similar to (3.70), once again sheds light on the benefit of using ResNet architecture for training.

We may proceed to state the main theorems of Du et al. [1, Theorem 5.1, Theorem 6.1].

**Theorem 2.5.1** (Convergence of Gradient Descent for Fully-connected Networks). *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$  and  $\mathbf{K}^{[L]}$  defined in Equation (2.15), with high probability w.r.t random initialization, for width*

$$m = \Omega \left( 2^{\mathcal{O}(L)} \max \left\{ \frac{n^4}{\lambda_{\min}^4(\mathbf{K}^{[L]})}, n, \frac{n^2 \log(Ln)}{\lambda_{\min}^2(\mathbf{K}^{[L]})} \right\} \right), \quad (2.24)$$

*and if we set the step size*

$$\eta = \mathcal{O} \left( \frac{\lambda_{\min}(\mathbf{K}^{[L]})}{n^2 2^{\mathcal{O}(L)}} \right),$$

*then for  $k = 0, 1, 2, \dots$ , the loss at each iteration satisfies*

$$R_S(\boldsymbol{\theta}(k)) \leq \left( 1 - \eta \frac{\lambda_{\min}(\mathbf{K}^{[L]})}{2} \right)^k R_S(\boldsymbol{\theta}(0)).$$

Main assumption of Equation (2.24) is that a large enough width for each layer is required. We notice that the requirement of  $m$  has three terms. The first term is used to show the Gram matrix remain stable during training. The second term is used to guarantee the output in each layer is approximately normalized at initial phase. The third term is used to show the perturbation of Gram matrix at initial phase is small. However, its dependency on depth  $L$  is exponential. Such exponentiality comes from the instability of the fully-connected architecture. In the next theorem, equipped with ResNet architecture, dependency on  $L$  can be reduced from  $2^L$  to  $\text{poly}(L)$ .

**Theorem 2.5.2** (Convergence of Gradient Descent for ResNet). *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$  and  $\mathbf{K}^{[L]}$  defined in Equation (2.19), with high probability w.r.t random initialization, for width*

$$m = \Omega \left( \max \left\{ \frac{n^4}{\lambda_{\min}^4(\mathbf{K}^{[L]})}, \frac{n^2}{\lambda_{\min}^2(\mathbf{K}^{[L]}) L^2}, n, \frac{n^2 \log(Ln)}{\lambda_{\min}^2(\mathbf{K}^{[L]})} \right\} \right), \quad (2.25)$$

*and if we set the step size*

$$\eta = \mathcal{O} \left( \frac{\lambda_{\min}(\mathbf{K}^{[L]}) L^2}{n^2} \right),$$

then for  $k = 0, 1, 2, \dots$ , the loss at each iteration satisfies

$$R_S(\boldsymbol{\theta}(k)) \leq \left(1 - \eta \frac{\lambda_{\min}(\mathbf{K}^{[L]})}{2}\right)^k R_S(\boldsymbol{\theta}(0)).$$

In contrast to Theorem 2.5.1, the required width  $m$  is fully polynomial in  $n$  and  $L$ . The main reason why the problematic exponential explosion can be circumvented is that the skip connection blocks enable the network structure to be more stable at both the initialization and the training phase. The requirement on  $m$  has four terms. The first two terms are used to show the stability of Gram matrix during training. The third term is used to assure that the output in each layer is approximately normalized at the initialization phase. The fourth term is used to bound the size of the perturbation of the Gram matrix at initial stage.

## 2.6 Main Results of Huang and Yau

Huang and Yau [3] proposed a framework in which an infinite hierarchy of ordinary differential equations, the neural tangent hierarchy (NTH), is derived.

**Theorem 2.6.1** (NTH for Fully-connected Networks). *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , with high probability w.r.t random initialization, there exists an infinite family of operators  $\mathcal{G}_t^{(r)} : \mathcal{X}^r \rightarrow \mathbb{R}$ , where  $r \geq 2$ , that describes the continuous time gradient descent:*

$$\partial_t(f_\alpha(t) - y_\alpha) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta)(f_\beta(t) - y_\beta), \quad (2.26)$$

and for any  $r \geq 2$ ,

$$\partial_t \mathcal{G}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r}) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{G}_t^{(r+1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r}, \mathbf{x}_\beta)(f_\beta(t) - y_\beta). \quad (2.27)$$

There also exists a deterministic family of operators (independent of  $m$ ):

$$\mathfrak{G}^{(r)} : \mathcal{X}^r \rightarrow \mathbb{R}, \quad 2 \leq r \leq p+1,$$

where  $\mathfrak{G}^{(r)} \equiv 0$  if  $r$  is odd, and some constants  $C, C^* > 0$ , such that with high probability w.r.t random initialization,

$$\left\| \left( \mathcal{G}_0^{(r)} - \frac{\mathfrak{G}^{(r)}}{m^{r/2-1}} \right) (\cdot) \right\|_\infty \lesssim \frac{(\ln m)^C}{m^{(r-1)/2}}, \quad (2.28)$$

and for time  $0 \leq t \leq m^{\frac{p}{2(p+1)}} / (\ln m)^{C^*}$ ,

$$\left\| \mathcal{G}_t^{(r)} (\cdot) \right\|_\infty \lesssim \frac{(\ln m)^C}{m^{r/2-1}}. \quad (2.29)$$

For fully-connected network defined in Equation (1.1), with width  $m = \Omega(2^{\mathcal{O}(L)} n^3)$ , gradient descent converges at a linear rate.

**Theorem 2.6.2.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , we further assume that there exists  $\lambda > 0$  (might depend on  $m$ )*

$$\lambda_{\min} \left( \left[ \mathcal{G}_0^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} \right) \geq \lambda, \quad (2.30)$$

and the width  $m$  of the neural network satisfies

$$m \geq C^* \left( \frac{n}{\lambda} \right)^3 (\ln m)^C \ln \left( \frac{n}{\varepsilon} \right)^2, \quad (2.31)$$

for some constants  $C, C^* > 0$ . Then with high probability w.r.t. random initialization, the training error decays exponentially,

$$\sum_{\alpha=1}^n \|f_\alpha(t) - y_\alpha\|_2^2 \leq \exp \left( -\frac{\lambda t}{n} \right) \sum_{\alpha=1}^n \|f_\alpha(0) - y_\alpha\|_2^2, \quad (2.32)$$

which reaches the training accuracy  $\varepsilon$  with time complexity

$$T = \mathcal{O} \left( \frac{n}{\lambda} \ln \left( \frac{1}{\varepsilon} \right) \right). \quad (2.33)$$

We remark that the operator  $\mathcal{G}_t^{(2)}(\cdot)$  by definition is the same as the NTK  $\mathcal{G}_{\theta_t}(\cdot)$  in Equation (2.7). Also in Theorem 2.6.2, a further assumption on the least eigenvalue of the NTK  $\mathcal{G}_0^{(2)}(\cdot)$  has been imposed directly, see Equation (2.30). However, as is shown in The-

orem 2.5.1, there shall be some additional requirements on width  $m$  in order for the NTK at initial phase to be positive-definite. We complement this technical issue in Chapter B, showing that the least eigenvalue of the NTK  $\mathcal{G}_0^{(2)}(\cdot)$  can be ensured for a finite threshold of the width  $m$ .



### 3. RESNET USING NTH

#### 3.1 Introduction

In this chapter, we study the ResNet model (1.2) with finite width using NTH. We rigorously prove that as long as no two training inputs from  $\mathcal{X}$  are parallel and the width  $m$  is large enough, gradient descent achieves zero training loss at a linear convergence rate, i.e., it finds a solution  $\boldsymbol{\theta}(t)$  with  $R_S(\boldsymbol{\theta}(t)) \leq \varepsilon$  in time  $t = \mathcal{O}\left(\ln\left(\frac{1}{\varepsilon}\right)\right)$ . Thus, our results give a quantitative convergence rate involving the desired accuracy.

Our exposition mainly follows the settings of [3]. However, different from [3] in analyzing the fully-connected network, we focus on the investigation of ResNet. We exploit further benefits of using ResNet architecture for training and advantages of choosing NTH over kernel regression.

#### 3.2 Main Results

**Theorem 3.2.1** (NTH for ResNet). *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , there exists an infinite family of operators  $\mathcal{G}_t^{(r)} : \mathcal{X}^r \rightarrow \mathbb{R}$ ,  $r \geq 2$  that describes the continuous time gradient descent*

$$\partial_t(f_\alpha(t) - y_\alpha) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta)(f_\beta(t) - y_\beta), \quad (3.1)$$

and for any  $r \geq 2$ ,

$$\partial_t \mathcal{G}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r}) = -\frac{1}{n} \sum_{\beta=1}^n \mathcal{G}_t^{(r+1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r}, \mathbf{x}_\beta)(f_\beta(t) - y_\beta), \quad (3.2)$$

then with high probability w.r.t random initialization, there exist some constants  $C, C^* > 0$ , such that for  $r \geq 2$ , time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

$$\begin{aligned} \|\mathcal{G}_t^{(2)}(\cdot)\|_\infty &\lesssim 1, \\ \|\mathcal{G}_t^{(r)}(\cdot)\|_\infty &\lesssim \frac{(\ln m)^C}{m^{r/2-1}}. \end{aligned} \quad (3.3)$$

Following Theorem 3.2.1, we shall remark that:

- The operator  $\mathcal{G}_t^{(2)}(\cdot)$  is the same as the NTK  $\mathcal{G}_{\theta_t}(\cdot)$  derived in Equation (2.10);
- Constant  $C$  depends linearly on  $r$ ;
- Pre-factors in Inequality (3.3) explode exponentially fast in  $r$ .

Even though the pre-factors explode exponentially, this does not exert influence on the convergence of gradient descent. Since the landscape of empirical loss  $R_S(\theta_t)$  is mainly affected by lower order kernels. As is shown in the proof of Theorem 3.2.2, we only need to analyze kernels up to order  $r = 4$ .

It has been proved in Theorem 3.2.1 and other literatures [7], [35], [38] that the change of NTK during the gradient descent dynamics for Deep Neural Network is bounded by  $\mathcal{O}\left(\frac{1}{\sqrt{m}}\right)$ . However, it was observed by Lee et al. [39] that time variation of the NTK is closer to  $\mathcal{O}\left(\frac{1}{m}\right)$ , 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 in [3, Corollary 2.4], and we present a different approach to obtain similar results.

Recall that the NTK  $\mathcal{G}_t^{(2)}(\cdot)$  consists of a series of kernels  $\left\{\mathcal{G}_t^{[l]}(\cdot)\right\}_{l=1}^{L+1}$ ,

$$\mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \sum_{l=1}^{L+1} \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta),$$

and for ResNet

$$\begin{aligned} \mathcal{G}_t^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \left\langle \mathbf{x}_\alpha^{[L]}, \mathbf{x}_\beta^{[L]} \right\rangle, \\ &\text{for } 2 \leq l \leq L \\ \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \left\langle \frac{c_{\text{res}}}{L\sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\alpha) \left(\mathbf{E}_{t,\alpha}^{[(l+1):L]}\right)^\top \mathbf{a}_t, \right. \\ &\quad \left. \frac{c_{\text{res}}}{L\sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) \left(\mathbf{E}_{t,\beta}^{[(l+1):L]}\right)^\top \mathbf{a}_t \right\rangle \left\langle \mathbf{x}_\alpha^{[l-1]}, \mathbf{x}_\beta^{[l-1]} \right\rangle, \\ \mathcal{G}_t^{[1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) &= \left\langle \sqrt{\frac{c_\sigma}{m}} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\alpha) \left(\mathbf{E}_{t,\alpha}^{[2:L]}\right)^\top \mathbf{a}_t, \right. \\ &\quad \left. \sqrt{\frac{c_\sigma}{m}} \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left(\mathbf{E}_{t,\beta}^{[2:L]}\right)^\top \mathbf{a}_t \right\rangle \left\langle \mathbf{x}_\alpha, \mathbf{x}_\beta \right\rangle. \end{aligned} \tag{3.4}$$

**Theorem 3.2.2.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , with high probability w.r.t random initialization, there exist some constants  $C, C^* > 0$ , such that for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,*

$$\left\| \partial_t \mathcal{G}_t^{[L+1]}(\cdot) \right\|_\infty \lesssim \frac{(1+t)(\ln m)^C}{m}, \quad (3.5)$$

where the constant  $C$  is independent of the depth  $L$ . Moreover, the pre-factor in the inequality (3.5) is at most of order  $\mathcal{O}(L^2)$ .

As a direct consequence of Theorem 3.2.2, for ResNet defined in Equation (1.2), with width  $m = \Omega(n^3 L^2)$ , gradient descent converges at a linear rate. Precise statements are given out below.

**Theorem 3.2.3.** *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$  and  $\mathbf{K}^{[L+1]}$  defined in Equation (2.18), set  $\lambda_0 > 0$  yielding  $\lambda_0 \leq \lambda_{\min}(\mathbf{K}^{[L+1]})$ , there exists a small constant  $\gamma_1 > 0$ , such that with high probability w.r.t random initialization, for  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\gamma_1}\right)$ ,*

$$\lambda_{\min}\left(\left[\mathcal{G}_0^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta)\right]_{1 \leq \alpha, \beta \leq n}\right) \geq \frac{3}{4}\lambda_0. \quad (3.6)$$

Furthermore, there exists a small constant  $\gamma_2 > 0$ , such that for  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{3+\gamma_2} L^2 \ln\left(\frac{1}{\varepsilon}\right)^2\right)$ , the quadratic training loss  $R_S(\boldsymbol{\theta}_t)$  decays exponentially

$$R_S(\boldsymbol{\theta}_t) \leq R_S(\boldsymbol{\theta}_0) \exp\left(-\frac{\lambda_0 t}{n}\right), \quad (3.7)$$

where  $\varepsilon > 0$  is the desired accuracy of  $R_S(\boldsymbol{\theta}_t)$ .

First of all, we note that positive-definiteness of  $\mathbf{K}^{[L+1]}$  is guaranteed from results in Chapter A. For convenience, we summarize Theorem 3.2.3 as follows. If width  $m$  satisfies that  $m = \max\{\Omega(n^2), \Omega(n^3 L^2)\}$ , then the continuous time gradient descent converges exponentially, and it reaches training accuracy  $\varepsilon$  with time complexity  $T = \mathcal{O}\left(n \ln\left(\frac{1}{\varepsilon}\right)\right)$ .

Before we end this section, we present a fair comparison of our results with others. First of all, Du et al. [1, Theorem 6.1] require  $m = \Omega\left(\frac{n^4}{\lambda_{\min}(\mathbf{K}^{[L]})^4 L^6}\right)$ . Since there is a scaling factor  $\frac{1}{L^2}$  in  $\lambda_{\min}(\mathbf{K}^{[L]})$ , this leads to  $m = \Omega(n^4 L^2)$ . Their iteration complexity for discrete

time gradient descent converges with  $T = \Omega\left(n^2 L^2 \ln\left(\frac{1}{\varepsilon}\right)\right)$ . Our Theorem 3.2.3 improves their result in two ways:

- (i) The quartic dependence on  $n$  is reduced directly to cubic dependence;
- (ii) Faster convergence of the training process of gradient descent.

Secondly, our work serves as an extension of the results established in [3] from fully-connected network to ResNet. We show that for ResNet, in one hand, it is possible for us to study directly the time variation of NTK using NTH. In the other hand, compared with fully-connected network, ResNet is more stable in many aspects. Our Theorem 3.2.3 improves the results in [3] in three ways: (i) With ResNet architecture, the dependency of the amount of over-parameterization on the depth  $L$  can be reduced from  $2^{\mathcal{O}(L)}$  to  $L^2$ ; (ii) While the time interval for the result in [3] 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  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ . Additionally, we are able to show even further that the results hold true for  $t = \infty$ ; (iii) In the proof of Corollary 2.5. of [3], further assumptions on least eigenvalue of the NTK at initial stage were imposed directly. We have rigorously shown in Chapter A and Chapter B that least eigenvalue of the NTK  $\mathcal{G}_0^{(2)}(\cdot)$  stays strictly positive as long as the width  $m$  satisfies  $m = \Omega(n^2)$ . Moreover, for fully-connected networks, Huang and Yau asserted that adding up the whole  $L + 1$  kernels would give rise to the convergence rate of  $R_S(\boldsymbol{\theta})$ , for the belief that the sum of the least eigenvalues of all the kernels  $\mathcal{G}_t^{[l]}$  is much larger than the counterpart of a single kernel, i.e.,

$$\begin{aligned} \lambda_{\min} \left( \sum_{l^*=1}^{L+1} [\mathcal{G}_t^{[l^*]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n} \right) &= \lambda_{\min} \left( [\mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n} \right) \\ &\gg \lambda_{\min} \left( [\mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)]_{1 \leq \alpha, \beta \leq n} \right). \end{aligned}$$

However, for ResNet, even if we assume straightforward that all kernels  $\mathcal{G}_t^{[l]}$  are positive definite, adding them up will not give substantial increase to the least eigenvalue. On account of the fact that there exists a scaling factor  $\frac{1}{L^2}$  for kernels  $\{\mathcal{G}_t^{[l]}\}_{l=2}^L$ , heuristically, the gap of the least eigenvalues between  $\mathcal{G}_t^{(2)}(\cdot)$  and  $\mathcal{G}_t^{[L+1]}(\cdot) + \mathcal{G}_t^{[1]}(\cdot)$  is at most of order  $\mathcal{O}\left(\frac{L-1}{L^2}\right) = \mathcal{O}\left(\frac{1}{L}\right)$ . Hence for ResNet, even if it goes really deep, the least eigenvalue of its

NTK still ‘concentrates’ on the kernels  $\mathcal{G}_t^{[L+1]}(\cdot)$  and  $\mathcal{G}_t^{[1]}(\cdot)$ . Thanks to that observation, we only need to bring  $\mathcal{G}_t^{[L+1]}(\cdot)$  to the spotlight. Analysis of  $\mathcal{G}_t^{[1]}(\cdot)$  is omitted because it is not needed for our proof.

### 3.3 Key Technique Number One: Kernel Structure

The core idea of this technique is to simply take derivatives.

#### 3.3.1 Replacement Rules

We revisit the NTK

$$\mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) = \sum_{l=1}^{L+1} \mathcal{G}_t^{[l]}(\mathbf{x}_\alpha, \mathbf{x}_\beta),$$

where  $\mathcal{G}_t^{(2)}(\cdot)$  is the sum of  $L + 1$  terms, with each term being the inner product of vectors containing components  $\mathbf{a}_t$ ,  $\mathbf{x}_\alpha^{[l]}$ ,  $\mathbf{E}_{t,\alpha}^{[l]}$ , and  $\boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\alpha)$ , as is shown in Equation (3.4). Following Equation (2.6), we are able to write down the gradient dynamics of  $\mathbf{a}_t$ ,  $\mathbf{x}_\alpha^{[l]}$ ,  $\mathbf{E}_{t,\alpha}^{[l]}$ , and  $\boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\alpha)$ . As is shown in Equation (2.6), whenever we take derivatives over a specific term, its expression heuristically reads

$$\partial_t (\text{Anti-Derivative}) = -\frac{1}{n} \sum_{\beta=1}^n (\text{Derivative}) (f_\beta(t) - y_\beta), \quad (3.8)$$

where ‘Anti-Derivative’ refers to the term we take derivative over. For instance, while the dynamics of  $\mathbf{a}_t$  is written into

$$\partial_t \mathbf{a}_t = -\frac{1}{n} \sum_{\beta=1}^n \frac{1}{\sqrt{m}} \sqrt{m} \mathbf{x}_\beta^{[L]} (f_\beta(t) - y_\beta), \quad (3.9)$$

we refer to  $\mathbf{a}_t$  as ‘Anti-Derivative’, and  $\frac{1}{\sqrt{m}}\sqrt{m}\mathbf{x}_\beta^{[L]}$  as its ‘Derivative’. For simplicity, we symbolize the dynamics (3.9) as  $\mathbf{a}_t \rightarrow \frac{1}{\sqrt{m}}\sqrt{m}\mathbf{x}_\beta^{[L]}$ . Similarly, for the dynamics of  $\mathbf{x}_\alpha^{[l]}$ ,

$$\begin{aligned} \sqrt{m}\mathbf{x}_\alpha^{[1]} &\rightarrow \frac{c_\sigma}{\sqrt{m}}\text{diag}\left(\boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\alpha)\boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta)\left(\mathbf{E}_{t,\beta}^{[2:L]}\right)^\top \mathbf{a}_t\right)\mathbf{1}\langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle, \\ \text{for } 2 \leq l \leq L, \\ \sqrt{m}\mathbf{x}_\alpha^{[l]} &\rightarrow \frac{c_\sigma}{\sqrt{m}}\text{diag}\left(\mathbf{E}_{t,\alpha}^{[2:l]}\boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\alpha)\boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta)\left(\mathbf{E}_{t,\beta}^{[2:L]}\right)^\top \mathbf{a}_t\right)\mathbf{1}\langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle \\ &\quad + \sum_{k=2}^l \frac{c_{\text{res}}^2}{L^2\sqrt{m}}\text{diag}\left(\mathbf{E}_{t,\alpha}^{[(k+1):l]}\right. \\ &\quad \left.\boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\alpha)\boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta)\left(\mathbf{E}_{t,\beta}^{[(k+1):L]}\right)^\top \mathbf{a}_t\right)\mathbf{1}\langle \mathbf{x}_\alpha^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle, \end{aligned} \tag{3.10}$$

of  $\mathbf{W}_t^{[l]}$ ,

$$\begin{aligned} \mathbf{W}_t^{[L]} &\rightarrow \frac{c_{\text{res}}}{L\sqrt{m}}\text{diag}\left(\boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_\beta)\mathbf{a}_t\right)\mathbf{1} \otimes (\mathbf{x}_\beta^{[L-1]})^\top, \\ \text{for } 2 \leq l \leq L-1, \\ \mathbf{W}_t^{[l]} &\rightarrow \frac{c_{\text{res}}}{L\sqrt{m}}\text{diag}\left(\boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta)\left(\mathbf{E}_{t,\beta}^{[(l+1):L]}\right)^\top \mathbf{a}_t\right)\mathbf{1} \otimes (\mathbf{x}_\beta^{[l-1]})^\top, \\ \mathbf{W}_t^{[1]} &\rightarrow \sqrt{\frac{c_\sigma}{m}}\text{diag}\left(\boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta)\left(\mathbf{E}_{t,\beta}^{[2:L]}\right)^\top \mathbf{a}_t\right)\mathbf{1} \otimes (\mathbf{x}_\beta)^\top, \end{aligned} \tag{3.11}$$

of  $\sigma_{[l]}^{(1)}(\mathbf{x}_\alpha)$ , where  $r \geq 1$ ,

$$\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, \\
&\text{for } 2 \leq l \leq L-1, \\
\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]} \right. \\
&\quad \left. \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) \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]} \right. \\
&\quad \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) \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle,
\end{aligned} \tag{3.12}$$

and finally of  $\mathbf{E}_{t,\alpha}^{[l]}$

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

for  $2 \leq l \leq L-1$ ,

$$\begin{aligned}
\mathbf{E}_{t,\alpha}^{[l+1]} &\rightarrow \frac{c_{\text{res}}^2}{L^2\sqrt{m}} \text{diag} \left( \boldsymbol{\sigma}_{[l+1]}^{(1)}(\mathbf{x}_\alpha) \boldsymbol{\sigma}_{[l+1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+2):L]} \right)^\top \mathbf{a}_t \right) \mathbf{1} \otimes \left( \frac{\sqrt{m} \mathbf{x}_\beta^{[l]}}{m} \right)^\top \\
&+ \frac{c_{\text{res}}^2}{L^2\sqrt{m}} \boldsymbol{\sigma}_{[l+1]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \boldsymbol{\sigma}_{[l+1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(l+2):L]} \right)^\top \mathbf{a}_t \right) \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \langle \mathbf{x}_\alpha^{[l]}, \mathbf{x}_\beta^{[l]} \rangle \\
&+ \sum_{k=2}^l \frac{c_{\text{res}}^3}{L^3\sqrt{m}} \boldsymbol{\sigma}_{[l+1]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \mathbf{E}_{t,\alpha}^{[(k+1):l]} \right. \\
&\quad \left. \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\alpha) \boldsymbol{\sigma}_{[k]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[(k+1):L]} \right)^\top \mathbf{a}_t \right) \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \langle \mathbf{x}_\alpha^{[k-1]}, \mathbf{x}_\beta^{[k-1]} \rangle \\
&+ \frac{c_\sigma}{\sqrt{m}} \frac{c_{\text{res}}}{L} \boldsymbol{\sigma}_{[l+1]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \mathbf{E}_{t,\alpha}^{[2:l]} \right. \\
&\quad \left. \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\alpha) \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta) \left( \mathbf{E}_{t,\beta}^{[2:L]} \right)^\top \mathbf{a}_t \right) \frac{\mathbf{W}_t^{[l+1]}}{\sqrt{m}} \langle \mathbf{x}_\alpha, \mathbf{x}_\beta \rangle.
\end{aligned} \tag{3.13}$$

Altogether Equation (3.10), Equation (3.11), Equation (3.12) and Equation (3.13) are termed the *Replacement Rules*. We instantly obtain the derivative for  $f_\alpha(t)$ , from Equation (1.2), we



have that  $f_\alpha(t) = f_{\text{res}}(\mathbf{x}_\alpha, \boldsymbol{\theta}_t) = \mathbf{a}_t^\top \mathbf{x}_\alpha^{[L]} = \langle \mathbf{a}_t, \mathbf{x}_\alpha^{[L]} \rangle$ , then by applying one of the replacement rules (Equation (3.10)) on  $f_\alpha(t)$ , we have

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

and we notice that sum of these terms altogether reads  $\mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$ , hence we obtain that

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

Applying the replacement rules once again, the derivative for NTK  $\mathcal{G}_t^{(2)}(\cdot)$  is obtained in the following form

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

with each term in  $\mathcal{G}_t^{(3)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_\beta)$  consisting of the summation of all the terms generated from  $\mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$  by performing the replacement procedure. In order to illustrate the idea, besides the above computational example where  $f_\alpha(t)$  is the main object, another example is given out in the proof of Theorem 3.2.2 in Section 3.6.

By same reasoning, we could obtain higher order kernels inductively by performing all the possible replacements. For kernel  $\mathcal{G}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ , where  $r \geq 2$ , the following Ordinary Differential Equation gives  $\mathcal{G}_t^{(r+1)}(\mathbf{x}_{\alpha_1}, \dots, \mathbf{x}_{\alpha_r}; \mathbf{x}_\beta)$ :

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

and it finishes the proof for Equation (3.1) and Equation (3.2) in Theorem 3.2.1.

Furthermore, in order to describe the members appearing in  $\mathcal{G}_t^{(p)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_p})$ ,  $p \geq 3$  more systematically, some notations shall be introduced.

### 3.3.2 Hierarchical Sets of Kernel Expressions

Firstly, we denote  $\mathbb{A}_0$  as the first set of expressions, which corresponds to the terms in  $\mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$ . We define  $\mathbb{A}_0$  as:

$$\mathbb{A}_0 \triangleq \{e_s e_{s-1} \dots e_1 e_0 : 0 \leq s \leq 4L\}, \quad (3.14)$$

where  $e_j$  is chosen following

$$\begin{aligned} e_0 &\in \left\{ \mathbf{a}_t, \{\sqrt{m} \mathbf{x}_\beta^{[0]}, \sqrt{m} \mathbf{x}_\beta^{[1]}, \sqrt{m} \mathbf{x}_\beta^{[2]}, \dots, \sqrt{m} \mathbf{x}_\beta^{[L]}\}_{1 \leq \beta \leq n} \right\}, \\ \text{for } 1 \leq j \leq s, & \\ e_j &\in \left\{ \left\{ \mathbf{E}_{t,\beta}^{[2]}, \left( \mathbf{E}_{t,\beta}^{[2]} \right)^\top, \dots, \mathbf{E}_{t,\beta}^{[L]}, \left( \mathbf{E}_{t,\beta}^{[L]} \right)^\top \right\}_{1 \leq \beta \leq n}, \left\{ \boldsymbol{\sigma}_{[1]}^{(1)}(\mathbf{x}_\beta), \dots, \boldsymbol{\sigma}_{[L]}^{(1)}(\mathbf{x}_\beta) \right\}_{1 \leq \beta \leq n} \right\}. \end{aligned} \quad (3.15)$$

From Equation (3.4), each term in  $\mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$  reads

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

with vectors  $\mathbf{u}_1(t), \mathbf{u}_2(t), \mathbf{u}_3(t), \mathbf{u}_4(t) \in \mathbb{A}_0$ .

We remark that compared with Huang and Yau [3], their  $\mathbf{e}_j$  ( $1 \leq j \leq s$ ) is chosen differently. The counterpart in [3] is selected 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 change arises from the difference in network structures, and it will be shown further that skip-connection matrices  $\mathbf{E}_{t,\beta}^{[l]}$  possess more stability than  $\mathbf{W}_t^{[l]}/\sqrt{m}$ .

Secondly, we shall investigate the set of expressions in higher orders. Given constructions of  $\mathbb{A}_0, \mathbb{A}_1, \dots, \mathbb{A}_r$ , we denote  $\mathbb{A}_{r+1}$  as the set of expressions in the following form:

$$\mathbb{A}_{r+1} \triangleq \{ \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_1 \mathbf{e}_0 : 0 \leq s \leq 4L \}, \quad (3.16)$$

where  $\mathbf{e}_0$  is chosen from

$$\mathbf{e}_0 \in \{ \mathbf{a}_t, \mathbf{1}, \{ \sqrt{m} \mathbf{x}_\beta^{[0]}, \sqrt{m} \mathbf{x}_\beta^{[1]}, \sqrt{m} \mathbf{x}_\beta^{[2]}, \dots, \sqrt{m} \mathbf{x}_\beta^{[L]} \}_{1 \leq \beta \leq n} \}, \quad (3.17)$$

while for  $1 \leq j \leq s$ , each  $\mathbf{e}_j$  is chosen from one of the three following sets

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

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

with  $2 \leq l \leq L$ ,  $1 \leq \beta \leq n$ ,  $1 \leq u \leq r$ ,

$$\left\{ \boldsymbol{\sigma}_{[l]}^{(u+1)}(\mathbf{x}_\beta) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \dots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \left( \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_{u+1}}, \right. \quad (3.18)$$

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

$$\left. \mathbf{g}_1, \mathbf{g}_2 \dots \mathbf{g}_u \in \mathbb{A}_0 \cup \mathbb{A}_1 \cup \dots \cup \mathbb{A}_r, \quad Q_1, Q_2 \dots Q_{u+1} \in \{0, 1\} \right\}.$$

We use Proposition 3.3.1 to shed light on the structures of the elements in  $\mathbb{A}_r$ , and consequently on the structures of each term in kernel  $\mathcal{G}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ .

**Proposition 3.3.1.** *For any vector  $\mathbf{u}(t) \in \mathbb{A}_r$ , the new vector obtained from  $\mathbf{u}(t)$  by applying the replacement rules is sum of the terms with the following forms*

$$\begin{aligned}
(a). & \frac{C}{\sqrt{m}} \mathbf{u}^*(t) : \mathbf{u}^*(t) \in \mathbb{A}_r, \\
(b). & \frac{C}{\sqrt{m}} \mathbf{u}^*(t) \frac{\langle \mathbf{p}, \mathbf{q} \rangle}{m} : \mathbf{u}^*(t) \in \mathbb{A}_{r+1}, \quad \mathbf{p}, \mathbf{q} \in \mathbb{A}_0, \\
(c). & \frac{C}{\sqrt{m}} \mathbf{u}^*(t) \frac{\langle \mathbf{p}, \mathbf{q} \rangle}{m} : \mathbf{u}^*(t) \in \mathbb{A}_{r-s+1}, \quad \mathbf{p} \in \mathbb{A}_s, \mathbf{q} \in \mathbb{A}_0, \text{ for some } s \geq 1, \\
(d). & \frac{C}{\sqrt{m}} \mathbf{u}^*(t) \frac{\langle \mathbf{p}, \mathbf{q} \rangle}{m} : \mathbf{u}^*(t) \in \mathbb{A}_s, \quad \mathbf{p} \in \mathbb{A}_{r-s+1}, \mathbf{q} \in \mathbb{A}_0, \text{ for some } s \geq 1.
\end{aligned} \tag{3.19}$$

*Proof.* We remark that the constant  $C$  in Equation (3.19) may change from term to term.

(i). Firstly, since  $\mathbf{a}_t$  appears only at the position  $\mathbf{e}_0$ , if  $\mathbf{u}(t) \in \mathbb{A}_r$ , from replacement rules

$$\mathbf{u}(t) = \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_1 \mathbf{a}_t \rightarrow \tilde{\mathbf{u}}(t) = \frac{1}{\sqrt{m}} \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_1 \sqrt{m} \mathbf{x}_\beta^{[L]} = \frac{1}{\sqrt{m}} \mathbf{u}^*(t),$$

then  $\mathbf{u}^*(t) = \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_1 \sqrt{m} \mathbf{x}_\beta^{[L]} \in \mathbb{A}_r$ .

(ii). Similarly,  $\sqrt{m} \mathbf{x}_\alpha^{[l]}$  also appears only at  $\mathbf{e}_0$ , then if  $\mathbf{u}(t) \in \mathbb{A}_r$ , from replacement rules

$$\begin{aligned}
\mathbf{u}(t) &= \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_1 \sqrt{m} \mathbf{x}_\alpha^{[l]} \rightarrow \tilde{\mathbf{u}}(t), \\
\tilde{\mathbf{u}}(t) &= \sum_k \frac{C}{\sqrt{m}} \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_1 \text{diag}(\mathbf{f}_k) \mathbf{1} \frac{\langle \sqrt{m} \mathbf{x}_\alpha^{[k]}, \sqrt{m} \mathbf{x}_\beta^{[k]} \rangle}{m} \\
&= \sum_k \frac{C}{\sqrt{m}} \mathbf{u}_k^*(t) \frac{\langle \sqrt{m} \mathbf{x}_\alpha^{[k]}, \sqrt{m} \mathbf{x}_\beta^{[k]} \rangle}{m},
\end{aligned}$$

with  $\mathbf{f}_k \in \mathbb{A}_0$ , then  $\mathbf{u}_k^*(t) \in \mathbb{A}_{r+1}$ .

(iii). Moreover, for  $\sigma_{[l]}^{(u)}(\mathbf{x}_\alpha)$ , which only appears at the starting or middle position. For  $u = 1$ ,  $\sigma_{[l]}^{(1)}(\mathbf{x}_\alpha)$  has no diag operations accompanied with it, and any vector  $\mathbf{u}(t) \in \mathbb{A}_r$  could contain  $\sigma_{[l]}^{(1)}(\mathbf{x}_\alpha)$  for  $r \geq 0$ . From replacement rules

$$\begin{aligned}\mathbf{u}(t) &= \mathbf{e}_s \dots \mathbf{e}_{j+1} \sigma_{[l]}^{(1)}(\mathbf{x}_\alpha) \mathbf{e}_{j-1} \dots \mathbf{e}_0 \rightarrow \tilde{\mathbf{u}}(t), \\ \tilde{\mathbf{u}}(t) &= \frac{C}{\sqrt{m}} \mathbf{e}_s \dots \mathbf{e}_{j+1} \sigma_{[l]}^{(2)}(\mathbf{x}_\alpha) \text{diag}(\mathbf{f}_1) \mathbf{e}_{j-1} \dots \mathbf{e}_0 \frac{\langle \mathbf{p}_1, \mathbf{q}_1 \rangle}{m} \\ &\quad + \sum_k \frac{C}{\sqrt{m}} \mathbf{e}_s \dots \mathbf{e}_{j+1} \sigma_{[l]}^{(2)}(\mathbf{x}_\alpha) \text{diag} \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \mathbf{f}_k \right) \mathbf{e}_{j-1} \dots \mathbf{e}_0 \frac{\langle \mathbf{p}_k, \mathbf{q}_k \rangle}{m} \\ &= \sum_l \frac{C}{\sqrt{m}} \mathbf{u}_l^*(t) \frac{\langle \mathbf{p}_l, \mathbf{q}_l \rangle}{m},\end{aligned}$$

with  $\mathbf{f}_k \in \mathbb{A}_0$ , then  $\mathbf{u}_l^*(t) \in \mathbb{A}_{r+1}$ , and  $\mathbf{p}_l, \mathbf{q}_l \in \mathbb{A}_0$ .

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

$$\begin{aligned}\mathbf{u}(t) &= \mathbf{e}_s \dots \mathbf{e}_{j+1} \mathbf{e}_j \mathbf{e}_{j-1} \dots \mathbf{e}_0 \rightarrow \tilde{\mathbf{u}}(t), \\ \exists \mathbf{e}_j &= \sigma_{[l]}^{(u)}(\mathbf{x}_\alpha) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \dots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_{u-1}} \mathbf{g}_{u-1} \right) \left( \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u}, \\ \text{or } \mathbf{e}_j &= \left( \frac{c_{\text{res}}}{L} \frac{(\mathbf{W}_t^{[l]})^\top}{\sqrt{m}} \right)^{Q_u} \sigma_{[l]}^{(u)}(\mathbf{x}_\alpha) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \dots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_{u-1}} \mathbf{g}_{u-1} \right),\end{aligned}$$

after applying replacement rules on  $e_j \rightarrow e_j^*$ ,

$$\begin{aligned}
e_j^* &= \frac{C}{\sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(u+1)}(\mathbf{x}_\alpha) \dots \left( \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \frac{\langle \mathbf{p}_1, \mathbf{q}_1 \rangle}{m} \\
&+ \sum_k \frac{C}{\sqrt{m}} \boldsymbol{\sigma}_{[l]}^{(u+1)}(\mathbf{x}_\alpha) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_0} \mathbf{f}_k \right) \dots \left( \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \frac{\langle \mathbf{p}_k, \mathbf{q}_k \rangle}{m}, \\
\text{or } e_j^* &= \frac{C}{\sqrt{m}} \left( \frac{c_{\text{res}}}{L} \frac{(\mathbf{W}_t^{[l]})^\top}{\sqrt{m}} \right)^{Q_u} \boldsymbol{\sigma}_{[l]}^{(u+1)}(\mathbf{x}_\alpha) \text{diag}(\mathbf{f}_1) \dots \frac{\langle \mathbf{p}_1, \mathbf{q}_1 \rangle}{m} \\
&+ \sum_k \frac{C}{\sqrt{m}} \left( \frac{c_{\text{res}}}{L} \frac{(\mathbf{W}_t^{[l]})^\top}{\sqrt{m}} \right)^{Q_u} \boldsymbol{\sigma}_{[l]}^{(u+1)}(\mathbf{x}_\alpha) \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_0} \mathbf{f}_k \right) \dots \frac{\langle \mathbf{p}_k, \mathbf{q}_k \rangle}{m},
\end{aligned}$$

hence

$$\tilde{\mathbf{u}}(t) = \sum_l \frac{C}{\sqrt{m}} \mathbf{u}_l^*(t) \frac{\langle \mathbf{p}_l, \mathbf{q}_l \rangle}{m},$$

with  $\mathbf{f}_k \in \mathbb{A}_0$ , then  $\mathbf{u}_l^*(t) \in \mathbb{A}_{r+1}$ , and  $\mathbf{p}_l, \mathbf{q}_l \in \mathbb{A}_0$ .

(iv). Since  $\frac{\mathbf{W}_t^{[l]}}{\sqrt{m}}$  only appears at the starting or the middle position, if  $\mathbf{u}(t) \in \mathbb{A}_r$ , from replacement rules

$$\begin{aligned}
\mathbf{u}(t) &= e_s e_{s-1} \dots e_{j+1} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} e_{j-1} \dots e_1 e_0 \rightarrow \tilde{\mathbf{u}}(t), \\
\tilde{\mathbf{u}}(t) &= \frac{C}{m} e_s e_{s-1} \dots e_{j+1} \text{diag}(\mathbf{g}) \mathbf{1} \otimes (\mathbf{x}_\beta^{[l-1]})^\top e_{j-1} \dots e_1 e_0 \\
&= \frac{C}{\sqrt{m}} e_s e_{s-1} \dots e_{j+1} \text{diag}(\mathbf{g}) \mathbf{1} \frac{\langle e_{j-1} \dots e_1 e_0, \sqrt{m} \mathbf{x}_\beta^{[l-1]} \rangle}{m} \\
&= \frac{C}{\sqrt{m}} \mathbf{u}^*(t) \frac{\langle \mathbf{p}, \mathbf{q} \rangle}{m},
\end{aligned}$$

with  $\mathbf{u}^*(t) \in \mathbb{A}_{r-s+1}$ , and  $\mathbf{p} \in \mathbb{A}_s$ ,  $\mathbf{q} \in \mathbb{A}_0$ , for some  $s \geq 1$ .

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

$$\begin{aligned}\mathbf{u}(t) &= \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_{j+1} \frac{(\mathbf{W}_t^{[l]})^\top}{\sqrt{m}} \mathbf{e}_{j-1} \dots \mathbf{e}_1 \mathbf{e}_0 \rightarrow \tilde{\mathbf{u}}(t) \\ \tilde{\mathbf{u}}(t) &= \frac{C}{m} \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_{j+1} \mathbf{x}_\beta^{[l-1]} \otimes \mathbf{1}^\top \text{diag}(\mathbf{g}) \mathbf{e}_{j-1} \dots \mathbf{e}_1 \mathbf{e}_0 \\ &= \frac{C}{\sqrt{m}} \mathbf{e}_s \mathbf{e}_{s-1} \dots \mathbf{e}_{j+1} \sqrt{m} \mathbf{x}_\beta^{[l-1]} \frac{\langle \text{diag}(\mathbf{g}) \mathbf{e}_{j-1} \dots \mathbf{e}_1 \mathbf{e}_0, \mathbf{1} \rangle}{m} \\ &= \frac{C}{\sqrt{m}} \mathbf{u}^*(t) \frac{\langle \mathbf{p}, \mathbf{q} \rangle}{m},\end{aligned}$$

with  $\mathbf{u}^*(t) \in \mathbb{A}_{r-s}$ , and  $\mathbf{p} \in \mathbb{A}_{r-s+1}$ ,  $\mathbf{q} \in \mathbb{A}_0$ , for some  $s \geq 1$ .

(v). Finally for  $\mathbf{E}_{t,\alpha}^{[l]}$ , whose scenario is the situations combined with  $\frac{\mathbf{w}_t^{[l]}}{\sqrt{m}}$  and  $\boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\alpha)$ , we shall skip the analysis and conclude the proof.  $\square$

From the discussion above, we apply Proposition 3.3.1 sequentially to  $\mathcal{G}_t^{(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2})$  for  $(r-1)$  many times to obtain kernel  $\mathcal{G}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ , whose individual components takes the form

$$\frac{1}{m^{r/2-1}} \prod_{j=1}^s \frac{\langle \mathbf{u}_{2j-1}(t), \mathbf{u}_{2j}(t) \rangle}{m}, \quad 1 \leq s \leq r, \quad \mathbf{u}_i(t) \in \mathbb{A}_0 \cup \mathbb{A}_1 \cup \dots \cup \mathbb{A}_{r-2}. \quad (3.20)$$

We remark that Equation (3.20) brings partial proof to Equation (3.3).

### 3.4 Key Technique Number Two: Apriori Estimates

Huang and Yau obtained Equation (3.20) in [3, Equation (3.8)], and they use the tensor program [38] to estimate the initial value of the kernel  $\mathcal{G}_0^{(r)}(\cdot)$ . They showed that for each vector  $\mathbf{u}_j(t)$  at  $t = 0$ , it is a linear combination of projections of independent Gaussian vectors. Hence, if we consider such quantity

$$\eta(t) = \{\|\mathbf{u}(t)\|_\infty : \mathbf{u}(t) \in \mathbb{A}_0 \cup \mathbb{A}_1 \dots \cup \mathbb{A}_r\}, \quad (3.21)$$

then with high probability,

$$\eta(0) \lesssim (\ln m)^C \quad (3.22)$$

holds, since any  $\mathbf{u}(0)$  is a linear combination of projections of independent Gaussian vectors.

For  $t > 0$ , Huang and Yau derived a self-consistent Ordinary Differential Inequality for  $\eta(t)$ ,

$$\partial_t^{(p+1)} \eta(t) \lesssim \frac{\eta(t)^{2p}}{m^{p/2}}, \quad (3.23)$$

$$\eta(0) \lesssim (\ln m)^C, \quad (3.24)$$

for some  $p \geq 2$ . Then for time  $0 \leq t \leq m^{\frac{p}{2p+1}} / (\ln m)^{C^*}$ ,

$$\eta(t) \lesssim (\ln m)^C.$$

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 Equation (3.22). Then we derive a Gronwall-type inequality for  $\eta(t)$ ,

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

Hence it follows that for time  $0 \leq t \leq \sqrt{m} / (\ln m)^{C^*}$ ,

$$\eta(t) \lesssim (\ln m)^C$$

holds, which brings the complete proof to Equation (3.3).

### 3.4.1 Apriori $L^2$ Bounds for Expressions in $\mathbb{A}_0$

We begin our proof with several lemmas.

**Lemma 3.4.1.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , for  $t \geq 0$ ,*

$$R_S(\boldsymbol{\theta}_t) \leq R_S(\boldsymbol{\theta}_0) \sim \mathcal{O}(1). \quad (3.25)$$



*Proof.* We get inequality (3.25) by non-negative definiteness of kernel  $\mathcal{G}_t^{(2)}(\cdot)$ . From Equation (3.1), we obtain that

$$\partial_t \sum_{\alpha=1}^n \|f_\alpha(t) - y_\alpha\|_2^2 = -\frac{2}{n} \sum_{\alpha,\beta=1}^n \mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta)(f_\alpha(t) - y_\alpha)(f_\beta(t) - y_\beta) \leq 0, \quad (3.26)$$

hence

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

which finishes the proof of the lemma.  $\square$

Our next lemma is mainly on the spectral property of random matrices, which was given out in [1, Lemma G.2.], also consequence of results in [40].

**Lemma 3.4.2.** *Given  $\mathbf{W} \in \mathbb{R}^{m \times m}$ , with entry  $W_{i,j} \sim \mathcal{N}(0, 1)$ , then with probability at least  $1 - \exp\left(-\frac{(c_{w,0}-2)^2 m}{2}\right)$ , the following holds*

$$\|\mathbf{W}\|_{2 \rightarrow 2} \leq c'_{w,0} \sqrt{m}, \quad (3.27)$$

where the constant  $c'_{w,0} > 2$ .

Our next lemma is on the tail bound of the chi-square distribution, whose proof can be found in [41].

**Lemma 3.4.3.** *If  $Z \sim \chi^2(m)$ , then we have a tail bound*

$$\mathbb{P}\left(Z \geq m + 2\sqrt{mx} + 2x\right) \leq e^{-x}. \quad (3.28)$$

Our next proposition is similar to Proposition B.1. in [3].

**Proposition 3.4.1.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , define  $\xi(t)$*

$$\xi(t) = \sup_{0 \leq t^* \leq t} \max \left\{ 1, \frac{1}{\sqrt{m}} \left\{ \left\| \mathbf{W}_{t^*}^{[2]} \right\|_{2 \rightarrow 2}, \left\| \left( \mathbf{W}_{t^*}^{[2]} \right)^\top \right\|_{2 \rightarrow 2}, \dots \right. \right. \\ \left. \left. \dots, \left\| \mathbf{W}_{t^*}^{[L]} \right\|_{2 \rightarrow 2}, \left\| \left( \mathbf{W}_{t^*}^{[L]} \right)^\top \right\|_{2 \rightarrow 2}, \left\| \mathbf{a}_{t^*} \right\|_2 \right\} \right\}, \quad (3.29)$$

then with high probability w.r.t the random initialization, for  $t \lesssim \sqrt{m}$

$$\xi(t) \leq c_{w,t}, \quad (3.30)$$

where the constant  $c_{w,t} > 2$  is independent of the depth  $L$ . Moreover for  $t \lesssim \sqrt{m}$ ,  $c_{w,t}$  has a uniform upper bound in  $t$ , i.e.,

$$c_{w,t} \leq \bar{c}, \quad (3.31)$$

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

*Proof.* (i). Firstly, set  $Z = \|\mathbf{a}\|_2^2$  in Lemma 3.4.3. Then if we write  $2tm = m + (2t - 1)m$ , with  $x = \frac{mt}{10}$ ,

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

and for  $t \geq 1$ , we have  $2t - 1 \geq 2 \left( \sqrt{t/10} + t/10 \right)$ . 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, where  $c_{w,0}$  is the constant in Lemma 3.4.2. Hence, combined with Lemma 3.4.2 for  $t = 0$ ,  $\xi(0) \leq \max \{1, c'_{w,0}\} = c_{w,0}$ .

(ii). Secondly, we derive the upper bound for  $\partial_t \xi(t)$ . In order for that, the  $L^2$  bound on each output layer shall be estimated. For  $l = 1$ ,

$$\begin{aligned} \|\mathbf{x}^{[1]}\|_2 &= \sqrt{\frac{c_\sigma}{m}} \|\sigma(\mathbf{W}_t^{[1]} \mathbf{x})\|_2 \leq \sqrt{c_\sigma} \left( |\sigma(0)| + \frac{1}{\sqrt{m}} \|\mathbf{W}_t^{[1]} \mathbf{x}\|_2 \right) \\ &\leq \sqrt{c_\sigma} (1 + \xi(t) \|\mathbf{x}\|_2) \leq C \xi(t), \end{aligned} \quad (3.32)$$

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

$$\begin{aligned} \|\mathbf{x}^{[l]}\|_2 &\leq \|\mathbf{x}^{[l-1]}\|_2 + \frac{c_{\text{res}}}{L\sqrt{m}} \|\sigma(\mathbf{W}_t^{[l]} \mathbf{x}^{[l-1]})\|_2 \\ &\leq \|\mathbf{x}^{[l-1]}\|_2 + \frac{c_{\text{res}}}{L} \left( |\sigma(0)| + \xi(t) \|\mathbf{x}^{[l-1]}\|_2 \right) \\ &\leq \|\mathbf{x}^{[l-1]}\|_2 + \frac{c_{\text{res}}}{L} (1 + \xi(t) \|\mathbf{x}^{[l-1]}\|_2) \\ &\leq \left( 1 + \frac{2c_{\text{res}}}{L} \xi(t) \right) \|\mathbf{x}^{[l-1]}\|_2. \end{aligned} \quad (3.33)$$

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

$$\|\mathbf{x}^{[l]}\|_2 \leq C \left(1 + \frac{2c_{\text{res}}}{L}\xi(t)\right)^{l-1} \xi(t). \quad (3.34)$$

Based on Equation (2.6), combined with Lemma 3.4.1

$$\begin{aligned} \partial_t \|\mathbf{W}_t^{[l]}\|_{2 \rightarrow 2} &\leq \frac{1}{n} \sum_{\beta=1}^n \frac{C}{\sqrt{m}} \|\boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) (\mathbf{E}_{t,\beta}^{[(l+1):L]})^\top \mathbf{a}_t\|_2 \|\mathbf{x}_\beta^{[l-1]}\|_2 |f_\beta(t) - y_\beta| \\ &\leq \frac{1}{n} \sum_{\beta=1}^n C \left(1 + \frac{c_{\text{res}}}{L}\xi(t)\right)^{L-l} \xi(t) \left(1 + \frac{2c_{\text{res}}}{L}\xi(t)\right)^{l-1} \xi(t) |f_\beta(t) - y_\beta| \\ &\leq C \left(1 + \frac{2c_{\text{res}}}{L}\xi(t)\right)^{L-1} \xi(t)^2 \sqrt{\frac{1}{n} \sum_{\beta=1}^n \|f_\beta(t) - y_\beta\|_2^2} \\ &\leq C \left(1 + \frac{2c_{\text{res}}}{L}\xi(t)\right)^{L-1} \xi(t)^2 \sqrt{R_S(\boldsymbol{\theta}_0)} \\ &\leq C \left(1 + \frac{2c_{\text{res}}}{L}\xi(t)\right)^{L-1} \xi(t)^2 \leq C \exp(2c_{\text{res}}\xi(t)) \xi(t)^2, \quad (3.35) \\ \partial_t \|\mathbf{a}_t\|_2 &\leq \frac{1}{n} \sum_{\beta=1}^n \|\mathbf{x}_\beta^{[L]}\|_2 |f_\beta(t) - y_\beta| \\ &\leq C \left(1 + \frac{2c_{\text{res}}}{L}\xi(t)\right)^{L-1} \xi(t) \sqrt{\frac{1}{n} \sum_{\beta=1}^n \|f_\beta(t) - y_\beta\|_2^2} \\ &\leq C \left(1 + \frac{2c_{\text{res}}}{L}\xi(t)\right)^{L-1} \xi(t) \sqrt{R_S(\boldsymbol{\theta}_0)} \\ &\leq C \left(1 + \frac{2c_{\text{res}}}{L}\xi(t)\right)^{L-1} \xi(t) \leq C \exp(2c_{\text{res}}\xi(t)) \xi(t). \end{aligned}$$

Consequently, we have

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

then an integration inequality can be obtained,

$$\int_{\xi(0)}^{\xi(t)} \frac{du}{\exp(2c_{\text{res}}u)u^2} \leq \frac{Ct}{\sqrt{m}}. \quad (3.36)$$

Hence the integration term on the LHS of Equation (3.36) is

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

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

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

the 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  $t \leq c\sqrt{m}$ ,  $c$  small enough, the following holds

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

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

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

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

$$\bar{c} = \frac{c_{w,0} + \sqrt{c_{w,0}^2 + 2c_{w,0}/c_{\text{res}}}}{2}, \tag{3.37}$$

which finishes the proof of Proposition 3.4.1. □

We state the inductive relation (3.34) as a proposition.

**Proposition 3.4.2.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , then with high probability w.r.t the random initialization, for each  $l$ , given time  $t \lesssim \sqrt{m}$ ,*

$$\|\mathbf{x}^{[l]}\|_2 \leq C, \quad (3.38)$$

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

We remark that the constant  $C$  in Proposition 3.4.2 only depends on  $c_{\text{res}}$  and  $c_{\sigma}$ . However, for fully-connected networks, Equation (3.38) in Proposition 3.4.2 become

$$\|\mathbf{x}^{[l]}\|_2 \leq C 2^l. \quad (3.39)$$

Hence the Euclidean norm of each output layer increases exponentially layer by layer for fully-connected networks, revealing that ResNet possesses more stability.

Finally, we make Apriori estimates on the Euclidean norm of arbitrary vector  $\mathbf{u}(t) \in \mathbb{A}_0$ .

**Proposition 3.4.3.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , with high probability w.r.t the random initialization, uniformly for any vector  $\mathbf{u}(t) \in \mathbb{A}_0$ , given time  $t \lesssim \sqrt{m}$ ,*

$$\|\mathbf{u}(t)\|_2 \leq c\sqrt{m}, \quad (3.40)$$

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

*Proof.* We shall start our analysis on the whole expressions in set  $\mathbb{A}_0$ . For any vector  $\mathbf{u}(t) \in \mathbb{A}_0$ , it can be written into

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

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

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \{\sqrt{m}\mathbf{x}_{\beta}^{[0]}, \sqrt{m}\mathbf{x}_{\beta}^{[1]}, \sqrt{m}\mathbf{x}_{\beta}^{[2]}, \dots, \sqrt{m}\mathbf{x}_{\beta}^{[L]}\}_{1 \leq \beta \leq n} \right\}.$$

- (a). If  $e_0 = \mathbf{a}_t$ , then by Lemma 3.4.3, for  $t \lesssim \sqrt{m}$ ,

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

- (b). If  $e_0 = \sqrt{m}\mathbf{x}_\beta^{[l]}$ , then based on Proposition 3.4.2, for  $t \lesssim \sqrt{m}$ ,

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

Now we proceed to  $e_j$ ,  $j \geq 1$ .

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

$$\begin{aligned} \|\mathbf{u}(t)\|_2 &= \|e_s e_{s-1} \dots e_1 e_0\|_2 \\ &= \|e_s\|_{2 \rightarrow 2} \|e_{s-1}\|_{2 \rightarrow 2} \dots \|e_1\|_{2 \rightarrow 2} \|e_0\|_2. \end{aligned}$$

Since  $\|\boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta)\|_{2 \rightarrow 2} \leq 1$ , thus for all  $j \geq 1$  with  $e_j = \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta)$

$$\|\mathbf{u}(t)\|_2 \leq c\sqrt{m}.$$

- (ii). If  $e_j = \mathbf{E}_{t,\beta}^{[l]}$  or  $e_j = (\mathbf{E}_{t,\beta}^{[l]})^\top$ , then based on Proposition 3.4.1

$$\|\mathbf{u}(t)\|_2 = \|e_s\|_{2 \rightarrow 2} \|e_{s-1}\|_{2 \rightarrow 2} \dots \|e_1\|_{2 \rightarrow 2} \|e_0\|_2.$$

Since

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

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

$$\|\mathbf{u}(t)\|_2 \leq \left(1 + \frac{c_{\text{res}} c_{w,t}}{L}\right)^s \|e_0\|_2, \quad (3.41)$$

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

$$\begin{aligned}\|\mathbf{u}(t)\|_2 &\leq \left(1 + \frac{c_{\text{res}}c_{w,t}}{L}\right)^{4L} \|\mathbf{e}_0\|_2 \\ &\leq c \exp(4c_{\text{res}}c_{w,t})\sqrt{m} \leq c\sqrt{m}.\end{aligned}$$

Combining these two observations, we finish the proof.  $\square$

To sum up, we define

$$\xi_{\infty,0}(t) = \sup_{0 \leq t^* \leq t} \{\|\mathbf{u}(t^*)\|_2 : \mathbf{u}(t^*) \in \mathbb{A}_0\}. \quad (3.42)$$

Then directly from Proposition 3.4.3, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

$$\xi_{\infty,0}(t) \leq c\sqrt{m}. \quad (3.43)$$

### 3.4.2 Apriori $L^\infty$ Bounds for Expressions in $\mathbb{A}_0$

In this part, we shall make estimate on such quantity

$$\eta_{\infty,0}(t) = \sup_{0 \leq t^* \leq t} \{\|\mathbf{u}(t^*)\|_\infty : \mathbf{u}(t^*) \in \mathbb{A}_0\}. \quad (3.44)$$

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

**Lemma 3.4.4.** *For any i.i.d. standard normal distribution  $X_1, X_2, \dots, X_m$ , it holds with high probability that  $L^\infty$ -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 large constant  $C > 0$ .

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

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

We shall optimize over  $\lambda$ ,

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

from symmetry of Gaussian variables

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

hence by taking 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

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

We remark that when  $C > 0$  is large enough,  $(\ln m)^C \approx m^\varepsilon$  for some small  $\varepsilon > 0$ .  $\square$

Our next lemma concerns the matrix two to infinity norm.

**Lemma 3.4.5.** *Given  $\mathbf{W} \in \mathbb{R}^{m \times m}$  with entry  $W_{ij} \sim \mathcal{N}(0, 1)$ , then with high probability,*

$$\|\mathbf{W}\|_{2 \rightarrow \infty} = \sup_{\|\mathbf{x}\|_2=1} \|\mathbf{W}\mathbf{x}\|_\infty \leq (\ln m)^C, \quad (3.45)$$

for some large constant  $C > 0$ .

*Proof.* Note that  $\mathbf{W}\mathbf{x}$  shares the same distribution as the Gaussian vector  $\mathbf{X}$  in Lemma 3.4.4, i.e.  $\mathbf{W}\mathbf{x} \sim \mathbf{X}$ . Then apply Lemma 3.4.4 directly, we finish the proof.  $\square$



To evaluate  $\eta_{\infty,0}(t)$ , another lemma shall be stated.

**Lemma 3.4.6.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , for vectors in  $\mathbb{A}_0$ , define*

$$\eta_{q,0}(t) := \sup_{0 \leq t^* \leq t} \left\{ \|\mathbf{u}_q(t^*)\|_\infty : \mathbf{u}_q(t^*) = \mathbf{e}_q \mathbf{e}_{q-1} \dots \mathbf{e}_1 \mathbf{e}_0, \mathbf{u}_q(t^*) \in \mathbb{A}_0 \right\}. \quad (3.46)$$

Moreover, define

$$\omega(t) := \sup_{0 \leq t^* \leq t} \max \left\{ \|\mathbf{W}_{t^*}^{[2]}\|_{2 \rightarrow \infty}, \left\| \left( \mathbf{W}_{t^*}^{[2]} \right)^\top \right\|_{2 \rightarrow \infty}, \dots, \|\mathbf{W}_{t^*}^{[L]}\|_{2 \rightarrow \infty}, \left\| \left( \mathbf{W}_{t^*}^{[L]} \right)^\top \right\|_{2 \rightarrow \infty} \right\},$$

then with high probability w.r.t the random initialization, for  $t \lesssim \sqrt{m}$ ,

$$\eta_{q,0}(t) \leq \eta_{0,0}(t) + c \omega(t) \left( 1 + \frac{c_{\text{res}} c_{w,t}}{L} \right)^q, \quad (3.47)$$

where constants  $c_{w,t} > 2, c > 0$  are independent of the depth  $L$ .

*Proof.* For any vector  $\mathbf{u}_q(t) \in \mathbb{A}_0$  of length  $q$ ,  $0 \leq q \leq 4L$ , it can be written into

$$\mathbf{u}_q(t) = \mathbf{e}_q \mathbf{e}_{q-1} \dots \mathbf{e}_1 \mathbf{e}_0.$$

We shall prove Equation (3.47) by performing induction on  $q$ . Firstly, for  $q = 0$ , Equation (3.47) is trivial. For  $q \geq 1$ , we shall investigate on the terms  $\mathbf{e}_j$  in the expression  $\mathbf{u}_q(t)$ .

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

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

since  $\|\boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta)\|_{\infty \rightarrow \infty} \leq 1$ , we have

$$\|\mathbf{u}_q(t)\|_\infty \leq \|\mathbf{e}_0\|_\infty,$$

then

$$\eta_{q,0}(t) \leq \eta_{0,0}(t).$$

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

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

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

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

Based on Proposition 3.4.1

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

then

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

inductively

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

where properties of geometric sums are used. By taking supreme, we have

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

□

These lemmas above enables us to state a proposition on the quantity  $\eta_{\infty,0}(0)$ .

**Proposition 3.4.4.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , then with high probability w.r.t the random initialization,*

$$\eta_{\infty,0}(0) \leq c(\ln m)^C, \quad (3.48)$$

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

*Proof.* As always, any vector  $\mathbf{u}(t) \in \mathbb{A}_0$  can be written into

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

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

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \{\sqrt{m}\mathbf{x}_\beta^{[0]}, \sqrt{m}\mathbf{x}_\beta^{[1]}, \sqrt{m}\mathbf{x}_\beta^{[2]}, \dots, \sqrt{m}\mathbf{x}_\beta^{[L]}\}_{1 \leq \beta \leq n} \right\}.$$

- (a). If  $\mathbf{e}_0 = \mathbf{a}_t$ , then by Lemma 3.4.4,

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

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

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

Moreover, for  $l \geq 1$ , from Proposition 3.4.2,

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

inductively,

$$\|\sqrt{m}\mathbf{x}_\beta^{[l]}\|_\infty \leq c \left(1 + \frac{l}{L}\right) (\ln m)^C \leq c(\ln m)^C, \tag{3.49}$$

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

Consequently,

$$\eta_{0,0}(0) \leq c(\ln m)^C. \tag{3.50}$$

Directly from Lemma 3.4.6

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

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

□

Our next proposition is on the estimate of  $\eta_{\infty,0}(t)$  for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ , it is one of the most important propositions in this thesis.

**Proposition 3.4.5.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , then with high probability w.r.t the random initialization, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,*

$$\eta_{\infty,0}(t) \leq c(\ln m)^C, \quad (3.51)$$

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

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

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \left\{ \sqrt{m} \mathbf{x}_\beta^{[0]}, \sqrt{m} \mathbf{x}_\beta^{[1]}, \sqrt{m} \mathbf{x}_\beta^{[2]}, \dots, \sqrt{m} \mathbf{x}_\beta^{[L]} \right\}_{1 \leq \beta \leq n} \right\}.$$

We observe that from the replacement rules given in Section 3.3.1

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

then for  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ , from Proposition 3.4.2

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

by taking supreme on time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

$$\eta_{0,0}(t) \leq c(\ln m)^C + \frac{C}{\sqrt{m}} \int_0^t \eta_{\infty,0}(s) \, ds. \quad (3.52)$$

Secondly, for the estimate of  $\omega(t)$ , from the replacement rules,

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

then from Proposition 3.4.3,

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

hence by taking supreme on time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

$$\omega(t) \leq (\ln m)^C + \frac{C}{\sqrt{m}} \int_0^t \eta_{\infty,0}(s) \, ds. \quad (3.53)$$

Directly from Lemma 3.4.6

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

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

$$\eta_{\infty,0}(t) \leq c(\ln m)^C + \frac{C}{\sqrt{m}} \int_0^t \eta_{\infty,0}(s) \, ds. \quad (3.54)$$

We notice that Equation (3.54) gives us a Gronwall-type inequality, hence

$$\eta_{\infty,0}(t) \leq c(\ln m)^C \exp \left( \frac{Ct}{\sqrt{m}} \right). \quad (3.55)$$

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

$$\eta_{\infty,0}(t) \leq c(\ln m)^C. \quad (3.56)$$

□

### 3.4.3 Apriori $L^2$ and $L^\infty$ Bounds for Expressions in $\mathbb{A}_r$ , $r \geq 1$

In this part, we shall make estimates on the norms  $\|\cdot\|_\infty$  and  $\|\cdot\|_2$  of vectors belonging to higher order sets  $\mathbb{A}_r$ ,  $r \geq 1$ . Naturally, several quantities for vectors in  $\mathbb{A}_r$  with length  $q$  shall be defined

$$\xi_{q,r}(t) := \sup_{0 \leq t^* \leq t} \left\{ \|\mathbf{u}_q(t^*)\|_2 : \mathbf{u}_q(t^*) = \mathbf{e}_q \mathbf{e}_{q-1} \dots \mathbf{e}_1 \mathbf{e}_0, \mathbf{u}_q(t^*) \in \mathbb{A}_r \right\}, \quad (3.57)$$

Specifically, from Proposition 3.4.1 and Proposition 3.4.3,

$$\xi_{q,0}(t) \leq c \left( 1 + \frac{c_{\text{res}} c_{w,t}}{L} \right)^q \sqrt{m}. \quad (3.58)$$

Moreover, we define

$$\xi_{\infty,r}(t) = \sup_{0 \leq q \leq 4L} \{ \xi_{q,r}(t) \}, \quad (3.59)$$

and we remark that  $\xi_{\infty,0}(t)$  is consistent with  $\xi_{\infty,r}(t)$  for  $r = 0$ .

Similarly we define

$$\eta_{q,r}(t) := \sup_{0 \leq t^* \leq t} \left\{ \|\mathbf{u}_q(t^*)\|_\infty : \mathbf{u}_q(t^*) = \mathbf{e}_q \mathbf{e}_{q-1} \dots \mathbf{e}_1 \mathbf{e}_0, \mathbf{u}_q(t^*) \in \mathbb{A}_r \right\}, \quad (3.60)$$

and

$$\eta_{\infty,r}(t) = \sup_{0 \leq q \leq 4L} \{ \eta_{q,r}(t) \}. \quad (3.61)$$

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

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

Since  $\mathbf{e}_0$  is chosen from

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \mathbf{1}, \{ \sqrt{m} \mathbf{x}_\beta^{[0]}, \sqrt{m} \mathbf{x}_\beta^{[1]}, \sqrt{m} \mathbf{x}_\beta^{[2]}, \dots, \sqrt{m} \mathbf{x}_\beta^{[L]} \}_{1 \leq \beta \leq n} \right\}.$$

Since  $\mathbf{e}_0$  is chosen from

$$\mathbf{e}_0 \in \left\{ \mathbf{a}_t, \mathbf{1}, \{ \sqrt{m} \mathbf{x}_\beta^{[0]}, \sqrt{m} \mathbf{x}_\beta^{[1]}, \sqrt{m} \mathbf{x}_\beta^{[2]}, \dots, \sqrt{m} \mathbf{x}_\beta^{[L]} \}_{1 \leq \beta \leq n} \right\}.$$

and  $\|\mathbf{1}\|_\infty = 1, \|\mathbf{1}\|_2 = \sqrt{m}$ , then from Proposition 3.4.1 and Proposition 3.4.5, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

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

Next we proceed to  $\mathbf{e}_j$ ,  $j \geq 1$ . For each  $\mathbf{e}_j$ ,

- (i)  $\mathbf{e}_j = \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta)$ ,  $\mathbf{e}_j = \mathbf{E}_{t,\beta}^{[l]}$  or  $\mathbf{e}_j = \left( \mathbf{E}_{t,\beta}^{[l]} \right)^\top$ ,  $2 \leq l \leq L$ .
- (ii)  $\mathbf{e}_j = \text{diag}(\mathbf{g}), \mathbf{g} \in \mathbb{A}_0 \cup \mathbb{A}_1 \cup \dots \cup \mathbb{A}_{r-1}$ .
- (iii)

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

or

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

We observe that the total number of  $\text{diag}$  operations in  $\mathbf{u}(t) \in \mathbb{A}_r$  is at most  $r$ , and that is how a vector belonging to different hierarchical sets is characterized. Especially if  $\exists \mathbf{e}_j$  belonging to case (iii), there are two scenarios:

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



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

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

or

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

we shall take advantage of the special structure of  $e_j$ . Define a new type of skip-connection matrix

$$\widetilde{\mathbf{E}}_{t,\beta}^{[l,r]} := \left( \mathbf{I}_m + \frac{c_{\text{res}}}{L} \boldsymbol{\sigma}_{[l]}^{(r)}(\mathbf{x}_\beta) \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right) \quad r \geq 2. \quad (3.62)$$

Then  $e_j$  can be written into

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

or

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

To illustrate such relation, if some vector  $\bar{\mathbf{u}}(t)$  contains  $\mathbf{e}_j$  belonging to case (iii), it can be decomposed into

$$\begin{aligned}
\bar{\mathbf{u}}(t) &= \mathbf{e}_s \mathbf{e}_{s-1} \cdots \mathbf{e}_{j+1} \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \boldsymbol{\sigma}_{[l]}^{(u+1)}(\mathbf{x}_\beta) \frac{c_{\text{res}}}{L} \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \mathbf{e}_{j-1} \cdots \mathbf{e}_0 \\
&= \mathbf{e}_s \mathbf{e}_{s-1} \cdots \mathbf{e}_{j+1} \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \widetilde{\mathbf{E}}_{t,\beta}^{[l,u+1]} \mathbf{e}_{j-1} \cdots \mathbf{e}_0 \\
&\quad - \mathbf{e}_s \mathbf{e}_{s-1} \cdots \mathbf{e}_{j+1} \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right) \cdots \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \right)^{Q_u} \mathbf{g}_u \right) \mathbf{e}_{j-1} \cdots \mathbf{e}_0.
\end{aligned}$$

From the analysis above, we are able to characterize an element in set  $\mathbb{A}_r$ . If  $\mathbf{u}(t) \in \mathbb{A}_r$ , then there exists  $\mathbf{e}_{j_1}, \mathbf{e}_{j_2}, \dots, \mathbf{e}_{j_k}$ , such that

$$\begin{aligned}
\mathbf{e}_{j_1} &= \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l_1]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right), \quad \mathbf{g}_1 \in \mathbb{A}_{r_1-1}, \\
\mathbf{e}_{j_2} &= \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l_2]}}{\sqrt{m}} \right)^{Q_2} \mathbf{g}_2 \right), \quad \mathbf{g}_2 \in \mathbb{A}_{r_2-1}, \\
&\vdots \\
\mathbf{e}_{j_k} &= \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l_k]}}{\sqrt{m}} \right)^{Q_k} \mathbf{g}_k \right), \quad \mathbf{g}_k \in \mathbb{A}_{r_k-1},
\end{aligned}$$

with

$$r_1 + r_2 + \cdots + r_k = r, \quad r_1, r_2, \dots, r_k \in \mathbb{N}^+, \quad (3.63)$$

then Equation (3.63) serves as the counting of the number of diag operations contained in  $\mathbf{u}(t)$ , while for other  $\mathbf{e}_j$  ( $j \notin \{j_1, j_2, \dots, j_k, 0\}$ ), chosen from

$$\left\{ \mathbf{E}_{t,\beta}^{[l]}, \left( \mathbf{E}_{t,\beta}^{[l]} \right)^\top : 2 \leq l \leq L \right\}_{1 \leq \beta \leq n}, \quad (3.64)$$

$$\left\{ \boldsymbol{\sigma}_{[l]}^{(1)}(\mathbf{x}_\beta) : 1 \leq l \leq L \right\}_{1 \leq \beta \leq n}, \quad (3.65)$$

$$\left\{ \widetilde{\mathbf{E}}_{t,\beta}^{[l,p]}, \left( \widetilde{\mathbf{E}}_{t,\beta}^{[l,p]} \right)^\top : 2 \leq l \leq L, p \geq 2 \right\}_{1 \leq \beta \leq n}, \quad (3.66)$$

note that the elements in set (3.64) and set (3.66) share the same matrix properties, thanks to the assumptions 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, which is shown in the proof of Proposition 3.4.6. We remark that Proposition 3.4.6 is also one of the most important proposition in this thesis.

**Proposition 3.4.6.** *Given  $\mathcal{X}$  and  $\sigma(\cdot)$ , then with high probability w.r.t the random initialization, for some finite  $r \geq 1$  and time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,*

$$\begin{aligned}\xi_{\infty,r}(t) &\leq c(\ln m)^C \sqrt{m}, \\ \eta_{\infty,r}(t) &\leq c(\ln m)^C,\end{aligned}\tag{3.67}$$

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

*Proof.* We recall the definition of  $\omega(t)$ ,  $\eta_{\infty,0}(t)$  and  $\xi_{\infty,0}(t)$ , then with high probability, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

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

(i). Let us start out induction with  $r = 1$ . For  $\mathbf{u}(t) \in \mathbb{A}_1$ , since there exists only one solution to Equation (3.63), then there is one and it is the only index  $i$ , such that  $\mathbf{e}_i = \text{diag}(\mathbf{g})$ , or  $\mathbf{e}_i = \text{diag}\left(\frac{\mathbf{w}_t^{[l]}}{\sqrt{m}}\mathbf{g}\right)$ , with  $\mathbf{g} \in \mathbb{A}_0$ . Consequently,

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

and for  $q > i$ ,

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

then inductively,

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

By taking supreme on  $q$  and  $i$ ,

$$\xi_{\infty,1}(t) \leq \exp(4c_{\text{res}} c_{w,t}) c(\ln m)^C \xi_{\infty,0}(t) \leq c(\ln m)^C \sqrt{m}. \quad (3.68)$$

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

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

and for  $q > i$ , inductively

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

by taking supreme on  $q$  and  $i$ , combined with Equation (3.68),

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

(iii). In the following we assume that Equation (3.67) holds for  $1, 2, \dots, r-1$  and we will prove it for  $r$ .

If  $\mathbf{u}(t) \in \mathbb{A}_r$ , there exists  $\mathbf{e}_{j_1}, \mathbf{e}_{j_2}, \dots, \mathbf{e}_{j_k}$ , such that

$$\begin{aligned} \mathbf{e}_{j_1} &= \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l_1]}}{\sqrt{m}} \right)^{Q_1} \mathbf{g}_1 \right), \quad \mathbf{g}_1 \in \mathbb{A}_{r_1-1}, \\ \mathbf{e}_{j_2} &= \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l_2]}}{\sqrt{m}} \right)^{Q_2} \mathbf{g}_2 \right), \quad \mathbf{g}_2 \in \mathbb{A}_{r_2-1}, \\ &\vdots \\ \mathbf{e}_{j_k} &= \text{diag} \left( \left( \frac{\mathbf{W}_t^{[l_k]}}{\sqrt{m}} \right)^{Q_k} \mathbf{g}_k \right), \quad \mathbf{g}_k \in \mathbb{A}_{r_k-1}, \end{aligned}$$

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  $\mathbf{e}_i = \text{diag}(\mathbf{g}_1)$ , or  $\mathbf{e}_i = \text{diag}\left(\frac{\mathbf{W}_t^{[l]}}{\sqrt{m}} \mathbf{g}_1\right)$  with  $\mathbf{g}_1 \in \mathbb{A}_{r_1-1}$ , then

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

inductively

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

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

$$\xi_{\infty,r}(t) \leq \exp(4c_{\text{res}} c_{w,t}) c(\ln m)^C \xi_{\infty,r-r_1}(t) \leq c(\ln m)^C \sqrt{m}. \quad (3.69)$$

(iv). For  $\eta_{i,r}(t)$ ,

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

and for  $q > i$ ,

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

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

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

□

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

### 3.5 Proof of Theorem 3.2.1

Since for each term in kernel  $\mathcal{G}_t^{(r)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \dots, \mathbf{x}_{\alpha_r})$ ,  $r \geq 2$ , it takes the form

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

Firstly, from Equation (3.43), for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

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

and for  $r \geq 3$ , from Proposition 3.4.6, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

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

### 3.6 Proof of Theorem 3.2.2

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. We need to investigate the order of  $C(r, L)$  for  $r = 2, 3, 4$  respectively.

Firstly,  $\mathcal{G}_t^{[L+1]}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) = \langle \mathbf{x}_{\alpha_1}^{[L]}, \mathbf{x}_{\alpha_2}^{[L]} \rangle$ , then  $C(2, L) = \mathcal{O}(1)$ , since there is only one term.

Secondly, from the replacement rule, all possible terms generated from  $\mathcal{G}_t^{[L+1]}(\cdot)$  are

$$\begin{aligned} \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 \rightarrow \mathcal{G}_t^{[L+1],(1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\beta}) \\ \mathcal{G}_t^{[L+1],(1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\beta}) &= \frac{c_{\sigma}}{m} \underbrace{\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}_{\text{I}} \\ &+ \sum_{k=2}^L \frac{c_{\text{res}}^2}{L^2 m} \underbrace{\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}_{\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}$$

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

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

In order to analyze the dynamics for  $\mathcal{G}_t^{[L+1],(1)}(\cdot)$ , we need the information of its derivative. Hence, we apply replacement rules once again to  $\mathcal{G}_t^{[L+1],(1)}(\cdot)$  to obtain  $\mathcal{G}_t^{[L+1],(2)}(\cdot)$ , i.e.,

$$\mathcal{G}_t^{[L+1],(1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \rightarrow \mathcal{G}_t^{[L+1],(2)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}, \mathbf{x}_{\beta}).$$

Finally for  $\mathcal{G}_t^{[L+1],(2)}(\cdot)$ , by symmetry, only terms I and II need to be analyzed.

- There are at most  $(2L+2)$  symbols in term I to be replaced, and each replacement operation 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

$$C(4, L) = \mathcal{O} \left( 2 \left( (2L+2)(L+1) + \frac{1}{L^2}(L-1)(2L+2)(L+1) \right) \right) = \mathcal{O} \left( L^2 \right). \quad (3.70)$$

Next, we turn to the proof of Equation (3.5). From Proposition 3.4.6, for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

$$\begin{aligned} \left\| \mathcal{G}_t^{[L+1],(2)}(\cdot) \right\|_{\infty} &\leq \frac{C(4, L)}{m} \left( \frac{\xi_{\infty,2}(t)^2}{m} \right)^4 \\ &\leq \frac{C(4, L)}{m} \left( \frac{(c(\ln m)^C \sqrt{m})^2}{m} \right)^4 \\ &\leq C(4, L) \frac{(\ln m)^C}{m}, \end{aligned}$$



then

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

moreover, for  $1 \leq \alpha_1, \alpha_2, \alpha_3 \leq n$ , and time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

$$\begin{aligned}
\left| \mathcal{G}_t^{[L+1],(1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| &\leq \left| \mathcal{G}_0^{[L+1],(1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| + t \sup_{0 \leq s \leq t} \left| \partial_s \mathcal{G}_s^{[L+1],(1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| \\
&\leq \left\| \mathcal{G}_0^{[L+1],(1)}(\cdot) \right\|_{\infty} + t C(4, L) \frac{(\ln m)^C}{m}.
\end{aligned}$$

Finally, estimates shall be made on  $\left\| \mathcal{G}_0^{[L+1],(1)}(\cdot) \right\|_{\infty}$ . We rewrite  $\mathcal{G}_t^{[L+1],(1)}(\cdot)$  into

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

We observe that each term in  $\mathcal{G}_t^{[L+1],(1)}(\cdot)$  is of the form

$$\frac{c}{m} \left\langle \mathbf{B} \mathbf{a}_t, \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \left\langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_{\beta}^{[l]} \right\rangle, \tag{3.71}$$

where  $\mathbf{B}$  is some matrix that varies from term to term,  $c$  is a constant that also changes term by term. Since

$$\left\langle \mathbf{a}_t, \mathbf{B}^{\top} \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \left\langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_{\beta}^{[l]} \right\rangle = \left\langle \mathbf{B} \mathbf{a}_t, \mathbf{x}_{\alpha_1}^{[L]} \right\rangle \left\langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_{\beta}^{[l]} \right\rangle \tag{3.72}$$

holds, then at  $t = 0$ , we shall focus on the term

$$\langle \mathbf{a}_0, \mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle. \quad (3.73)$$

Note that  $\mathbf{a}_0$  is a standard Gaussian vector, from Proposition 3.4.1 and Proposition 3.4.2, with high probability w.r.t random initialization, there exists a uniform constant  $C > 0$ , such that

$$\|\mathbf{B}^\top\|_{2 \rightarrow 2}, \quad \|\mathbf{x}_{\alpha_1}^{[L]}\|_2, \quad \|\mathbf{x}_{\alpha_2}^{[l]}\|_2, \quad \|\mathbf{x}_\beta^{[l]}\|_2 \leq C,$$

after taking conditional expectation,

$$\langle \mathbf{a}_0, \mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle \sim \mathcal{N}\left(0, \left(\langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle\right)^2 \|\mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]}\|_2^2\right), \quad (3.74)$$

apply Lemma 3.4.4 directly, with high probability

$$\frac{c}{m} \langle \mathbf{a}_0, \mathbf{B}^\top \mathbf{x}_{\alpha_1}^{[L]} \rangle \langle \mathbf{x}_{\alpha_2}^{[l]}, \mathbf{x}_\beta^{[l]} \rangle \leq c \frac{(\ln m)^C}{m}. \quad (3.75)$$

Consequently,

$$\|\mathcal{G}_0^{[L+1],(1)}(\cdot)\|_\infty \leq C(3, L) \frac{(\ln m)^C}{m}, \quad (3.76)$$

then for  $1 \leq \alpha_1, \alpha_2, \alpha_3 \leq n$ , and time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

$$\begin{aligned} \left| \mathcal{G}_t^{[L+1],(1)}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}, \mathbf{x}_{\alpha_3}) \right| &\leq \left\| \mathcal{G}_0^{[L+1],(1)}(\cdot) \right\|_\infty + tC(4, L) \frac{(\ln m)^C}{m} \\ &\leq C(3, L) \frac{(\ln m)^C}{m} + tC(4, L) \frac{(\ln m)^C}{m}. \end{aligned}$$

Set  $\mathbf{x}_\beta = \mathbf{x}_{\alpha_3}$ ,

$$\begin{aligned} \left| \partial_t \mathcal{G}_t^{[L+1]}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) \right| &\leq \left( C(3, L) \frac{(\ln m)^C}{m} + tC(4, L) \frac{(\ln m)^C}{m} \right) \sqrt{\frac{\sum_{\beta=1}^n |f_\beta(t) - y_\beta|^2}{n}} \\ &\leq (C(3, L) + tC(4, L)) \frac{(\ln m)^C}{m}, \end{aligned} \quad (3.77)$$

which finishes the proof of Theorem 3.2.2.

### 3.7 Proof of Theorem 3.2.3

Firstly, based on Proposition B.3.2, for  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\varepsilon}\right)$ , then with high probability w.r.t random initialization,

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

set  $\gamma_1 = \varepsilon$ , we finish the first part of the proof.

We move on to the change of the least eigenvalue of the NTK. Recall Equation (3.77) in the proof of Theorem 3.2.2 (Section 3.6), for time  $0 \leq t \leq \sqrt{m}/(\ln m)^{C^*}$ ,

$$\left| \partial_t \mathcal{G}_t^{[L+1]}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) \right| \leq (C(3, L) + tC(4, L)) \frac{(\ln m)^C}{m},$$

consequently,

$$\left| \mathcal{G}_t^{[L+1]}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) - \mathcal{G}_0^{[L+1]}(\mathbf{x}_{\alpha_1}, \mathbf{x}_{\alpha_2}) \right| \leq t(C(3, L) + tC(4, L)) \frac{(\ln m)^C}{m}.$$

The inequality above can be used to derive an upper bound of the change of the least eigenvalue of the  $\mathcal{G}_t^{[L+1]}(\cdot)$ .

$$\begin{aligned} \left\| \left( \mathcal{G}_t^{[L+1]} - \mathcal{G}_0^{[L+1]} \right) (\cdot) \right\|_{2 \rightarrow 2} &\leq \left\| \left( \mathcal{G}_t^{[L+1]} - \mathcal{G}_0^{[L+1]} \right) (\cdot) \right\|_{\text{F}} \leq n \left\| \left( \mathcal{G}_t^{[L+1]} - \mathcal{G}_0^{[L+1]} \right) (\cdot) \right\|_{\infty} \\ &\leq nt(C(3, L) + tC(4, L)) \frac{(\ln m)^C}{m}, \end{aligned}$$

set  $t^*$  satisfying

$$nt^*(C(3, L) + t^*C(4, L)) \frac{(\ln m)^C}{m} = \frac{\lambda_0}{4},$$

rewrite the equation above,

$$C(4, L)(t^*)^2 + C(3, L)t^* = \frac{\lambda_0 m}{4(\ln m)^C n}, \quad (3.78)$$

solving (3.78), we obtain that

$$t^* = \frac{-C(3, L) + \sqrt{(C(3, L))^2 + C(4, L)\lambda_0 \frac{m}{(\ln m)^C n}}}{2C(4, L)}, \quad (3.79)$$

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

$$t^* \geq \frac{1}{4} \sqrt{\frac{\lambda_0 m}{C(4, L)(\ln m)^C n}}. \quad (3.80)$$

Moreover

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

set  $\bar{t} := \inf \left\{ t : \lambda_{\min} \left[ \mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} \geq \lambda_0/2 \right\}$ , then

$$t^* \leq \bar{t}, \quad (3.81)$$

for any  $0 \leq t \leq \bar{t}$ ,

$$\begin{aligned} \partial_t \sum_{\alpha=1}^n \|f_\alpha(t) - y_\alpha\|_2^2 &= -\frac{2}{n} \sum_{\alpha, \beta=1}^n \mathcal{G}_t^{(2)}(\mathbf{x}_\alpha, \mathbf{x}_\beta) (f_\alpha(t) - y_\alpha)(f_\beta(t) - y_\beta) \\ &\leq -\frac{\lambda_0}{n} \sum_{\gamma=1}^n \|f_\gamma(t) - y_\gamma\|_2^2, \end{aligned}$$

then

$$R_S(\boldsymbol{\theta}_t) \leq \exp \left( -\frac{\lambda_0 t}{n} \right) R_S(\boldsymbol{\theta}_0). \quad (3.82)$$

Set  $R_S(\boldsymbol{\theta}_t) = \varepsilon$ , it takes time

$$t \leq \frac{n}{\lambda_0} \ln\left(\frac{C^*}{\varepsilon}\right)$$

for  $R_S(\boldsymbol{\theta}_t)$  to reach accuracy  $\varepsilon$ , hence if

$$t \leq \frac{n}{\lambda} \ln\left(\frac{C^*}{\varepsilon}\right) \leq t^* \leq \bar{t}, \quad (3.83)$$

then width  $m$  is required to yield the lower bound for  $t^*$  derived in Equation (3.80),

$$\frac{n}{\lambda_0} \ln\left(\frac{C^*}{\varepsilon}\right) \leq \frac{1}{4} \sqrt{\frac{\lambda_0 m}{C(4, L)(\ln m)^C n}}. \quad (3.84)$$

then for

$$m \geq C(4, L) \left(\frac{n}{\lambda_0}\right)^3 (\ln m)^C \ln\left(\frac{C^*}{\varepsilon}\right)^2,$$

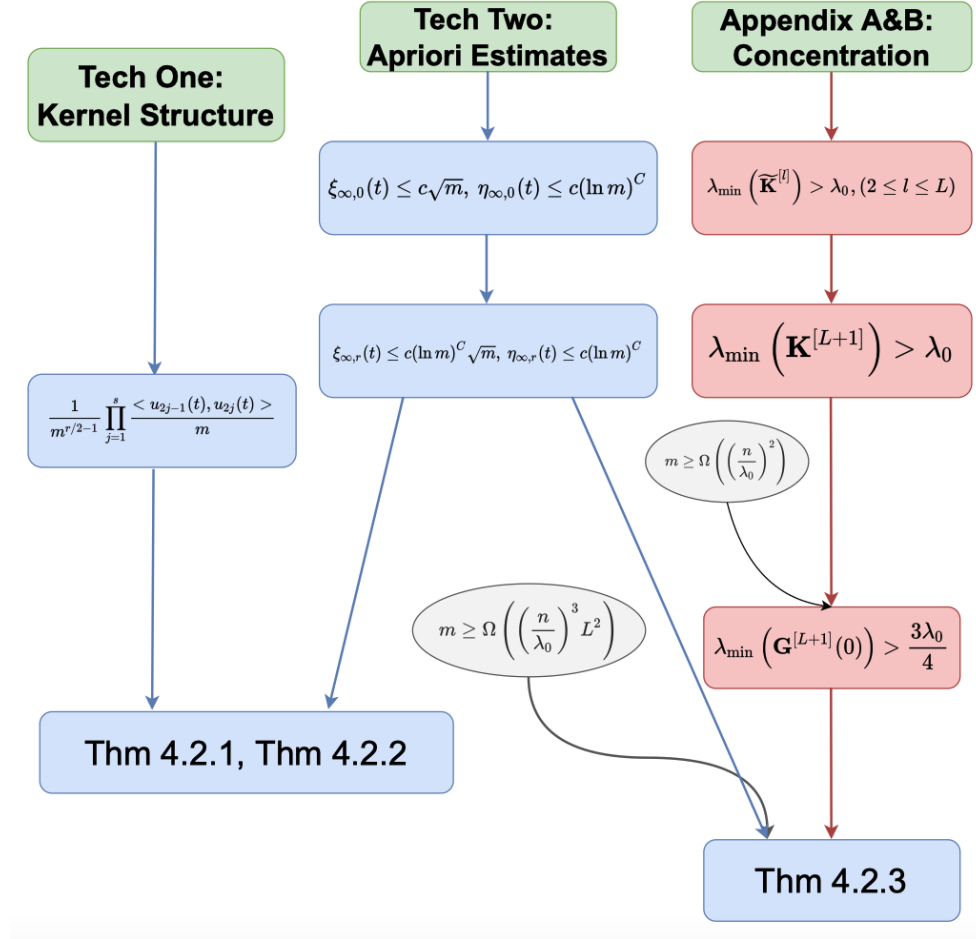
since  $C(4, L) = \mathcal{O}(L^2)$ , we conclude that the required width  $m$  should be

$$m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^3 L^2 (\ln m)^C \ln\left(\frac{C^*}{\varepsilon}\right)^2\right), \quad (3.85)$$

where  $\varepsilon$  is the desired training accuracy.

## 4. SUMMARY AND FUTURE WORK

In this thesis, we described a framework for analyzing the training behavior of ResNet in continuous time gradient descent dynamics. At a high level, our approach is illustrated using Figure 4.1, which builds upon sharp analysis of the least eigenvalue of randomly initialized Gram matrix, and upon uniform estimates for kernels of all orders in the NTH. Finally, we would like to re-emphasize the significance of the  $\frac{c_{\text{res}}}{L}$  scaling placed in ResNet, which has been proven to be successful in the stabilization of random propagation (Proposition B.3.2, Equation (B.36)), in the uniform estimate of  $\xi(t)$  (Proposition 3.4.1),  $\xi_{\infty,0}(t)$  (Proposition 3.4.3, Equation (3.43)), and  $\eta_{\infty,r}(t)$  (Proposition 3.4.6) with  $r \geq 1$ , and most importantly, in lowering the magnitude of  $C(4, L)$  (Equation (3.70)).



**Figure 4.1.** Diagram for the Proof of Main Theorems

Some future directions are listed out for our research:

- The NTH is an infinite sequence of relationship. Huang and Yau showed that under certain conditions, the NTH can be truncated and its truncated version is still able to approximate the original dynamics up to any precision. We have faith in that for ResNet, such technical conditions can be loosened.
- In Theorem 3.2.3, the dependence of  $m$  on the depth  $L$  is quadratic, we believe that the dependence can be reduced even further. We conjecture that  $m$  depends linearly in  $L$ .
- In this thesis, we focus on gradient descent, and we believe that it can be extended to stochastic gradient descent, while maintaining the linear convergence rate.
- The test loss has not been addressed in our work. To further investigate the generalization power of ResNet, some Apriori estimates for its generalization error would be useful [42].

## REFERENCES

- [1] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai, “Gradient descent finds global minima of deep neural networks,” in *Proceedings of the 36th International Conference on Machine Learning*, K. Chaudhuri and R. Salakhutdinov, Eds., ser. Proceedings of Machine Learning Research, vol. 97, PMLR, Sep. 2019, pp. 1675–1685. [Online]. Available: <http://proceedings.mlr.press/v97/du19c.html>.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- [3] J. Huang and H.-T. Yau, “Dynamics of deep neural networks and neural tangent hierarchy,” in *Proceedings of the 37th International Conference on Machine Learning*, H. D. III and A. Singh, Eds., ser. Proceedings of Machine Learning Research, vol. 119, PMLR, 13–18 Jul 2020, pp. 4542–4551. [Online]. Available: <http://proceedings.mlr.press/v119/huang20l.html>.
- [4] A. Jacot, F. Gabriel, and C. Hongler, “Neural tangent kernel: Convergence and generalization in neural networks,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31, Curran Associates, Inc., 2018, pp. 8571–8580. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf>.
- [5] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT press, 2016.
- [6] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding deep learning requires rethinking generalization,” in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, OpenReview.net, 2017. [Online]. Available: <https://openreview.net/forum?id=Sy8gdB9xx>.
- [7] S. S. Du, X. Zhai, B. Póczos, and A. Singh, “Gradient descent provably optimizes over-parameterized neural networks,” in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=S1eK3i09YQ>.
- [8] D. Zou, Y. Cao, D. Zhou, and Q. Gu, “Gradient descent optimizes over-parameterized deep relu networks,” *Machine Learning*, vol. 109, no. 3, pp. 467–492, 2020.



- [9] S. Arora, R. Ge, B. Neyshabur, and Y. Zhang, “Stronger generalization bounds for deep nets via a compression approach,” in *Proceedings of the 35th International Conference on Machine Learning*, J. Dy and A. Krause, Eds., ser. Proceedings of Machine Learning Research, vol. 80, Stockholmsmässan, Stockholm Sweden: PMLR, Oct. 2018, pp. 254–263. [Online]. Available: <http://proceedings.mlr.press/v80/arora18b.html>.
- [10] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, *et al.*, “Mastering the Game of Go with Deep Neural Networks and Tree Search,” *Nature*, vol. 529, no. 7587, p. 484, 2016.
- [11] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, *et al.*, “Mastering the Game of Go without Human Knowledge,” *Nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [12] S. Zagoruyko and N. Komodakis, “Wide residual networks,” in *Proceedings of the British Machine Vision Conference (BMVC)*, E. R. H. Richard C. Wilson and W. A. P. Smith, Eds., BMVA Press, Sep. 2016, pp. 87.1–87.12, ISBN: 1-901725-59-6. DOI: [10.5244/C.30.87](https://dx.doi.org/10.5244/C.30.87). [Online]. Available: <https://dx.doi.org/10.5244/C.30.87>.
- [13] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, “Inception-v4, inception-resnet and the impact of residual connections on learning,” in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, ser. AAAI’17, San Francisco, California, USA: AAAI Press, 2017, pp. 4278–4284.
- [14] M. Hardt and T. Ma, “Identity matters in deep learning,” *ICLR*, 2017. [Online]. Available: <https://arxiv.org/abs/1611.04231>.
- [15] G. Cybenko, “Approximation by superpositions of a sigmoidal function,” *Mathematics of control, signals and systems*, vol. 2, no. 4, pp. 303–314, 1989.
- [16] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, vol. 15, no. 56, pp. 1929–1958, 2014. [Online]. Available: <http://jmlr.org/papers/v15/srivastava14a.html>.
- [17] S. Mei, A. Montanari, and P.-M. Nguyen, “A mean field view of the landscape of two-layer neural networks,” *Proceedings of the National Academy of Sciences*, vol. 115, no. 33, E7665–E7671, 2018, ISSN: 0027-8424. DOI: [10.1073/pnas.1806579115](https://doi.org/10.1073/pnas.1806579115). eprint: <https://www.pnas.org/content/115/33/E7665.full.pdf>. [Online]. Available: <https://www.pnas.org/content/115/33/E7665>.

- [18] G. Rotskoff and E. Vanden-Eijnden, “Parameters as interacting particles: Long time convergence and asymptotic error scaling of neural networks,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31, Curran Associates, Inc., 2018, pp. 7146–7155. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/file/196f5641aa9dc87067da4ff90fd81e7b-Paper.pdf>.
- [19] L. Chizat and F. Bach, “On the global convergence of gradient descent for over-parameterized models using optimal transport,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31, Curran Associates, Inc., 2018, pp. 3036–3046. [Online]. Available: <https://arxiv.org/pdf/1805.09545.pdf>.
- [20] T. Luo, Z.-Q. J. Xu, Z. Ma, and Y. Zhang, “Phase diagram for two-layer relu neural networks at infinite-width limit,” *CoRR*, vol. abs/2007.07497, 2020. arXiv: [2007.07497](https://arxiv.org/abs/2007.07497). [Online]. Available: <https://arxiv.org/abs/2007.07497>.
- [21] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, “Deep networks with stochastic depth,” in *European conference on computer vision*, Springer, 2016, pp. 646–661.
- [22] R. Eldan and O. Shamir, “The power of depth for feedforward neural networks,” in *29th Annual Conference on Learning Theory*, V. Feldman, A. Rakhlin, and O. Shamir, Eds., ser. Proceedings of Machine Learning Research, vol. 49, Columbia University, New York, New York, USA: PMLR, 23–26 Jun 2016, pp. 907–940. [Online]. Available: <http://proceedings.mlr.press/v49/eldan16.html>.
- [23] Y. Bengio, P. Simard, and P. Frasconi, “Learning long-term dependencies with gradient descent is difficult,” *IEEE Transactions on Neural Networks*, vol. 5, no. 2, pp. 157–166, 1994. DOI: [10.1109/72.279181](https://doi.org/10.1109/72.279181).
- [24] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *Proceedings of the 32nd International Conference on Machine Learning*, F. Bach and D. Blei, Eds., ser. Proceedings of Machine Learning Research, vol. 37, Lille, France: PMLR, Jul. 2015, pp. 448–456. [Online]. Available: <http://proceedings.mlr.press/v37/ioffe15.html>.
- [25] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, Y. W. Teh and M. Titterton, Eds., ser. Proceedings of Machine Learning Research, vol. 9, Chia Laguna Resort, Sardinia, Italy: JMLR Workshop and Conference Proceedings, 13–15 May 2010, pp. 249–256. [Online]. Available: <http://proceedings.mlr.press/v9/glorot10a.html>.

- [26] K. He and J. Sun, “Convolutional neural networks at constrained time cost,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 5353–5360. DOI: [10.1109/CVPR.2015.7299173](https://doi.org/10.1109/CVPR.2015.7299173).
- [27] S. S. Du, C. Jin, J. D. Lee, M. I. Jordan, A. Singh, and B. Poczos, “Gradient descent can take exponential time to escape saddle points,” in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30, Curran Associates, Inc., 2017, pp. 1067–1077. [Online]. Available: <https://proceedings.neurips.cc/paper/2017/file/f79921bbae40a577928b76d2fc3edc2a-Paper.pdf>.
- [28] C. Jin, R. Ge, P. Netrapalli, S. M. Kakade, and M. I. Jordan, “How to escape saddle points efficiently,” in *Proceedings of the 34th International Conference on Machine Learning*, D. Precup and Y. W. Teh, Eds., ser. Proceedings of Machine Learning Research, vol. 70, International Convention Centre, Sydney, Australia: PMLR, Jun. 2017, pp. 1724–1732. [Online]. Available: <http://proceedings.mlr.press/v70/jin17a.html>.
- [29] R. Ge, F. Huang, C. Jin, and Y. Yuan, “Escaping from saddle points — online stochastic gradient for tensor decomposition,” in *Proceedings of The 28th Conference on Learning Theory*, P. Grünwald, E. Hazan, and S. Kale, Eds., ser. Proceedings of Machine Learning Research, vol. 40, Paris, France: PMLR, Mar. 2015, pp. 797–842. [Online]. Available: <http://proceedings.mlr.press/v40/Ge15.html>.
- [30] Q. Nguyen and M. Hein, “The loss surface of deep and wide neural networks,” in *Proceedings of the 34th International Conference on Machine Learning*, D. Precup and Y. W. Teh, Eds., ser. Proceedings of Machine Learning Research, vol. 70, International Convention Centre, Sydney, Australia: PMLR, Jun. 2017, pp. 2603–2612. [Online]. Available: <http://proceedings.mlr.press/v70/nguyen17a.html>.
- [31] K. Kawaguchi, “Deep learning without poor local minima,” in *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, Eds., vol. 29, Curran Associates, Inc., 2016, pp. 586–594. [Online]. Available: <https://proceedings.neurips.cc/paper/2016/file/f2fc990265c712c49d51a18a32b39f0c-Paper.pdf>.
- [32] S. Arora, N. Cohen, N. Golowich, and W. Hu, “A convergence analysis of gradient descent for deep linear neural networks,” in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=SkMQg3C5K7>.

- [33] P. Bartlett, D. Helmbold, and P. Long, “Gradient descent with identity initialization efficiently learns positive definite linear transformations by deep residual networks,” in *Proceedings of the 35th International Conference on Machine Learning*, J. Dy and A. Krause, Eds., ser. Proceedings of Machine Learning Research, vol. 80, Stockholmsmässan, Stockholm Sweden: PMLR, Oct. 2018, pp. 521–530. [Online]. Available: <http://proceedings.mlr.press/v80/bartlett18a.html>.
- [34] Z. Allen-Zhu, Y. Li, and Z. Song, “A convergence theory for deep learning via overparameterization,” in *Proceedings of the 36th International Conference on Machine Learning*, K. Chaudhuri and R. Salakhutdinov, Eds., ser. Proceedings of Machine Learning Research, vol. 97, PMLR, Sep. 2019, pp. 242–252. [Online]. Available: <http://proceedings.mlr.press/v97/allen-zhu19a.html>.
- [35] S. Arora, S. S. Du, W. Hu, Z. Li, R. R. Salakhutdinov, and R. Wang, “On exact computation with an infinitely wide neural net,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32, Curran Associates, Inc., 2019, pp. 8141–8150. [Online]. Available: <https://proceedings.neurips.cc/paper/2019/file/dbc4d84bfcfe2284ba11beffb853a8c4-Paper.pdf>.
- [36] Y. Li and Y. Liang, “Learning overparameterized neural networks via stochastic gradient descent on structured data,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31, Curran Associates, Inc., 2018, pp. 8157–8166. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/file/54fe976ba170c19ebae453679b362263-Paper.pdf>.
- [37] R. M. Neal, *Bayesian learning for neural networks*. Springer Science & Business Media, 2012, vol. 118.
- [38] G. 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.
- [39] J. Lee, L. Xiao, S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, and J. Pennington, “Wide neural networks of any depth evolve as linear models under gradient descent,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32, Curran Associates, Inc., 2019. [Online]. Available: <https://proceedings.neurips.cc/paper/2019/file/0d1a9651497a38d8b1c3871c84528bd4-Paper.pdf>.
- [40] R. Vershynin, “Introduction to the Non-asymptotic Analysis of Random Matrices,” *arXiv preprint arXiv:1011.3027*, 2010.

- [41] B. Laurent and P. Massart, “Adaptive Estimation of a Quadratic Functional by Model Selection,” *Annals of Statistics*, pp. 1302–1338, 2000.
- [42] C. Ma, Q. Wang, and E. Weinan, “A Priori Estimates of the Population Risk for Residual Networks,” *arXiv preprint arXiv:1903.02154*, 2019.
- [43] S. Boucheron, G. Lugosi, and P. Massart, *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.

## A. LEAST EIGENVALUE OF GRAM MATRICES

### A.1 Introduction

We shall recall the series of matrices  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$  and vectors  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$  given in Definition 2.4.4:

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

Given expressions above, we shall recall definitions of  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$  (Definition 2.4.5):

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

We shall state two lemmas concerning full rankness of  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$ , which have been stated as [1, Lemma F.1, Lemma F.2].

**Lemma A.1.1** (Lemma F.1 in [1]). Assume  $\sigma(\cdot)$  is analytic and not a polynomial function. Consider input data set  $\mathcal{U} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$  of  $n$  non-parallel points, i.e.,  $\mathbf{u}_j \notin \text{span}(\mathbf{u}_k)$  for any  $j \neq k$ . Define

$$\mathbf{G}(\mathcal{U})_{ij} := \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\sigma(\mathbf{w}^\top \mathbf{u}_i) \sigma(\mathbf{w}^\top \mathbf{u}_j)], \quad (\text{A.1})$$

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

**Lemma A.1.2** (Lemma F.2 in [1]). Assume  $\sigma(\cdot)$  is analytic and not a polynomial function. Consider input data set  $\mathcal{U} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$  of  $n$  non-parallel points, i.e.,  $\mathbf{u}_j \notin \text{span}(\mathbf{u}_k)$  for any  $j \neq k$ . Define

$$\mathbf{G}(\mathcal{U})_{ij} := \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\sigma^{(1)}(\mathbf{w}^\top \mathbf{u}_i) \sigma^{(1)}(\mathbf{w}^\top \mathbf{u}_j) (\mathbf{u}_i^\top \mathbf{u}_j)], \quad (\text{A.2})$$

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

Given Lemma A.1.1 and Lemma A.1.2, we may proceed to quantify the least eigenvalues of  $\mathbf{K}^{[L+1]}$  and  $\mathbf{K}^{[L]}$ .

## A.2 Full Rankness for $(L + 1)$ -th Gram matrix

**Lemma A.2.1.** Given  $\mathcal{X}$ ,  $\sigma(\cdot)$ ,  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ , and  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , then for each  $l$ , each diagonal entry of  $\widetilde{\mathbf{K}}^{[l]}$  is the same with each other, and each element of the vector  $\widetilde{\mathbf{b}}^{[l]}$  is also the same, i.e.,

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

Moreover

$$\left(1 - \frac{l}{L} \frac{c}{\sqrt{c_\sigma}}\right)^2 \leq \widetilde{\mathbf{K}}_{ii}^{[l]} \leq \left(1 + \frac{l}{L} \frac{c}{\sqrt{c_\sigma}}\right)^2, \quad (\text{A.3})$$

and

$$(\widetilde{\mathbf{b}}_i^{[l]})^2 < \widetilde{\mathbf{K}}_{ii}^{[l]}, \quad (\text{A.4})$$

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

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

Next we need to show that (A.3), (A.4) hold true for  $l = 1$ . Based on definition, we recall that  $c_\sigma = \left(\mathbb{E}_{x \sim \mathcal{N}(0,1)} [\sigma(x)^2]\right)^{-1}$ ,

$$\begin{aligned}\mathbf{K}_{ii}^{[1]} &= c_\sigma \mathbb{E}_{u \sim \mathcal{N}\left(0, \widetilde{\mathbf{K}}_{ii}^{[0]}\right)} \left(\sigma(u)^2\right) = c_\sigma \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left(\sigma(u)^2\right) = 1, \\ \widetilde{\mathbf{b}}_i^{[1]} &= \sqrt{c_\sigma} \mathbb{E}_{u \sim \mathcal{N}\left(0, \widetilde{\mathbf{K}}_{ii}^{[0]}\right)} [\sigma(u)] = \sqrt{c_\sigma} \mathbb{E}_{u \sim \mathcal{N}(0,1)} [\sigma(u)],\end{aligned}$$

then

$$\left(\widetilde{\mathbf{b}}_i^{[1]}\right)^2 = c_\sigma \left(\mathbb{E}_{u \sim \mathcal{N}\left(0, \widetilde{\mathbf{K}}_{ii}^{[0]}\right)} [\sigma(u)]\right)^2 < 1, \quad (\text{A.5})$$

Equation (A.5) 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 choice of the index  $i$ , then  $\widetilde{\mathbf{K}}_{ii}^{[1]} = \widetilde{\mathbf{K}}_{jj}^{[1]}$ ,  $\widetilde{\mathbf{b}}_i^{[1]} = \widetilde{\mathbf{b}}_j^{[1]}$ , for any  $i \neq j$ .

Now we assume that (A.3), (A.4) hold for  $1, 2, \dots, l-1$ , and we want to show that it is also the case for  $l$ . First of all, since

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

then such quantities  $\widetilde{\mathbf{K}}_{ii}^{[l]}$ ,  $\widetilde{\mathbf{b}}_i^{[l]}$  are independent of choice of the index  $i$ .



Secondly, concerning  $\tilde{\mathbf{b}}_i^{[l]}$ , we assume Equation (A.4) holds for  $l = 1, 2, \dots, l-1$ , then

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

Finally, for  $\tilde{\mathbf{K}}_{\text{ii}}^{[l]}$ , we note that the following inequality holds

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

Since  $\sigma(\cdot)$  is 1-Lipschitz, then for any  $1/2 \leq \alpha \leq 2$ , we have

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

then

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

Based on our induction hypothesis,

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

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

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

we shall choose  $c$  wisely, for instance, set  $c$  as

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

by our choice of  $c$ , combined with Equation (A.6), we have

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

then

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

which finishes our proof.  $\square$

Our next lemma is crucial in that it reveals a covariance type structure for the series of matrices  $\left\{ \widetilde{\mathbf{K}}^{[l]} \right\}_{l=1}^L$ . A standard notation related to matrices shall be introduced. We denote that  $\mathbf{A} \succeq \mathbf{B}$  if and only if  $\mathbf{A} - \mathbf{B}$  is a semi-positive definite matrix, and  $\mathbf{A} \succ \mathbf{B}$  if and only if  $\mathbf{A} - \mathbf{B}$  is strictly positive definite.

**Proposition A.2.1.** *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$ ,  $\left\{ \widetilde{\mathbf{K}}^{[l]} \right\}_{l=1}^L$ , and  $\left\{ \widetilde{\mathbf{b}}^{[l]} \right\}_{l=1}^L$ , then for each  $l$ ,*

$$\widetilde{\mathbf{K}}^{[l]} - \widetilde{\mathbf{b}}^{[l]} \otimes \left( \widetilde{\mathbf{b}}^{[l]} \right)^\top \succeq \widetilde{\mathbf{K}}^{[l-1]} - \widetilde{\mathbf{b}}^{[l-1]} \otimes \left( \widetilde{\mathbf{b}}^{[l-1]} \right)^\top. \quad (\text{A.6})$$

Set  $\lambda_0 > 0$  as

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

then for all  $2 \leq l \leq L$ ,

$$\lambda_{\min} \left( \widetilde{\mathbf{K}}^{[l]} \right) \geq \lambda_0. \quad (\text{A.7})$$

*Proof.* Since for  $1 \leq i, j \leq n$ , and  $1 \leq l \leq L$ ,

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

We define another series of covariance matrices  $\{\mathbf{P}^{[s]} : 1 \leq s \leq L\}$ ,

$$\mathbf{P}_{ij}^{[s]} := \frac{c_{\text{res}}^2}{L^2} \text{Cov}_{(u,v) \sim \mathcal{N}(0, \widetilde{\mathbf{A}}_{ij}^{[s+1]})} [\sigma(u) \sigma(v)], \quad 1 \leq s \leq L.$$

We notice that  $\mathbf{P}^{[s]}$  are covariance matrices, hence naturally  $\mathbf{P}^{[s]} \succeq 0$ . Inductively, we have for all  $l$ ,

$$\begin{aligned}
\widetilde{\mathbf{K}}^{[l]} &\succeq \widetilde{\mathbf{K}}^{[l]} - \widetilde{\mathbf{b}}^{[l]} \otimes (\widetilde{\mathbf{b}}^{[l]})^\top \\
&= \widetilde{\mathbf{K}}^{[l-1]} - \widetilde{\mathbf{b}}^{[l-1]} \otimes (\widetilde{\mathbf{b}}^{[l-1]})^\top + \mathbf{P}^{[l-1]} \\
&\succeq \widetilde{\mathbf{K}}^{[l-1]} - \widetilde{\mathbf{b}}^{[l-1]} \otimes (\widetilde{\mathbf{b}}^{[l-1]})^\top \\
&= \widetilde{\mathbf{K}}^{[l-2]} - \widetilde{\mathbf{b}}^{[l-2]} \otimes (\widetilde{\mathbf{b}}^{[l-2]})^\top + \mathbf{P}^{[l-2]} \\
&\quad \vdots \\
&\succeq \widetilde{\mathbf{K}}^{[1]} - \widetilde{\mathbf{b}}^{[1]} \otimes (\widetilde{\mathbf{b}}^{[1]})^\top,
\end{aligned} \tag{A.8}$$

the last line brings us to notice of the entry of  $\widetilde{\mathbf{K}}^{[1]} - \widetilde{\mathbf{b}}^{[1]} \otimes (\widetilde{\mathbf{b}}^{[1]})^\top$ , which reads respectively

$$\left[ \widetilde{\mathbf{K}}^{[1]} - \widetilde{\mathbf{b}}^{[1]} \otimes (\widetilde{\mathbf{b}}^{[1]})^\top \right]_{ij} = c_\sigma \text{Cov}_{(u,v) \sim \mathcal{N}(0, \widetilde{\mathbf{A}}_{ij}^{[1]})} [\sigma(u) \sigma(v)].$$

We notice that  $\left[\widetilde{\mathbf{K}}^{[1]} - \widetilde{\mathbf{b}}^{[1]} \otimes (\widetilde{\mathbf{b}}^{[1]})^\top\right]$  is a covariance matrix, hence naturally

$$\left[\widetilde{\mathbf{K}}^{[1]} - \widetilde{\mathbf{b}}^{[1]} \otimes (\widetilde{\mathbf{b}}^{[1]})^\top\right] \succeq 0,$$

and  $\left[\widetilde{\mathbf{K}}^{[1]} - \widetilde{\mathbf{b}}^{[1]} \otimes (\widetilde{\mathbf{b}}^{[1]})^\top\right] \succ 0$  except when one of the samples is an exact linear combination of the others. From Lemma A.1.1, we see that  $\widetilde{\mathbf{K}}^{[1]} \succ 0$ , consequently

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

then  $\lambda_0 > 0$ , and it depends solely on the input samples and activation function.  $\square$

**Proposition A.2.2.** *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$ ,  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ , and  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , if we set  $\lambda_0 > 0$  as*

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

then

$$\lambda_{\min} \left( \mathbf{K}^{[L+1]} \right) \geq \lambda_0. \quad (\text{A.9})$$

*Proof.* Proof of Proposition A.2.2 is quite similar to the proof of Proposition A.2.1. We recall that

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

and we define further that

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

then

$$\mathbf{K}_{ij}^{[L+1]} - \mathbf{b}_i^{[L+1]} \mathbf{b}_j^{[L+1]} = \widetilde{\mathbf{K}}_{ij}^{[L]} - \widetilde{\mathbf{b}}_i^{[L]} \widetilde{\mathbf{b}}_j^{[L]} + \frac{c_{\text{res}}^2}{L^2} \text{Cov}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{A}}_{ij}^{[L+1]})} [\sigma(u) \sigma(v)],$$

hence by same reasoning in Equation (A.8),

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

apply Proposition A.2.1 directly, we are able to finish the proof.  $\square$

From Proposition A.2.2, we observe that  $\lambda_{\min}(\mathbf{K}^{[L+1]}) \sim \Omega(1)$ .

### A.3 Full Rankness for the $L$ -th Gram matrix

Our next proposition concerns the least eigenvalue of  $\mathbf{K}^{[L]}$ , which has been stated as Proposition F.2 in Du et al. [1].

**Proposition A.3.1.** *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$ ,  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ , and  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , then*

$$\lambda_{\min}(\mathbf{K}^{[L]}) \geq \frac{c_{\text{res}}^2}{L^2} \kappa, \quad (\text{A.10})$$

where  $\kappa$  is independent of depth  $L$ .

*Proof.* Based on Lemma A.2.1, there exists a constant  $c > 0$ , such that

$$1/c \leq \widetilde{\mathbf{K}}_{\text{ii}}^{[L]} \leq c.$$

We define function  $\mathbf{G} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$ :

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

from which a scalar function  $g(\lambda)$  could be defined as

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

and Lemma A.1.2 guarantees that

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

From Proposition A.2.1, we observe that

$$\lambda_{\min} \left( \widetilde{\mathbf{K}}^{[L-1]} \right) \geq \lambda_0,$$

hence

$$\lambda_{\min} \left( \mathbf{K}^{[L]} \right) \geq \frac{c_{\text{res}}^2}{L^2} g(\lambda_0). \tag{A.11}$$

Set  $\kappa = g(\lambda_0) > 0$ , since  $\kappa$  is independent of depth  $L$ , we finish our proof.  $\square$

From Proposition A.3.1, we remark that  $\lambda_{\min} \left( \mathbf{K}^{[L]} \right) \sim \Omega(\frac{1}{L^2})$ .

## B. RANDOM INITIALIZATION OF NTK

In this chapter, we are going to show that

$$\lambda_{\min} \left( \left[ \mathcal{G}_0^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta) \right]_{1 \leq \alpha, \beta \leq n} \right) \geq \frac{3\lambda_0}{4}, \quad (\text{B.1})$$

holds with high probability, where  $\lambda_0 > 0$  is set to be the least eigenvalue of matrix  $\widetilde{\mathbf{K}}^{[1]} - \widetilde{\mathbf{b}}^{[1]} \otimes (\widetilde{\mathbf{b}}^{[1]})^\top$ , i.e.,

$$\lambda_0 = \lambda_{\min} \left( \widetilde{\mathbf{K}}^{[1]} - \widetilde{\mathbf{b}}^{[1]} \otimes (\widetilde{\mathbf{b}}^{[1]})^\top \right).$$

### B.1 Several Lemmas on Gaussian Concentration and Other Aspects

We investigate the concentration properties of Lipschitz functions of Gaussian variables. Let us say that a function  $f(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}$  is  $C_L$ -Lipschitz function with respect to the Euclidean norm  $\|\cdot\|_2$ , if for any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ ,

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq C_L \|\mathbf{x} - \mathbf{y}\|_2.$$

Our next lemma reveals that any Lipschitz function of Gaussian variables is itself a sub-Gaussian variable.

**Lemma B.1.1** (Gaussian Concentration Inequality). *Let  $\mathbf{X} = (X_1, \dots, X_p)^\top \in \mathbb{R}^p$ , whose components  $X_1, X_2, \dots, X_p$  are i.i.d. Gaussian variables drawn from  $\mathcal{N}(0, \sigma^2)$ , and  $f(\cdot) : \mathbb{R}^p \rightarrow \mathbb{R}$  is a  $C_L$ -Lipschitz function with respect to the Euclidean norm  $\|\cdot\|_2$ , then for all  $t \geq 0$ :*

$$\mathbb{P}(|f(\mathbf{X}) - \mathbb{E} f(\mathbf{X})| \geq t) \leq 2 \exp\left(-\frac{t^2}{2C_L^2 \sigma^2}\right). \quad (\text{B.2})$$

We refer to [43, Theorem 5.6] for the proof of Lemma B.1.1, and remarkably, this is a dimension free inequality.

Next we shall state two lemmas, which have been stated as [1, Lemma G.3, Lemma G.4]

**Lemma B.1.2** (Lemma G.3 in [1]). *If  $\sigma(\cdot)$  is  $C_L$ -Lipschitz, then for  $a, b > 0$ , satisfying*

$$1/c \leq \min(a, b), \max(a, b) \leq c,$$

for some  $c > 0$ , we have

$$\left| \mathbb{E}_{z \sim \mathcal{N}(0,1)} [\sigma(az)] - \mathbb{E}_{z \sim \mathcal{N}(0,1)} [\sigma(bz)] \right| \leq C |a - b|, \quad (\text{B.3})$$

where  $C > 0$  only depends on  $c$  and Lipschitz constant  $C_L$ .

**Lemma B.1.3** (Lemma G.4 in [1]). *If  $\sigma(\cdot)$  is  $C_L$ -Lipschitz, define  $F(\mathbf{K}) : \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}$*

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

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

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

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

whose entries satisfy

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

$$-1 \leq \rho_1, \rho_2 \leq 1$$

for some  $c > 0$ . Then, we have

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

where  $C > 0$  only depends on  $c$  and Lipschitz constant  $C_L$ .



## B.2 Analysis of Random Propagation

We recall readers once again the series of matrices  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ ,  $\{\widetilde{\mathbf{A}}^{[l]}\}_{l=1}^{L+1}$  and vectors  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$  given in Definition 2.4.4:

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

We shall begin with a proposition on the tail probabilities concerning the outputs of each layer at initial state, denoted by  $\mathbf{x}_i^{[l]}(0)$  for all  $1 \leq i \leq n$  and  $1 \leq l \leq L$ .

**Proposition B.2.1.** *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$ ,  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ , and  $\{\widetilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , then for all  $t > 0$ ,  $1 \leq i \leq n$ ,  $1 \leq l \leq L$ :*

$$\mathbb{P} \left( \left| \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 - \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \right| \geq t \right) \leq \exp(-cmt^2), \quad (\text{B.4})$$

$$\mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \widetilde{\mathbf{b}}_i^{[l]} \right| \geq t \right) \leq \exp(-cmt^2), \quad (\text{B.5})$$

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

*Proof.* (i). For  $l = 1$ , we have

$$\left\| \mathbf{x}_i^{[1]}(0) \right\|_2^2 = \frac{c_\sigma}{m} \sum_{j=1}^m \sigma \left( (\mathbf{W}_0^{[1]} \mathbf{x}_i)_j \right)^2,$$

with

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

Since for all  $i$ ,  $\|\mathbf{x}_i\|_2 = 1$ , then for each  $j = 1, 2, \dots, m$ ,  $(\mathbf{W}_0^{[1]} \mathbf{x}_i)_j$  is a standard Gaussian, i.e.,  $(\mathbf{W}_0^{[1]} \mathbf{x}_i)_j \sim \mathcal{N}(0, 1)$ . Moreover,  $\sigma(\cdot)$  is 1-Lipschitz, then  $\sigma((\mathbf{W}_0^{[1]} \mathbf{x}_i)_j)$  is sub-Gaussian, hence  $\left\{ \sigma((\mathbf{W}_0^{[1]} \mathbf{x}_i)_j)^2 \right\}_{j=1}^m$  is a collection of i.i.d. sub-exponential variables. Consequently, for all  $t > 0$ ,

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

hence,

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

We have shown already for  $l = 1$ , Equation (B.4) holds. We need to show further that Equation (B.5) holds.

We shall note that

$$\left\langle \frac{\mathbf{x}_i^{[1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle = \frac{\sqrt{c_\sigma}}{m} \sum_{j=1}^m \sigma((\mathbf{W}_0^{[1]} \mathbf{x}_i)_j),$$

and

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

Since  $\mathbf{x}_i^{[1]}(0)$  can be written into the form

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

where  $\mathbf{X}$  is a standard normal Gaussian vector.

We shall focus on the inner product function  $g^{[1]}(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$ ,

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

We show that  $g^{[1]}(\cdot)$  is also Lipschitz with respect to the Euclidean norm  $\|\cdot\|_2$ :

For any  $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{R}^m$ :

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

Hence  $g^{[1]}(\cdot)$  is  $\frac{C}{\sqrt{m}}$ -Lipschitz. Apply Lemma B.1.1 directly, we have

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

then

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

(ii). Our next step is to prove that Equation (B.4) and Equation (B.5) hold for  $l \geq 2$ , and we will prove it by induction. Assume that Equation (B.4) and Equation (B.5) hold true for  $1, 2, 3, \dots, l$ , and we show further that it is also the case for  $l+1$ , i.e.,

$$\mathbb{P} \left( \left| \left\| \mathbf{x}_i^{[l+1]}(0) \right\|_2 - \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l+1]}} \right| \geq t \right) \leq \exp(-cmt^2), \quad (\text{B.6})$$

$$\mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l+1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \tilde{\mathbf{b}}_i^{[l+1]} \right| \geq t \right) \leq \exp(-cmt^2). \quad (\text{B.7})$$

The structure of outputs of each layer at initial state shall be recalled as follows

$$\mathbf{x}_i^{[l+1]}(0) = \mathbf{x}_i^{[l]}(0) + \frac{c_{\text{res}}}{L\sqrt{m}} \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right),$$

and the counterpart of  $\widetilde{\mathbf{K}}_{\text{ii}}^{[l]}$  and  $\widetilde{\mathbf{b}}_{\text{i}}^{[l]}$

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

then

$$\begin{aligned}\|\mathbf{x}_{\text{i}}^{[l+1]}(0)\|_2^2 &= \|\mathbf{x}_{\text{i}}^{[l]}(0)\|_2^2 + \underbrace{\frac{2c_{\text{res}}}{L} \left\langle \frac{\mathbf{x}_{\text{i}}^{[l]}(0)}{\sqrt{m}}, \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_{\text{i}}^{[l]}(0)) \right\rangle}_{\text{I}} \\ &\quad + \underbrace{\frac{c_{\text{res}}^2}{L^2} \frac{1}{m} \left\langle \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_{\text{i}}^{[l]}(0)), \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_{\text{i}}^{[l]}(0)) \right\rangle}_{\text{II}}.\end{aligned}$$

Based on our induction hypothesis, there exists estimates on  $\|\mathbf{x}_{\text{i}}^{[l]}(0)\|_2^2$ . We need to focus on terms I and II. Firstly for term I, there is a  $\frac{1}{\sqrt{m}}$  scaling factor contained in  $\mathbf{x}_{\text{i}}^{[l]}(0)$ , and  $\sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_{\text{i}}^{[l]}(0))$  has distribution

$$\sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_{\text{i}}^{[l]}(0)) \sim \sigma(\|\mathbf{x}_{\text{i}}^{[l]}(0)\|_2 \mathbf{Y}),$$

with  $\mathbf{Y}$  being a standard normal Gaussian vector, then

$$\mathbb{E} \left[ \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_{\text{i}}^{[l]}(0), \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_{\text{i}}^{[l]}(0)) \right\rangle \right] = \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_{\text{i}}^{[l]}(0), \mathbb{E} [\sigma(\|\mathbf{x}_{\text{i}}^{[l]}(0)\|_2 \mathbf{Y})] \right\rangle.$$

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

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

then for any  $\mathbf{Y}_1, \mathbf{Y}_2 \in \mathbb{R}^m$ :

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

Based on our induction hypothesis, with high probability, there exists a uniform constant  $C > 0$ , such that for  $1, 2, \dots, l$ ,  $\|\mathbf{x}_i^{[l]}(0)\|_2 \leq C$ . Hence  $g^{[l]}(\cdot)$  is  $\frac{C}{\sqrt{m}}$ -Lipschitz, apply Lemma B.1.1 once again, we have

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[l]}(0), \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)) \rangle - \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[l]}(0), \mathbb{E} [\sigma(\|\mathbf{x}_i^{[l]}(0)\|_2 \mathbf{Y})] \rangle \right| \geq t \right) \\ & \leq \exp(-cmt^2). \end{aligned} \quad (\text{B.8})$$

From our induction hypothesis

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

then differ by a multiplication of constant, the inequality

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[l]}(0), \mathbb{E} [\sigma(\|\mathbf{x}_i^{[l]}(0)\|_2 \mathbf{Y})] \rangle - \tilde{\mathbf{b}}_i^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} [\sigma(\|\mathbf{x}_i^{[l]}(0)\|_2 u)] \right| \geq t \right) \\ & \leq \exp(-cmt^2) \end{aligned} \quad (\text{B.9})$$

holds. From Lemma B.1.2,

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

naturally, we have

$$\begin{aligned} & \left| \tilde{\mathbf{b}}_i^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} [\sigma(\|\mathbf{x}_i^{[l]}(0)\|_2 u)] - \tilde{\mathbf{b}}_i^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} [\sigma(\sqrt{\widehat{\mathbf{K}}_{\text{ii}}^{[l]}} u)] \right| \\ & \leq C \left| \tilde{\mathbf{b}}_i^{[l]} \right| \left| \|\mathbf{x}_i^{[l]}(0)\|_2 - \sqrt{\widehat{\mathbf{K}}_{\text{ii}}^{[l]}} \right|. \end{aligned} \quad (\text{B.10})$$

Once again from our induction hypothesis,

$$\mathbb{P} \left( \left| \|\mathbf{x}_i^{[l]}(0)\|_2 - \sqrt{\widehat{\mathbf{K}}_{\text{ii}}^{[l]}} \right| \geq t \right) \leq \exp(-cmt^2)$$

holds, then from Equation (B.10), we obtain that

$$\begin{aligned} & \mathbb{P} \left( \left| \tilde{\mathbf{b}}_{\mathbf{i}}^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 u \right) \right] - \tilde{\mathbf{b}}_{\mathbf{i}}^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\mathbf{ii}}^{[l]}} u \right) \right] \right| \geq t \right) \\ & \leq \mathbb{P} \left( \left| C \left| \tilde{\mathbf{b}}_{\mathbf{i}}^{[l]} \right| \left| \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 - \sqrt{\widetilde{\mathbf{K}}_{\mathbf{ii}}^{[l]}} \right| \right| \geq t \right) \leq \exp(-cmt^2). \end{aligned} \quad (\text{B.11})$$

Combine altogether Equation (B.8), Equation (B.9) and Equation (B.11),

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_{\mathbf{i}}^{[l]}(0), \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right) \right\rangle - \tilde{\mathbf{b}}_{\mathbf{i}}^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\mathbf{ii}}^{[l]}} u \right) \right] \right| \geq t \right) \\ & \leq \mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_{\mathbf{i}}^{[l]}(0), \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right) \right\rangle - \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_{\mathbf{i}}^{[l]}(0), \mathbb{E} \left[ \sigma \left( \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right] \right\rangle \right| \geq \frac{t}{3} \right) \\ & + \mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_{\mathbf{i}}^{[l]}(0), \mathbb{E} \left[ \sigma \left( \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 \mathbf{Y} \right) \right] \right\rangle - \tilde{\mathbf{b}}_{\mathbf{i}}^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 u \right) \right] \right| \geq \frac{t}{3} \right) \\ & + \mathbb{P} \left( \left| \tilde{\mathbf{b}}_{\mathbf{i}}^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 u \right) \right] - \tilde{\mathbf{b}}_{\mathbf{i}}^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\mathbf{ii}}^{[l]}} u \right) \right] \right| \geq \frac{t}{3} \right) \\ & \leq \exp(-cmt^2). \end{aligned} \quad (\text{B.12})$$

Secondly for term II, notice that  $\sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right)$  has distribution

$$\sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right) \sim \sigma \left( \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 \mathbf{Y} \right),$$

with  $\mathbf{Y}$  being a standard normal Gaussian vector, then

$$\mathbb{E} \left[ \frac{1}{m} \left\langle \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right), \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right) \right\rangle \right] = \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 u \right)^2 \right],$$

we recall once again that  $\left\{ \sigma \left( (\mathbf{W}_0^{[1]} \mathbf{x}_{\mathbf{i}})_j \right) \right\}_{j=1}^m$  is a collection of i.i.d. sub-Gaussian variables,  $\left\{ \sigma \left( (\mathbf{W}_0^{[1]} \mathbf{x}_{\mathbf{i}})_j \right)^2 \right\}_{j=1}^m$  sub-exponential. Consequently, for all  $t > 0$ ,

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{m} \left\langle \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right), \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right) \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_{\mathbf{i}}^{[l]}(0) \right\|_2 u \right)^2 \right] \right| \geq t \right) \\ & \leq \exp(-cmt^2) \end{aligned} \quad (\text{B.13})$$

holds. From Lemma B.1.3,

$$\begin{aligned} & \left| \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 u \right)^2 \right] - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} u \right)^2 \right] \right| \\ & \leq C \left| \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 - \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \right|. \end{aligned} \quad (\text{B.14})$$

Based on our induction hypothesis,

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

holds, then from Equation (B.14), we obtain that

$$\begin{aligned} & \mathbb{P} \left( \left| \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 u \right)^2 \right] - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} u \right)^2 \right] \right| \geq t \right) \\ & \leq \mathbb{P} \left( \left| C \left| \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 - \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \right| \right| \geq t \right) \leq \exp(-cmt^2). \end{aligned} \quad (\text{B.15})$$

Combine altogether Equation (B.13) and Equation (B.15), since

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{m} \left\langle \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right), \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right) \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} u \right)^2 \right] \right| \geq t \right) \\ & \leq \mathbb{P} \left( \left| \frac{1}{m} \left\langle \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right), \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right) \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 u \right)^2 \right] \right| \geq \frac{t}{2} \right) \\ & + \mathbb{P} \left( \left| \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 u \right)^2 \right] - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} u \right)^2 \right] \right| \geq \frac{t}{2} \right), \end{aligned}$$

then

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{m} \left\langle \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right), \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right) \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} u \right)^2 \right] \right| \geq t \right) \\ & \leq \exp(-cmt^2). \end{aligned} \quad (\text{B.16})$$

We recall once again that

$$\begin{aligned} \|\mathbf{x}_i^{[l+1]}(0)\|_2^2 &= \|\mathbf{x}_i^{[l]}(0)\|_2^2 + \underbrace{\frac{2c_{\text{res}}}{L} \left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)) \right\rangle}_{\text{I}} \\ &\quad + \underbrace{\frac{c_{\text{res}}^2}{L^2} \frac{1}{m} \left\langle \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)), \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)) \right\rangle}_{\text{II}}, \end{aligned}$$

where term I is close to  $\tilde{\mathbf{b}}_i^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\text{ii}}^{[l]}} u \right) \right]$ , as is shown in Equation (B.12)

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \left\langle \mathbf{x}_i^{[l]}(0), \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)) \right\rangle - \tilde{\mathbf{b}}_i^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\text{ii}}^{[l]}} u \right) \right] \right| \geq t \right) \\ &\leq \exp(-cmt^2), \end{aligned} \tag{B.17}$$

and term II is close to  $\mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\text{ii}}^{[l]}} u \right)^2 \right]$ , as is shown in Equation (B.16)

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{m} \left\langle \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)), \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)) \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\text{ii}}^{[l]}} u \right)^2 \right] \right| \geq t \right) \\ &\leq \exp(-cmt^2), \end{aligned} \tag{B.18}$$

then

$$\begin{aligned} &\mathbb{P} \left( \left| \|\mathbf{x}_i^{[l+1]}(0)\|_2^2 - \widetilde{\mathbf{K}}_{\text{ii}}^{[l+1]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right)^2 \right) \\ &\leq \mathbb{P} \left( \left| \|\mathbf{x}_i^{[l]}(0)\|_2^2 - \widetilde{\mathbf{K}}_{\text{ii}}^{[l]} \right| \geq t \right) \\ &\quad + \mathbb{P} \left( \left| 2 \frac{c_{\text{res}}}{L} \left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)) \right\rangle - \tilde{\mathbf{b}}_i^{[l]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\text{ii}}^{[l]}} u \right) \right] \right| \geq 2 \frac{c_{\text{res}}}{L} t \right) \\ &\quad + \mathbb{P} \left( \left| \frac{c_{\text{res}}^2}{L^2} \left| \frac{1}{m} \left\langle \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)), \sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0)) \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\text{ii}}^{[l]}} u \right)^2 \right] \right| \geq \frac{c_{\text{res}}^2}{L^2} t \right) \\ &\leq \exp(-cmt^2). \end{aligned}$$



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

$$\mathbb{P} \left( \left| \left\| \mathbf{x}_i^{[l+1]}(0) \right\|_2^2 - \widetilde{\mathbf{K}}_{ii}^{[l+1]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right)^2 \right) \leq \exp(-cmt^2),$$

if we choose  $c_{\text{res}}$  and  $t$  smartly, set  $t^* = t \left( 1 + \frac{c_{\text{res}}}{L} \right)^2$ , then  $t = t^* \left( 1 + \frac{c_{\text{res}}}{L} \right)^{-2}$ , there exists a uniform constant  $c^*$ , such that

$$1/c^* \leq \left( 1 + \frac{c_{\text{res}}}{L} \right)^{-2L} \leq c^*,$$

then

$$\mathbb{P} \left( \left| \left\| \mathbf{x}_i^{[l+1]}(0) \right\|_2^2 - \widetilde{\mathbf{K}}_{ii}^{[l+1]} \right| \geq t^* \right) \leq \exp(-cmt^{*2}), \quad (\text{B.19})$$

which finishes the proof of Equation (B.6).

Finally, we need to prove Equation (B.7) for  $\tilde{\mathbf{b}}_i^{[l+1]}$ . The structure of outputs of each layer at initial state shall be recalled once again,

$$\mathbf{x}_i^{[l+1]}(0) = \mathbf{x}_i^{[l]}(0) + \frac{c_{\text{res}}}{L\sqrt{m}} \sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right),$$

and the counterpart of  $\tilde{\mathbf{b}}_i^{[l]}$

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

Apply Lemma B.1.1,

$$\mathbb{P} \left( \left| \left\langle \frac{\sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right)}{m}, \mathbf{1} \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} [\sigma(\|\mathbf{x}_i^{[l]}(0)\|_2 u)] \right| \geq t \right) \leq \exp(-cmt^2). \quad (\text{B.20})$$

Apply Lemma B.1.3,

$$\begin{aligned} & \left| \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 u \right)^2 \right] - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} u \right)^2 \right] \right| \\ & \leq C \left| \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 - \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \right|, \end{aligned} \quad (\text{B.21})$$

from induction hypothesis, since

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

holds, then from Equation (B.21), we obtain that

$$\begin{aligned} & \mathbb{P} \left( \left| \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 u \right)^2 \right] - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} u \right)^2 \right] \right| \geq t \right) \\ & \leq \mathbb{P} \left( \left| C \left| \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 - \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \right| \right| \geq t \right) \leq \exp(-cmt^2). \end{aligned} \quad (\text{B.22})$$

combine Equation (B.20) and Equation (B.22),

$$\begin{aligned} & \mathbb{P} \left( \left| \left\langle \frac{\sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right)}{m}, \mathbf{1} \right\rangle - \mathbb{E} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} \mathbf{Y} \right) \right] \right| \geq t \right) \\ & \leq \mathbb{P} \left( \left| \left\langle \frac{\sigma \left( \mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0) \right)}{m}, \mathbf{1} \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 u \right) \right] \right| \geq \frac{t}{2} \right) \\ & + \mathbb{P} \left( \left| \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \left\| \mathbf{x}_i^{[l]}(0) \right\|_2 u \right)^2 \right] - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[l]}} u \right)^2 \right] \right| \geq \frac{t}{2} \right) \\ & \leq \exp(-cmt^2), \end{aligned} \quad (\text{B.23})$$

then

$$\begin{aligned}
& \mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l+1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \tilde{\mathbf{b}}_i^{[l+1]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right) \right) \\
& \leq \mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \tilde{\mathbf{b}}_i^{[l]} \right| \geq t \right) \\
& + \mathbb{P} \left( \frac{c_{\text{res}}}{L} \left| \left\langle \frac{\sigma(\mathbf{W}_0^{[l+1]} \mathbf{x}_i^{[l]}(0))}{m}, \mathbf{1} \right\rangle - \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{\text{ii}}^{[l]}} u \right)^2 \right] \right| \geq \frac{c_{\text{res}}}{L} t \right) \\
& \leq \exp(-cmt^2).
\end{aligned} \tag{B.24}$$

We shall see once again that thanks to the  $\frac{c_{\text{res}}}{L}$  structure, with high probability, the difference  $\left| \left\langle \frac{\mathbf{x}_i^{[l]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \tilde{\mathbf{b}}_i^{[l]} \right|$  only has linear increment with respect to the number of layer. If we choose  $c_{\text{res}}$  and  $t$  smartly, set  $t^* = t \left( 1 + \frac{c_{\text{res}}}{L} \right)$ , then  $t = t^* \left( 1 + \frac{c_{\text{res}}}{L} \right)^{-1}$ , there exists a uniform constant  $c^*$ , such that

$$1/c^* \leq \left( 1 + \frac{c_{\text{res}}}{L} \right)^{-L} \leq c^*,$$

then

$$\mathbb{P} \left( \left| \left\langle \frac{\mathbf{x}_i^{[l+1]}(0)}{\sqrt{m}}, \mathbf{1} \right\rangle - \tilde{\mathbf{b}}_i^{[l+1]} \right| \geq t^* \right) \leq \exp(-cmt^{*2}), \tag{B.25}$$

which finishes the proof of Equation (B.7).  $\square$

### B.3 Analysis on Random Initialization

Our next proposition is on the least eigenvalue of the randomly initialized Gram matrix. First, we shall denote the series of randomly initialized Gram matrices by  $\{\mathbf{G}^{[l]}(0)\}_{l=1}^{L+1}$ , whose components read:

$$\mathbf{G}_{ij}^{[l]}(0) = \left\langle \mathbf{x}_i^{[1]}(0), \mathbf{x}_j^{[1]}(0) \right\rangle. \tag{B.26}$$

**Proposition B.3.1.** *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$ ,  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ , and  $\{\tilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , set  $\lambda_0 > 0$  as*

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

there exists a small constant  $\varepsilon > 0$ , such that with high probability w.r.t random initialization, for  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\varepsilon}\right)$ , then

$$\lambda_{\min}\left(\mathbf{G}^{[1]}(0)\right) \geq \frac{3\lambda_0}{4}. \quad (\text{B.27})$$

*Proof.* We have that

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

Now we need to apply Lemma B.1.1 again, except that we are going to apply it to the inner product function  $h^{[1]}(\cdot) : \mathbb{R}^{2m} \rightarrow \mathbb{R}$ ,

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

where  $\mathbf{Z}^\top = (\mathbf{X}^\top, \mathbf{Y}^\top)$ ,  $\mathbf{X}, \mathbf{Y}$  are standard normal Gaussian vectors, and  $-1 \leq \rho \leq 1$ .

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 B.2.1, with high probability, there exists a uniform constant  $C > 0$ , such that

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

So we have

$$|h^{[1]}(\mathbf{Z}_1) - h^{[1]}(\mathbf{Z}_2)| \leq 4C \sqrt{\frac{c_\sigma}{m}} \|\mathbf{Z}_1 - \mathbf{Z}_2\|_2,$$

hence  $h^{[1]}(\mathbf{Z})$  is  $4C\sqrt{\frac{c_\sigma}{m}}$ -Lipschitz. Then set  $\rho = \widetilde{\mathbf{K}}_{ij}^{[0]}$ , for all  $1 \leq i, j \leq m$ ,

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

Note that based on Proposition A.2.1,  $\lambda_{\min}(\widetilde{\mathbf{K}}^{[1]}) \geq \lambda_0$ , also we have

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

then if we choose  $t = \frac{\lambda_0}{4n}$  and with a union of  $m^2$  such events, we have with probability  $1 - m^2 \exp(-cm\lambda_0^2/n^2)$ ,

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

holds. Hence if there exists a small constant  $\varepsilon > 0$ , such that  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\varepsilon}\right)$ , then with probability  $1 - \exp(-m^\varepsilon)$ ,

$$\lambda_{\min}(\mathbf{G}^{[1]}(0)) \geq \lambda_{\min}(\widetilde{\mathbf{K}}^{[1]}) - \left\| \mathbf{G}^{[1]}(0) - \widetilde{\mathbf{K}}^{[1]} \right\|_{2 \rightarrow 2} \geq \frac{3\lambda_0}{4}. \quad (\text{B.28})$$

□

Our next proposition is on the least eigenvalue of other randomly initialized Gram matrices  $\mathbf{G}^{[l]}(0)$ ,  $l \neq 1$ .

**Proposition B.3.2.** *Given  $\mathcal{X}$ ,  $\sigma(\cdot)$ ,  $\{\widetilde{\mathbf{K}}^{[l]}\}_{l=1}^L$ , and  $\{\tilde{\mathbf{b}}^{[l]}\}_{l=1}^L$ , set  $\lambda_0 > 0$  as*

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

*there exists a small constant  $\varepsilon > 0$ , such that with high probability w.r.t random initialization, for  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\varepsilon}\right)$ , then*

$$\lambda_{\min}(\mathbf{G}^{[l]}(0)) \geq \frac{3\lambda_0}{4}, \quad 2 \leq l \leq L. \quad (\text{B.29})$$

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

$$\begin{aligned} \mathbf{G}_{ij}^{[2]}(0) &= \langle \mathbf{x}_i^{[2]}(0), \mathbf{x}_j^{[2]}(0) \rangle = \mathbf{G}_{ij}^{[1]}(0) + \underbrace{\frac{c_{\text{res}}}{L} \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle}_{\text{I}} \\ &\quad + \underbrace{\frac{c_{\text{res}}}{L} \frac{1}{\sqrt{m}} \langle \mathbf{x}_j^{[1]}(0), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_i^{[1]}(0)) \rangle}_{\text{II}} + \underbrace{\frac{c_{\text{res}}^2}{L^2} \frac{1}{m} \langle \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_i^{[1]}(0)), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle}_{\text{III}} \\ \tilde{\mathbf{b}}_i^{[1]} &= \sqrt{c_\sigma} \mathbb{E}_{u \sim \mathcal{N}(0, \tilde{\mathbf{K}}_{ii}^{[0]})} [\sigma(u)], \\ \widetilde{\mathbf{K}}_{ij}^{[2]} &= \widetilde{\mathbf{K}}_{ij}^{[1]} + \mathbb{E}_{(u,v)^\top \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[2]})} \left[ \underbrace{\frac{c_{\text{res}}}{L} \tilde{\mathbf{b}}_i^{[1]} \sigma(v)}_{\text{I}'} + \underbrace{\frac{c_{\text{res}}}{L} \tilde{\mathbf{b}}_j^{[1]} \sigma(u)}_{\text{II}'} + \underbrace{\frac{c_{\text{res}}^2}{L^2} \sigma(u) \sigma(v)}_{\text{III}'} \right]. \end{aligned}$$

We need to tackle the difference between I, II, III and I', II', III'. For I and I', we need to write the difference into:

$$\begin{aligned} &\left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle - \tilde{\mathbf{b}}_i^{[1]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{jj}^{[1]}} u \right) \right] \right| \\ &\leq \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle - \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \mathbb{E} [\sigma(\|\mathbf{x}_j^{[1]}(0)\|_2 \mathbf{Y})] \rangle \right| \\ &\quad + \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \mathbb{E} [\sigma(\|\mathbf{x}_j^{[1]}(0)\|_2 \mathbf{Y})] \rangle - \tilde{\mathbf{b}}_i^{[1]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} [\sigma(\|\mathbf{x}_j^{[1]}(0)\|_2 u)] \right| \\ &\quad + \left| \tilde{\mathbf{b}}_i^{[1]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} [\sigma(\|\mathbf{x}_j^{[1]}(0)\|_2 u)] - \tilde{\mathbf{b}}_i^{[1]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{jj}^{[1]}} u \right) \right] \right|, \end{aligned}$$

similar to the proof in Proposition B.2.1, with  $\mathbf{Y}$  being a standard normal Gaussian vector,

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle - \tilde{\mathbf{b}}_i^{[1]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{jj}^{[1]}} u \right) \right] \right| \geq t \right) \\ &\leq \exp(-cmt^2). \end{aligned} \tag{B.30}$$

By symmetry, for II and II',

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_j^{[1]}(0), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_i^{[1]}(0)) \rangle - \tilde{\mathbf{b}}_j^{[1]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\widetilde{\mathbf{K}}_{ii}^{[1]}} u \right) \right] \right| \geq t \right) \\ &\leq \exp(-cmt^2). \end{aligned} \tag{B.31}$$

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

$$h^{[2]}(\mathbf{Z}) = \frac{1}{m} \left\langle \sigma \left( \left\| \mathbf{x}_i^{[1]}(0) \right\|_2 \mathbf{X} \right), \sigma \left( \left( \left\| \mathbf{x}_i^{[1]}(0) \right\|_2 \left( \rho \mathbf{X} + \sqrt{1 - \rho^2} \mathbf{Y} \right) \right) \right) \right\rangle,$$

where  $\mathbf{Z}^\top = (\mathbf{X}^\top, \mathbf{Y}^\top)$ ,  $\mathbf{X}, \mathbf{Y}$  are standard normal Gaussian vectors, and  $-1 \leq \rho \leq 1$ .

For any  $\mathbf{Z}_1, \mathbf{Z}_2 \in \mathbb{R}^m$ ,

$$\begin{aligned} |h^{[2]}(\mathbf{Z}_1) - h^{[2]}(\mathbf{Z}_2)| &\leq \frac{1}{\sqrt{m}} \left\| \sigma \left( \left( \left\| \mathbf{x}_i^{[1]}(0) \right\|_2 \left( \rho \mathbf{X}_1 + \sqrt{1 - \rho^2} \mathbf{Y}_1 \right) \right) \right) \right\|_2 \frac{1}{\sqrt{m}} \|\mathbf{X}_1 - \mathbf{X}_2\|_2 \\ &\quad + \frac{1}{\sqrt{m}} \left\| \sigma \left( \left\| \mathbf{x}_i^{[1]}(0) \right\|_2 \mathbf{X}_2 \right) \right\|_2 \frac{\left\| \mathbf{x}_i^{[1]}(0) \right\|_2}{\sqrt{m}} \left( |\rho| \|\mathbf{X}_1 - \mathbf{X}_2\|_2 + \sqrt{1 - \rho^2} \|\mathbf{Y}_1 - \mathbf{Y}_2\|_2 \right), \end{aligned}$$

combined with Proposition B.2.1, with high probability, there exists a uniform constant  $C > 0$ , such that

$$\left\| \mathbf{x}_i^{[1]}(0) \right\|_2 \leq C,$$

consequently,

$$\begin{aligned} \frac{1}{\sqrt{m}} \left\| \sigma \left( \left( \left\| \mathbf{x}_i^{[1]}(0) \right\|_2 \left( \rho \mathbf{X}_1 + \sqrt{1 - \rho^2} \mathbf{Y}_1 \right) \right) \right) \right\|_2 &\leq 10C, \\ \frac{1}{\sqrt{m}} \left\| \sigma \left( \left\| \mathbf{x}_i^{[1]}(0) \right\|_2 \mathbf{X}_2 \right) \right\|_2 &\leq 10C, \end{aligned}$$

with abuse of notations,  $h^{[2]}(\cdot)$  is  $\frac{C}{\sqrt{m}}$ -Lipschitz, then

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

hence we have

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{m} \left\langle \sigma \left( \mathbf{W}_0^{[2]} \mathbf{x}_i^{[1]}(0) \right), \sigma \left( \mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0) \right) \right\rangle - \mathbb{E}_{(u,v)^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_{ij}^{[2]})} [\sigma(u)\sigma(v)] \right| \geq t \right) \\ &\leq \exp(-cmt^2), \end{aligned} \tag{B.32}$$

with

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

Combine Lemma B.1.3 and Proposition B.2.1

$$\begin{aligned} & \mathbb{P} \left( \left| \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_{ij}^{[2]})} [\sigma(u)\sigma(v)] - \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[2]})} [\sigma(u)\sigma(v)] \right| \geq t \right) \\ & \leq \exp(-cmt^2), \end{aligned} \quad (\text{B.33})$$

consequently, we have

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{m} \langle \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_i^{[1]}(0)), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle - \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[2]})} [\sigma(u)\sigma(v)] \right| \geq t \right) \\ & \leq \mathbb{P} \left( \left| \frac{1}{m} \langle \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_i^{[1]}(0)), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle - \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_{ij}^{[2]})} [\sigma(u)\sigma(v)] \right| \geq \frac{t}{2} \right) \\ & + \mathbb{P} \left( \left| \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_{ij}^{[2]})} [\sigma(u)\sigma(v)] - \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[2]})} [\sigma(u)\sigma(v)] \right| \geq \frac{t}{2} \right) \\ & \leq \exp(-cmt^2). \end{aligned} \quad (\text{B.34})$$

To sum up, we have

$$\begin{aligned} & \mathbb{P} \left( \left| \mathbf{G}_{ij}^{[2]}(0) - \tilde{\mathbf{K}}_{ij}^{[2]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right)^2 \right) \\ & \leq \mathbb{P} \left( \left| \mathbf{G}_{ij}^{[1]}(0) - \tilde{\mathbf{K}}_{ij}^{[1]} \right| \geq t \right) \\ & + \mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_i^{[1]}(0), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle - \tilde{\mathbf{b}}_i^{[1]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\tilde{\mathbf{K}}_{jj}^{[1]}} u \right) \right] \right| \geq t \right) \\ & + \mathbb{P} \left( \left| \frac{1}{\sqrt{m}} \langle \mathbf{x}_j^{[1]}(0), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_i^{[1]}(0)) \rangle - \tilde{\mathbf{b}}_j^{[1]} \mathbb{E}_{u \sim \mathcal{N}(0,1)} \left[ \sigma \left( \sqrt{\tilde{\mathbf{K}}_{ii}^{[1]}} u \right) \right] \right| \geq t \right) \\ & + \mathbb{P} \left( \left| \frac{1}{m} \langle \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_i^{[1]}(0)), \sigma(\mathbf{W}_0^{[2]} \mathbf{x}_j^{[1]}(0)) \rangle - \mathbb{E}_{(u,v) \sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{A}}_{ij}^{[2]})} [\sigma(u)\sigma(v)] \right| \geq t \right) \end{aligned} \quad (\text{B.35})$$



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

$$\mathbb{P} \left( \left| \mathbf{G}_{ij}^{[l]}(0) - \widetilde{\mathbf{K}}_{ij}^{[l]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right)^{2l-2} \right) \leq \exp(-cmt^2), \quad (\text{B.36})$$

specifically, for  $l = L$ ,

$$\mathbb{P} \left( \left| \mathbf{G}_{ij}^{[L+1]}(0) - \mathbf{K}_{ij}^{[L+1]} \right| \geq t \left( 1 + \frac{c_{\text{res}}}{L} \right)^{2L} \right) \leq \exp(-cmt^2). \quad (\text{B.37})$$

(ii). We recall the matrix inequality once again,

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

based on Proposition A.2.1,  $\lambda_{\min}(\mathbf{K}^{[L+1]}) > \lambda_0$ . Then if we choose  $t = \frac{\lambda_0}{4n \exp(2c_{\text{res}})}$ , then for all  $l$ , with probability  $1 - \exp(-cm\lambda_0^2/n^2)$ ,

$$\left\| \mathbf{G}^{[l]}(0) - \widetilde{\mathbf{K}}^{[l]} \right\|_{2 \rightarrow 2} \leq \frac{\lambda_0}{4} \quad (\text{B.38})$$

holds. Hence if there exists a small constant  $\varepsilon > 0$ , such that  $m = \Omega \left( \left( \frac{n}{\lambda_0} \right)^{2+\varepsilon} \right)$ , then with probability  $1 - \exp(-m^\varepsilon)$ ,

$$\lambda_{\min}(\mathbf{G}^{[l]}(0)) \geq \lambda_{\min}(\widetilde{\mathbf{K}}^{[l]}) - \left\| \mathbf{G}^{[l]}(0) - \widetilde{\mathbf{K}}^{[l]} \right\|_{2 \rightarrow 2} > \frac{3\lambda_0}{4}. \quad (\text{B.39})$$

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

$$\left\| \mathbf{G}^{[L+1]}(0) - \mathbf{K}^{[L+1]} \right\|_{2 \rightarrow 2} \leq \frac{\lambda_0}{4}, \quad (\text{B.40})$$

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

$$\lambda_{\min}(\mathbf{G}^{[L+1]}(0)) \geq \lambda_{\min}(\mathbf{K}^{[L+1]}) - \left\| \mathbf{G}^{[L+1]}(0) - \mathbf{K}^{[L+1]} \right\|_{2 \rightarrow 2} > \frac{3\lambda_0}{4}. \quad (\text{B.41})$$

□

We remark that

$$\left[\mathcal{G}_0^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)\right]_{1 \leq \alpha, \beta \leq n} = \mathbf{G}^{[L+1]}(0),$$

then directly from Proposition B.3.2, for  $m = \Omega\left(\left(\frac{n}{\lambda_0}\right)^{2+\varepsilon}\right)$ ,

$$\lambda_{\min}\left(\left[\mathcal{G}_0^{[L+1]}(\mathbf{x}_\alpha, \mathbf{x}_\beta)\right]_{1 \leq \alpha, \beta \leq n}\right) \geq \frac{3\lambda_0}{4} \quad (\text{B.42})$$

holds with high probability w.r.t random initialization.