

# Global Convergence Rate of Deep Equilibrium Models with General Activations

Anonymous authors

Paper under double-blind review

## Abstract

In a recent paper, Ling et al. investigated the over-parametrized Deep Equilibrium Model (DEQ) with ReLU activation. They proved that the gradient descent converges to a globally optimal solution at a linear convergence rate for the quadratic loss function. This paper shows that this fact still holds for DEQs with any general activation that has bounded first and second derivatives. Since the new activation function is generally non-homogeneous, bounding the least eigenvalue of the Gram matrix of the equilibrium point is particularly challenging. To accomplish this task, we need to create a novel population Gram matrix and develop a new form of dual activation with Hermite polynomial expansion.

## 1 Introduction

Deep learning is a class of machine learning algorithms that uses multiple layers to progressively extract higher-level features from the raw input. For example, in image processing, lower layers may identify edges, while higher layers may identify the concepts relevant to a human such as digits or letters or faces. Deep neural networks have underpinned state of the art empirical results in numerous applied machine learning tasks (Krizhevsky et al., 2012). Understanding neural network learning, particularly its recent successes, commonly decomposes into the two main themes: (i) studying generalization capacity of the deep neural networks and (ii) understanding why efficient algorithms, such as stochastic gradient, find good weights. Though still far from being complete, previous work provides some understanding on generalization capability of deep neural networks. However, question (ii) is rather poorly understood. While learning algorithms succeed in practice, theoretical analysis is overly pessimistic. Direct interpretation of theoretical results suggests that when going slightly deeper beyond single layer networks, e.g. to depth-two networks with very few hidden units, it is hard to predict even marginally better than random (Daniely et al., 2013; Kearns & Valiant, 1994).

The standard approach to develop generalization bounds on deep learning (and machine learning) was developed in seminal papers by (Vapnik, 1998), and it is based on bounding the difference between the generalization error and the training error. These bounds are expressed in terms of the so called VC-dimension of the class. However, these bounds are very loose when the VC-dimension of the class can be very large, or even infinite. In 1998, several authors (Bartlett & Shawe-Taylor, 1999; Bartlett et al., 1998) suggested another class of upper bounds on generalization error that are expressed in terms of the empirical distribution of the margin of the predictor (the classifier). Later, Koltchinskii and Panchenko proposed new probabilistic upper bounds on generalization error of the combination of many complex classifiers such as deep neural networks (Koltchinskii & Panchenko, 2002). These bounds were developed based on the general results of the theory of Gaussian, Rademacher, and empirical processes in terms of general functions of the margins, satisfying a Lipschitz condition. They improved previously known bounds on generalization error of convex combination of classifiers. (Truong, 2022a) and Truong (2022b) have recently provided generalization bounds for learning with Markov dataset based on Rademacher and Gaussian complexity functions. The development of new symmetrization inequalities and contraction lemmas in high-dimensional probability for Markov chains is a key element in these works. Several recent works have focused on gradient descent based PAC-Bayesian algorithms, aiming to minimise a generalisation bound for stochastic classifiers (Biggs & Guedj, 2021; Dziugaite & Roy., 2017). Most of these studies use a surrogate loss to avoid dealing with

the zero-gradient of the misclassification loss. There were some other works which use information-theoretic approach to find PAC-bounds on generalization errors for machine learning (Esposito et al., 2021; Xu & Raginsky, 2017) and deep learning (Jakubovitz et al., 2018).

Recently, deep equilibrium model (DEQ)(Bai et al., 2019) was introduced as a new approach to modelling sequential data. In many many existing deep sequence models, the hidden layers converge toward some fixed points. DEQ directly finds these equilibrium points via root-finding of implicit equations. Such a model is equivalent to an infinite-depth weight-tied model with input-injection. DEQ has emerged as an important model in various applications such as computer vision (Bai et al., 2020; Xie et al., 2022), natural language processing (Bai et al., 2019), and inverse problems (Gilton et al., 2021). This model has been shown to achieve performance competitive with the state-of-the-art deep networks while using significantly less memory. Despite of the empirical success of DEQ, theoretical understanding of this model is still limited. The effectiveness of over-parameterization in optimizing feedforward neural networks has been validated in many research literature (Arora et al., 2019; Du et al., 2018; Li & Liang, 2018). A recent work (Nguyen, 2021) showed that the convergence of gradient descent (GD) to a global optimum can be guaranteed when the width of the last hidden layer exceeds the number of training samples. The main idea is to investigate the property at initialization and bound the traveling distance of GD from the initialization.

However, it remains unknown whether the above results can be directly applied to DEQs. Due to the implicit weight-sharing, the initial random weights and features are dependent, which causes the standard concentration approaches in the existing research literature fail in DEQs. Recently, Ling et al. (2022) investigated the training dynamics of over-parameterized DEQs with ReLU activation. More specifically, they proposed a novel probabilistic framework to overcome the challenge arising from the weight-sharing and the infinite depth. By supposing a condition on the initial equilibrium point, they proved that the gradient descent converges to a globally optimal solution at a linear convergence rate for the quadratic loss function. To achieve this target, they developed a lower bound on the least eigenvalue of the Gram matrix for the DEQs with ReLU activation. One interesting open question is whether the gradient descent algorithm still converge at a linear rate for DEQs with non-linear activation functions? In this paper, we show that this fact still holds for DEQs with a general activation function which has bounded first and second derivatives. Many popular activation functions such as  $1/(1 + e^{-x})$ ,  $\text{erf}(x)$ ,  $x/\sqrt{1+x^2}$ ,  $\sin(x)$ ,  $\tanh(x)$  satisfy the boundedness requirements. In general, the new activation function does not have homogeneous property as ReLU, hence a novel population Gram matrix is designed for DEQs with general activations, and a new form of dual activation with Hermite polynomial expansion is developed in our work.

## 2 Problem settings

We consider the same model as Ling et al. (2022). However, different from Ling et al. (2022), we assume that the activation function,  $\varphi$ , satisfies some constraints in the first and second derivatives. These properties can be observed in many common activation functions. More specifically, we define a vanilla deep equilibrium model (DEQ) with the transform of the  $l$ -th layer as

$$\mathbf{T}^{(l)} = \varphi(\mathbf{W}\mathbf{T}^{(l-1)} + \mathbf{U}\mathbf{X}) \quad (1)$$

where  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  denotes the training inputs,  $\mathbf{U} \in \mathbb{R}^{m \times d}$  and  $\mathbf{W} \in \mathbb{R}^{m \times m}$  are trainable weight matrices, and  $\mathbf{T}^{(l)} \in \mathbb{R}^{m \times n}$  is the output feature at the  $l$ -th hidden layer. The output of the last hidden layer is defined by  $\mathbf{T}^* := \lim_{l \rightarrow \infty} \mathbf{T}^{(l)}$  under the condition that this limit exists uniquely. Therefore, instead of running infinitely deep layer-by-layer forward propagation,  $\mathbf{T}^*$  can be calculated by directly solving the equilibrium point of the following equation

$$\mathbf{T}^* = \varphi(\mathbf{W}\mathbf{T}^* + \mathbf{U}\mathbf{X}). \quad (2)$$

Let  $\mathbf{y} = [y_1, y_2, \dots, y_n] \in \mathbb{R}^n$  denote the labels, and  $\hat{\mathbf{y}}(\boldsymbol{\theta}) = \mathbf{a}^T \mathbf{T}^*$  be the prediction function with  $\mathbf{a} \in \mathbb{R}^m$  being a trainable vector and  $\boldsymbol{\theta} = \text{vec}(\mathbf{W}, \mathbf{U}, \mathbf{a})$ . Our target is to minimize the empirical risk with the quadratic loss function:

$$\Phi(\boldsymbol{\theta}) = \frac{1}{2} \|\hat{\mathbf{y}}(\boldsymbol{\theta}) - \mathbf{y}\|_2^2. \quad (3)$$

To optimize this loss function, we use the gradient descent update  $\boldsymbol{\theta}(\tau+1) = \boldsymbol{\theta}(\tau) - \eta \nabla \Phi(\boldsymbol{\theta}(\tau))$ , where  $\eta$  is the learning rate and  $\boldsymbol{\theta}(\tau) = \text{vec}(\mathbf{W}(\tau), \mathbf{U}(\tau), \mathbf{a}(\tau))$ . For notational simplicity, we omit the superscript and denote  $\mathbf{T}$  to be the equilibrium  $\mathbf{T}^*$  when it is clear from the context. Moreover, the Gram matrix of the equilibrium point is defined by  $\mathbf{G}(\tau) := \mathbf{T}^T(\tau)\mathbf{T}(\tau)$  and we define its least eigenvalue by  $\lambda_\tau = \lambda_{\min}(\mathbf{G}(\tau))$ . In this paper, for brevity we denote by  $\mathbf{G} = \mathbf{G}(0)$ .

**Definition 1.** An activation  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  is  $L$ -bounded if it is twice continuously differentiable and  $\|\varphi\|_\infty, \|\varphi'\|_\infty, \|\varphi''\|_\infty \leq L$ .

In this paper, we assume that  $\varphi(\cdot)$  is  $L$ -bounded. In addition, the following holds:

$$q := \sqrt{\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \varphi^2(z) \exp\left(-\frac{z^2}{2}\right) dz} > 0. \quad (4)$$

Many popular activation functions such as  $1/(1+e^{-x})$ ,  $\text{erf}(x)$ ,  $x/\sqrt{1+x^2}$ ,  $\sin(x)$ ,  $\tanh(x)$  satisfy the boundedness requirements.

Besides, we use similar assumptions on the random initialization and input data as Ling et al. (2022):

- **Assumption 1** (Random initialization). Assume that  $\sigma_w^2 < \frac{1}{48L^2}$ . In addition,  $\mathbf{W}$  is initialized with an  $m \times m$  matrix with i.i.d. entries  $\mathbf{W}_{ij} \sim \mathcal{N}(0, 2\sigma_w^2/m)$ ,  $\mathbf{U}$  is initialized with an  $m \times d$  matrix with i.i.d. entries  $\mathbf{U}_{ij} \sim \mathcal{N}(0, 2/d)$ , and  $\mathbf{a}$  is initialized with a random vector with i.i.d. entries  $\sim \mathcal{N}(0, 1/m)$ .
- **Assumption 2** (Input data). We assume that (i)  $\|\mathbf{x}_i\|_2 = \sqrt{d}$  for all  $i \in [n]$  and  $\mathbf{x}_i \not\parallel \mathbf{x}_j$  for all  $i \neq j$ ; (ii) the labels satisfy  $|y_i| = O(1)$  for all  $i \in [n]$ .

### 3 Main Results

In this paper, we show that if the learning rate is small enough, the loss converges to a global minimum at linear rate. The result is as follows.

**Theorem 2.** Consider a DEQ. Let  $\delta$  be a constant such that  $\|\mathbf{W}(0)\| + \delta < 1$ . Denote by  $\bar{\rho}_w = \|\mathbf{W}(0)\|_2 + \delta$ ,  $\bar{\rho}_u = \|\mathbf{U}(0)\|_2 + \delta$ ,  $\bar{\rho}_a = \|\mathbf{a}(0)\|_2 + \delta$  and define

$$c_a = \frac{L\bar{\rho}_u}{1 - L\bar{\rho}_w}, \quad c_u = \frac{L\bar{\rho}_a}{1 - L\bar{\rho}_w}, \quad c_m = \frac{m^2\sigma(0)}{1 - L\bar{\rho}_w}. \quad (5)$$

In addition, assume at initialization that

$$\lambda_0 \geq \frac{4}{\delta} \max \left\{ c_u(c_a\|\mathbf{X}\|_F + c_m), c_u\|\mathbf{X}\|_F, c_a\|\mathbf{X}\|_F + c_m \right\} \|\hat{\mathbf{y}}(0) - \mathbf{y}\|, \quad (6)$$

$$\lambda_0^{3/2} \geq \frac{4(2+\sqrt{2})L}{(1-L\bar{\rho}_w)\lambda_0} \left[ (c_a\|\mathbf{X}\|_F + c_m)^2 + c_u\|\mathbf{X}\|_F^2 \right] \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2, \quad (7)$$

$$\lambda_0 \geq 8c_u^2(c_a\|\mathbf{X}\|_F + c_m)^2 + c_u^2\|\mathbf{X}\|_F^2 \quad (8)$$

where  $\lambda_0$  is the least eigenvalue of  $\mathbf{G}(0) = \mathbf{T}(0)^T\mathbf{T}(0)$ . Then, if the learning rate satisfies

$$\eta < \min \left( \frac{2}{\lambda_0}, \frac{2[c_u^2(c_a\|\mathbf{X}\|_F + c_m)^2 + c_u^2\|\mathbf{X}\|_F^2]}{c_u^2(c_a\|\mathbf{X}\|_F + c_m)^2 + c_u^2\|\mathbf{X}\|_F^2 + (c_a\|\mathbf{X}\|_F + c_m)^2} \right), \quad (9)$$

for every  $\tau \geq 0$ , the following hold:

- $\|\mathbf{W}(\tau)\|_2 \leq 1$ , i.e., the equilibrium points always exists,
- $\lambda_\tau > \frac{1}{2}\lambda_0$ , and

$$\|\nabla_\theta \Phi(\boldsymbol{\theta}(\tau))\|_2^2 \geq \lambda_0 \Phi(\boldsymbol{\theta}(\tau)). \quad (10)$$

- The loss converges to a global minimum as

$$\Phi(\theta(\tau)) \leq \left(1 - \eta \frac{\lambda_0}{2}\right)^\tau \Phi(\theta(0)). \quad (11)$$

The main challenge now is to find some initializations such that  $\lambda_0$  satisfies all the conditions in Theorem 2. To lower bound  $\lambda_0$ , we need to design a population Gram matrix  $\mathbf{K}$  and compare  $\lambda_0$  with the least eigenvalue of  $\mathbf{K}$  Ling et al. (2022). However, since the new activation function,  $\varphi$ , is non-linear in general, bounding  $\lambda_0$  is more challenging than the ReLU network in Ling et al. (2022). The non-homogeneity of activation functions causes the techniques to design  $\mathbf{K}$  in (Ling et al., 2022, Definition 1) can not be applied. For example, (Ling et al., 2022, Eq. 11) only holds for ReLU.

In Section 4, we propose a new method to create the population Gram matrix  $\mathbf{K}$  for DEQs with general Lipschitz activation function. By using our new form of dual activation and Hermite polynomial expansion, we can prove that  $\mathbf{K}$  is symmetric positive definite. In addition, we show that with probability at least  $1 - t$ ,  $\lambda_0 \geq \frac{m}{2} \lambda_*$  provided that  $m = \Omega\left(\frac{n}{\lambda_*^2} \log \frac{n}{t}\right)$  where  $\lambda_*$  is the least eigenvalue of  $\mathbf{K}$ . This fact indicates that all the conditions of Theorem 2 at least hold for over-parametrized DEQs (or  $m$  sufficiently large) with  $\varphi(0) = 0$ . Hence, by (11) in Theorem 2, the gradient descent algorithm converges to a global optimum at a linear rate for the over-parametrized DEQs. This fascinating fact is reaffirmed by our numerical experiments on real datasets such as MNIST and CFAR10 in Section 7.

## 4 A novel design of the population Gram matrix $\mathbf{K}$

The key approach in lower bounding  $\lambda_0$  is to design a population Gram matrix  $\mathbf{K}$  in such a way that we can lower bound  $\lambda_0$  by the least eigenvalue of  $\mathbf{K}$  and that  $\mathbf{K}$  is symmetric positive definite. This novel population Gram matrix is developed through our introduction of a new form of dual activation.

First, we define a new class of dual activation functions  $\tilde{Q}_{\alpha,\beta} : [-1, 1] \rightarrow \mathbb{R}$  for all pairs  $(\alpha, \beta) \in \mathbb{R}_+^2$ .

**Definition 3.** Recall the definition of  $q$  in (4). For each pair  $(\alpha, \beta)$ , define

$$\tilde{Q}_{\alpha,\beta}(x) := \frac{1}{\alpha\beta q^2} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} [\varphi(\alpha a) \varphi(\beta b)], \quad \forall |x| \leq 1. \quad (12)$$

If  $\varphi(x) = \max\{x, 0\}$  (ReLU), then  $\tilde{Q}_{\alpha,\beta}(x) = \bar{Q}(x)$  for all  $(\alpha, \beta) \in \mathbb{R}_+^2$ , where

$$\bar{Q}(x) := \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} [\varphi(a) \varphi(b)]$$

is the dual activation defined in (Daniely et al., 2016, Sec. 3.2).

Now, we provide a novel design of the population Gram matrix  $\mathbf{K}$  based on this new dual activation function.

**Definition 4.** Given the training input  $\mathbf{X} := [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$  satisfying Assumption 2. Let

$$Q_{ij}(x) := \tilde{Q}_{\sqrt{2\left(\frac{\sigma_w^2}{m} \mathbf{G}_{ii} + 1\right)}, \sqrt{2\left(\frac{\sigma_w^2}{m} \mathbf{G}_{jj} + 1\right)}}(x), \quad \forall x \in \mathbb{R}. \quad (13)$$

We define the population Gram matrices  $\mathbf{K}^{(l)}$  of each layer recursively as

$$\rho_{ij}^{(0)} = 0, \quad (14)$$

$$\rho_{ii}^{(l)} = 2q^2 \sigma_w^2 \rho_{ii}^{(l-1)} Q_{ii}(1) + 1, \quad (15)$$

$$\rho_{ij}^{(l)} = \sqrt{\rho_{ii}^{(l)} \rho_{jj}^{(l)}}, \quad i \neq j \quad (16)$$

$$\mathbf{K}^{(0)} = 0, \quad (17)$$

$$\nu_{ij}^{(l)} = \frac{\sigma_w^2 \mathbf{K}_{ij}^{(l-1)} + d^{-1} \mathbf{x}_i^T \mathbf{x}_j}{\sqrt{(\sigma_w^2 \mathbf{K}_{ii}^{(l-1)} + 1)(\sigma_w^2 \mathbf{K}_{jj}^{(l-1)} + 1)}} \quad (18)$$

$$\mathbf{K}_{ij}^{(l)} = 2q^2 \rho_{ij}^{(l)} Q_{ij}(\nu_{ij}^{(l)}) \quad (19)$$

for all  $l \geq 1$  and  $i, j \in [n] \times [n]$ .

The next result shows that  $\lambda_0$  can be lower bounded via the least eigenvalue of the population matrix  $\mathbf{K}$ .

**Theorem 5.** *If  $m = \Omega\left(\frac{n^2}{\lambda_*^2} \log \frac{n}{t}\right)$ , with probability at least  $1 - t$ , it holds that*

$$\lambda_0 \geq \frac{m}{2} \lambda_*. \quad (20)$$

Finally, the following result shows sufficient conditions such that  $\mathbf{K}$  is strictly positive definite.

**Theorem 6.** *Assume that there exists a polynomial expansion of  $\tilde{Q}_{\alpha, \alpha}$  satisfying:*

$$\tilde{Q}_{\alpha, \alpha}(x) = \sum_{r=0}^{\infty} \mu_{r, \alpha}^2(\varphi) x^r \quad (21)$$

for all  $\alpha > 0$  such that  $\sup\{r : \mu_{r, \alpha}^2(\varphi) > 0\} = \infty$ . Then,  $\mathbf{K}$  is strictly positive definite (or  $\lambda_* > 0$ ).

## 5 Proof of Theorem 5

To prove Theorem 5, we first state some auxiliary results based on the population Gram matrix  $\mathbf{K}$  in Definition 4. The proofs of these lemmas and prepositions can be found in Supplement Material.

**Lemma 7.** *Recall the definition of  $\tilde{Q}_{\alpha, \beta}$  in Definition 3. Then, the following hold for all  $\alpha \geq 1, \beta \geq 1$  and  $x \in \mathbb{R}$ :*

$$|\tilde{Q}_{\alpha, \beta}(x)| \leq \sqrt{\tilde{Q}_{\alpha, \alpha}(1) \tilde{Q}_{\beta, \beta}(1)}, \quad (22)$$

$$|\tilde{Q}_{\alpha, \beta}(x)| \leq \frac{4L^2}{q^2}, \quad \forall |x| \leq 1. \quad (23)$$

In addition,  $\tilde{Q}_{\alpha, \beta}(\cdot)$  is  $\frac{2L^2}{q^2}$ -Lipchitz for any fixed positive pair  $(\alpha, \beta)$ .

**Lemma 8.** *(Ling et al., 2022, Proof of Lemma 4) For  $l \geq 1$ ,  $\mathbf{G}_{ij}^{(l+1)}$  can be reconstructed as  $\mathbf{G}_{ij}^{(l+1)} = \varphi(\mathbf{M} \mathbf{h}_{l+1})^T \varphi(\mathbf{M} \mathbf{h}_{l'+1})$  such that*

- (i)  $\mathbf{h}_{l+1}^T \mathbf{h}_{l'+1} = \frac{\sigma_w^2}{m} \mathbf{G}_{ij}^{(l)} + \frac{1}{d} \mathbf{x}_i^T \mathbf{x}_j$ ,
- (ii)  $\mathbf{M} \in \mathbb{R}^{m \times (2l+d+2)}$  is a rectangle matrix, and the entries of  $\mathbf{M}$  are i.i.d. from  $\mathcal{N}(0, 2)$  conditioning on previous layers.

**Lemma 9.** *For the given setting, we have*

$$\rho_{ii}^{(l)} = \sigma_w^2 \mathbf{K}_{ii}^{(l-1)} + 1, \quad (24)$$

$$\rho_{ij}^{(l)} \nu_{ij}^{(l)} = \sigma_w^2 \mathbf{K}_{ij}^{(l-1)} + d^{-1} \mathbf{x}_i^T \mathbf{x}_j, \quad \forall i, j, \quad (25)$$

and

$$\nu_{ij}^{(l)} = \begin{cases} \frac{Q_{ij}(\nu_{ij}^{(l-1)}) / \sqrt{Q_{ii}(1) Q_{jj}(1)} \sqrt{(\rho_{ii}^{(l)} - 1)(\rho_{jj}^{(l)} - 1) + d^{-1} \mathbf{x}_i^T \mathbf{x}_j}}{\sqrt{\rho_{ii}^{(l)} \rho_{jj}^{(l)}}}, & i \neq j \\ 1, & i = j \end{cases}. \quad (26)$$

In addition, we also have

$$|\nu_{ij}^{(l)}| \leq 1 \quad (27)$$

for all  $i, j \in [n] \times [n]$  and  $l \geq 0$ .

**Proposition 10.** *Under the Assumptions 1 and 2, with probability at least  $1 - m \exp(-\Omega(m))$ , we have  $\|\mathbf{K} - \mathbf{K}^{(l)}\| = O(n(8L^2\sigma_w^2)^l)$  which implies that, for  $l \rightarrow \infty$ ,  $\mathbf{K}^{(l)} \rightarrow \mathbf{K}$  with entries*

$$\mathbf{K}_{ij} = 2q^2 Q_{ij}(\nu_{ij}) \sqrt{\rho_{ii} \rho_{jj}} \quad (28)$$

where

$$\nu_{ij} = \begin{cases} \frac{Q_{ij}(\nu_{ij}) / \sqrt{Q_{ii}(1)Q_{jj}(1)} \sqrt{(\rho_{ii}-1)(\rho_{jj}-1)} + d^{-1} \mathbf{x}_i^T \mathbf{x}_j}{\sqrt{\rho_{ii} \rho_{jj}}}, & i \neq j \\ 1, & i = j \end{cases}. \quad (29)$$

Here,

$$\rho_{ii} = \frac{1}{1 - 2q^2 \sigma_w^2 Q_{ii}(1)}. \quad (30)$$

**Proposition 11.** *Under Assumptions 1 and 2 with probability at least  $1 - n^2 \exp(-\Omega(m))$ , it holds that*

$$\frac{1}{m} \left\| \mathbf{G} - \mathbf{G}^{(l)} \right\|_F = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (31)$$

**Proposition 12.** *Under Assumptions 1 and 2, with probability at least  $1 - n^2 \exp\{-\Omega(8^l L^{2l} \sigma_w^{2l} mn L^2) + O(l^2)\}$ , it holds that*

$$\left\| \frac{1}{m} \mathbf{G}^{(l)} - \mathbf{K}^{(l)} \right\|_F = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (32)$$

By combining Propositions 10–12, we can bound  $\lambda_0$  via the least eigenvalue of the population matrix  $\mathbf{K}$  as follows.

*Proof of Theorem 5.* From Propositions 10–12, with probability at least  $1 - n^2 \exp(-\Omega(m8^l L^{2l} \sigma_w^{2l}) + O(l^2))$ , it holds that

$$\left\| \frac{1}{m} \mathbf{G} - \mathbf{K} \right\|_F \leq \frac{1}{m} \left\| \mathbf{G} - \mathbf{G}^{(l)} \right\|_F + \left\| \frac{1}{m} \mathbf{G}^{(l)} - \mathbf{K}^{(l)} \right\|_F + \left\| \mathbf{K} - \mathbf{K}^{(l)} \right\|_F \quad (33)$$

$$= O\left(n(2L\sqrt{2}\sigma_w)^l\right) + O\left(n(2L\sqrt{2}\sigma_w)^l\right) + O\left(n(8L^2\sigma_w^2)^l\right) \quad (34)$$

$$= O\left(n(2L\sqrt{2}\sigma_w)^l\right), \quad (35)$$

where (35) follows from  $\sigma_w^2 < 1/(8L^2)$ .

Next, we fix  $l$  to omit the explicit dependence on  $l$ . Specifically, let

$$l = \Theta(\log(2\lambda_*^{-1}n) / \log(\sqrt{2}/(4L\sigma_w))),$$

then from (35), we have

$$\left\| \frac{1}{m} \mathbf{G} - \mathbf{K} \right\|_F \leq \frac{\lambda_*}{2}. \quad (36)$$

Therefore, by Weyl's inequality (Ling et al., 2022, Lemma 5), it holds that

$$\max_{i \in [r]} \left| \lambda_i \left( \frac{1}{m} \mathbf{G} \right) - \lambda_i(\mathbf{K}) \right| \leq \left\| \frac{1}{m} \mathbf{G} - \mathbf{K} \right\|_2 \leq \left\| \frac{1}{m} \mathbf{G} - \mathbf{K} \right\|_F \leq \frac{\lambda_*}{2} \quad (37)$$

Now, by choosing  $i_0 := \arg \min_i \lambda_i(\mathbf{K})$ , we have

$$\lambda_{i_0}(\mathbf{K}) = \lambda_* \quad (38)$$

and

$$\left| \frac{1}{m} \lambda_{\min}(\mathbf{G}) - \lambda_* \right| \leq \frac{\lambda_*}{2}. \quad (39)$$

It follows from (38) and (39) that

$$\lambda_0 = \lambda_{\min}(\mathbf{G}) \geq \frac{m}{2} \lambda_*. \quad (40)$$

Consequently, w.p.  $\geq 1 - t$ , we have  $\lambda_0 \geq \frac{m}{2} \lambda_*$  provided that  $m = \Omega\left(\frac{n^2}{\lambda_*^2} \log \frac{n}{t}\right)$ .  $\square$

## 6 Checking the conditions of Theorem 6

In this section, we will show how the condition in Theorem 6 holds for some common activation functions. We first recall the definition of a traditional dual activation function, say  $\hat{\varphi}$ , associate with  $\varphi$  in (Daniely et al., 2016, Sect. 4.2):

$$\hat{\varphi}(x) = \mathbb{E}_{(u,v) \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} [\varphi(u)\varphi(v)]. \quad (41)$$

Then, by using a similar proof as (Daniely et al., 2016, Lemma 11), it can be shown that the new activation function (see Definition 3) satisfies

$$\tilde{Q}_{\alpha, \alpha}(x) = \frac{1}{q^2 \alpha^2} \sum_{n=1}^{\infty} a_n^2 \alpha^{2n} x^n \quad (42)$$

if  $\varphi(x) = \sum_{n=1}^{\infty} a_n h_n(x)$  (Hermite polynomial expansion) or  $\hat{\varphi}(x) = \sum_{n=1}^{\infty} a_n^2 x^n$ .

In the following, we apply (42) and show how the condition in Theorem 6 is fulfilled.

**Example 13.** Consider the sine activation,  $\varphi(x) = \sin(ax)$ . By (Daniely et al., 2016, Sect. 8), we have

$$\hat{\varphi}(x) = e^{-a^2} \sinh(a^2 x). \quad (43)$$

By Taylor's expansion of  $\sinh$  function, i.e.,

$$\sinh(x) = \sum_{r=0}^{\infty} \frac{1}{(2r+1)!} x^{2r+1}. \quad (44)$$

Hence, from (42) and (Daniely et al., 2016, Lemma 11), we have

$$Q_{\alpha, \alpha}(x) = \frac{1}{q^2 \alpha^2} e^{-a^2} \sum_{r=0}^{\infty} \frac{a^{4r+2} \alpha^{4r+2}}{(2r+1)!} x^{2r+1}, \quad (45)$$

which leads to

$$\mu_{r, \alpha}^2(\varphi) = \begin{cases} \frac{1}{q^2 \alpha^2} e^{-a^2} \frac{a^{2r} \alpha^{2r}}{r!} & r \bmod 2 = 1 \\ 0 & \text{otherwise} \end{cases}. \quad (46)$$

This means that the condition in Theorem 6 is satisfied.

**Example 14.** Consider the tanh activation function,  $\varphi(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ . By (Szego, 1959, Eq. 8.23.4),  $\varphi(x)$  can be uniquely described in the basis of Hermite polynomials,

$$\varphi(x) = \sum_{n=1}^{\infty} a_n h_n(x) \quad (47)$$

where

$$|a_n| = \frac{1}{\sqrt{\pi 2^n n!}} \frac{\Gamma(\frac{n}{2} + 1)}{\Gamma(n + 1)} \exp\left(-\frac{\pi \sqrt{2n}}{2}\right). \quad (48)$$

Hence, from (42), we obtain

$$Q_{\alpha, \alpha}(x) = \frac{1}{q^2 \alpha^2} \sum_{n=1}^{\infty} a_n^2 \alpha^{2n} x^n, \quad (49)$$

so we have

$$\mu_{r, \alpha}^2(\varphi) = \frac{1}{q^2 \alpha^2} a_n^2 \alpha^{2n} \quad (50)$$

This means that the condition in Theorem 6 is satisfied.

**Example 15.** Consider the sigmoid activation function  $\varphi(x) = \frac{1}{1+e^{-x}}$ . It is known that

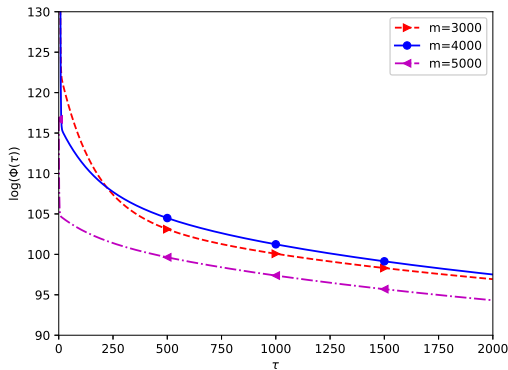
$$\varphi(x) = \frac{1 + \tanh(x/2)}{2}. \quad (51)$$

Hence, by using similar arguments as Example 14, we can prove that the condition in Theorem 6 is also satisfied.

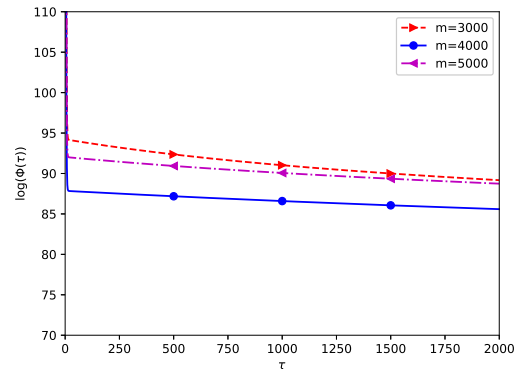
## 7 Numerical Results

In this section, we implement some experiments to verify Theorem 2. We evaluate the DEQ model on MNIST and CIFAR-10 datasets. For each dataset, the training dataset is generated by randomly sampling 500 images from the first and second classes. We use Gaussian initialization as Assumption 1 and normalize each data point as Assumption 2.

In the first experiment, we variate  $m$  and plot the training dynamic for MNIST and CIFAR-10 when  $\varphi$  is the sigmoid function ( $L = 1$ ). It can be seen from Fig. 1 that as  $m$  big enough and  $\tau$  sufficient large, the curves become straight lines. This fact re-affirms that (11) holds.



(a) MNIST



(b) CIFAR-10

Figure 1: Training dynamics at different values of  $m$ .



In the second experiment, we variate the activation function and plot the training dynamic for MNIST and CIFAR-10 at  $m = 3000$ . It can be seen from Fig. 2 that as  $m$  big enough and  $\tau$  sufficient large, the tanh network converges faster than the sigmoid or ReLU one for both datasets.

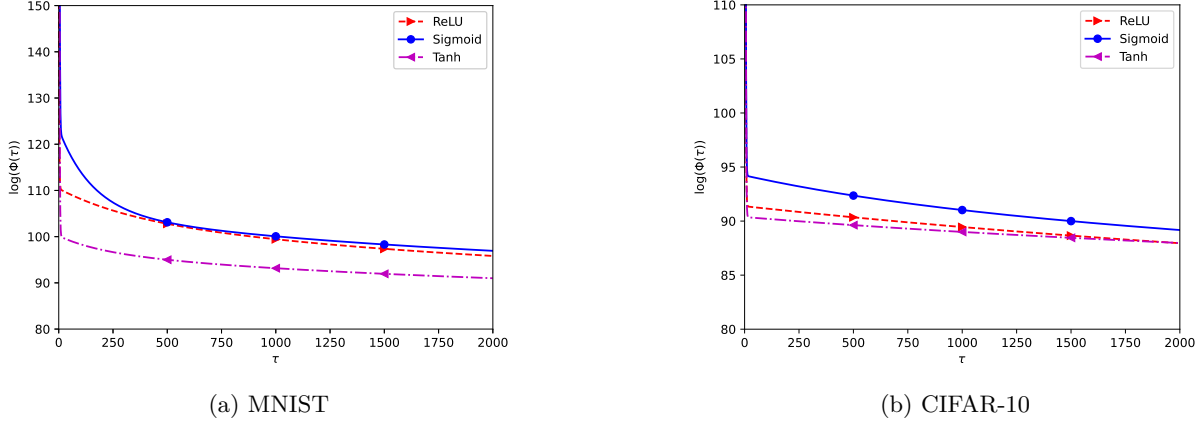


Figure 2: Training dynamics for different activation functions.

## 8 Conclusion

In this paper, we proved that the gradient descent converges to a globally optimal solution at a linear convergence rate for the quadratic loss function for the over-parametrized DEQ with  $L$ -bounded activation functions. This fascinating fact is also re-affirmed by our numerical experiments on MNIST and CFAR-10 datasets. To overcome new technical challenges caused by the non-linearity of activation functions, a novel population Gram matrix is introduced and a new form of dual activation with Hermite polynomial expansion is developed. An interesting future research direction is to study whether the linear convergence rate property still holds for other classes of activation functions.

## References

- Sanjeev Arora, Simon Shaolei Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Neural Information Processing Systems*, 2019.
- Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. Deep equilibrium models. *ArXiv*, abs/1909.01377, 2019.
- Shaojie Bai, Vladlen Koltun, and J. Zico Kolter. Multiscale deep equilibrium models. *ArXiv*, abs/2006.08656, 2020.
- Peter Bartlett and John Shawe-Taylor. *Generalization Performance of Support Vector Machines and Other Pattern Classifiers*, pp. 43–54. MIT Press, 1999.
- Peter Bartlett, Yoav Freund, Wee Sun Lee, and Robert E. Schapire. Boosting the margin: a new explanation for the effectiveness of voting methods. *The Annals of Statistics*, 26(5):1651 – 1686, 1998.
- F. Biggs and B. Guedj. Differentiable PAC-Bayes objectives with partially aggregated neural networks. *Entropy*, 23, 2021.
- Amit Daniely, Nathan Linial, and Shai Shalev-Shwartz. From average case complexity to improper learning complexity. *Proceedings of the forty-sixth annual ACM symposium on Theory of computing*, 2013.
- Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- Simon Shaolei Du, J. Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, 2018.
- G. K. Dziugaite and D. M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In *Uncertainty in Artificial Intelligence (UAI)*, 2017.
- Amedeo Roberto Esposito, Michael Gastpar, and Ibrahim Issa. Generalization error bounds via Rényi-f-divergences and maximal leakage. *IEEE Transactions on Information Theory*, 67(8):4986–5004, 2021.
- Davis Gilton, Greg Ongie, and Rebecca M. Willett. Deep equilibrium architectures for inverse problems in imaging. *IEEE Transactions on Computational Imaging*, 7:1123–1133, 2021.
- R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, 1985.
- D. Jakubovitz, R. Giryes, and M. R. D. Rodrigues. Generalization Error in Deep Learning. *Arxiv: 1808.01174*, 30, 2018.
- Michael Kearns and Leslie G. Valiant. Cryptographic limitations on learning boolean formulae and finite automata. In *JACM*, 1994.
- V. Koltchinskii and D. Panchenko. Empirical Margin Distributions and Bounding the Generalization Error of Combined Classifiers. *The Annals of Statistics*, 30(1):1 – 50, 2002.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60:84 – 90, 2012.
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. *ArXiv*, abs/1808.01204, 2018.
- Zenan Ling, Xingyu Xie, Qiuhao Wang, Zongpeng Zhang, and Zhouchen Lin. Global convergence of over-parameterized deep equilibrium models. *ArXiv*, abs/2205.13814, 2022.
- Quynh N. Nguyen. On the proof of global convergence of gradient descent for deep relu networks with linear widths. *ArXiv*, abs/2101.09612, 2021.

G. Szego. Orthogonal polynomials. *American Mathematical Society*, 1959.

T. Tao. *Topics in Random Matrix Theory*. American Mathematical Society, 2012.

Lan V. Truong. Generalization error bounds on deep learning with markov datasets. *Thirty-Sixth Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022a.

Lan V. Truong. On rademacher complexity-based generalization bounds for deep learning. *ArXiv*, abs/2208.04284, 2022b.

V. N. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.

Xingyu Xie, Qiuhaio Wang, Zenan Ling, Xia Li, Guangcan Liu, and Zhouchen Lin. Optimization induced equilibrium networks: An explicit optimization perspective for understanding equilibrium models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45:3604–3616, 2022.

A. Xu and M. Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. In *Advances of Neural Information Processing Systems (NIPS)*, 2017.

## A Appendix

### B Proof of Lemma 7

By Cauchy–Schwarz inequality, we have

$$|\tilde{Q}_{\alpha,\beta}(x)| \leq \frac{1}{\alpha\beta q^2} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} |\varphi(\alpha a)\varphi(\beta b)| \quad (52)$$

$$= \frac{1}{\alpha\beta q^2} \mathbb{E}_{(u,v)^T \sim \mathcal{N}\left(0, \begin{bmatrix} \alpha^2 & x\alpha\beta \\ x\alpha\beta & \beta^2 \end{bmatrix}\right)} |\varphi(u)\varphi(v)| \quad (53)$$

$$\leq \frac{1}{\alpha\beta} \sqrt{\frac{1}{q^2} \mathbb{E}_{a \sim \mathcal{N}(0,\alpha^2)} [\varphi^2(a)]} \sqrt{\frac{1}{q^2} \mathbb{E}_{b \sim \mathcal{N}(0,\beta^2)} [\varphi^2(b)]} \quad (54)$$

$$= \sqrt{\tilde{Q}_{\alpha,\alpha}(1)\tilde{Q}_{\beta,\beta}(1)}, \quad (55)$$

where (54) follows from Cauchy–Schwarz inequality. The equality in (54) holds if and only if  $\alpha = \beta$  and  $x = 1$ .

In addition, by the  $L$ -bounded property of  $\varphi$ , we also have

$$|\varphi(\alpha z) - \varphi(0)| \leq L|\alpha z|. \quad (56)$$

Hence, for any  $\alpha \geq 1$ , it holds that

$$|\varphi(\alpha z)| \leq |\varphi(0)| + L|\alpha||z| \quad (57)$$

$$\leq L(1 + |\alpha||z|) \quad (58)$$

$$\leq L|\alpha|\sqrt{2(1 + z^2)}. \quad (59)$$

From (59), we obtain

$$\mathbb{E}_{a \sim \mathcal{N}(0,\alpha^2)} [\varphi^2(a)] = \int_{-\infty}^{\infty} \frac{1}{\alpha\sqrt{2\pi}} \varphi^2(z) \exp\left(-\frac{z^2}{2\alpha^2}\right) dz \quad (60)$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \varphi^2(\alpha z) \exp\left(-\frac{z^2}{2}\right) dz \quad (61)$$

$$\leq 2L^2\alpha^2 \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} (1 + z^2) \exp\left(-\frac{z^2}{2}\right) dz \quad (62)$$

$$= 4L^2\alpha^2. \quad (63)$$

Similarly, we also have

$$\mathbb{E}_{b \sim \mathcal{N}(0, \beta^2)}[\varphi^2(b)] \leq 4L^2\beta^2. \quad (64)$$

From (54), (63) and (64), we obtain  $|\tilde{Q}_{\alpha, \beta}(x)| \leq 4L^2/q^2$  for all  $\alpha \geq 1$ ,  $\beta \geq 1$ , and  $x \in \mathbb{R}$ .

Now, for a fixed pair  $(\alpha \geq 1, \beta \geq 1)$ , define  $z := (u, v)$ ,  $\phi(z) := \varphi(u)\varphi(v)$ , and

$$\Sigma_x := \begin{bmatrix} \alpha^2 & x\alpha\beta \\ x\alpha\beta & \beta^2 \end{bmatrix}. \quad (65)$$

Then, by (Daniely et al., 2016, Lemma 12) we have

$$\frac{\partial \tilde{Q}_{\alpha, \beta}}{\partial \Sigma_x} = -\frac{1}{2q^2\alpha\beta} \mathbb{E}_{(u, v) \sim \mathcal{N}(0, \Sigma_x)} \left[ \frac{\partial \phi^2(z)}{\partial^2 z}(u, v) \right]. \quad (66)$$

On the other hand, we note that

$$\frac{\partial \phi^2(z)}{\partial^2 z}(u, v) = \begin{bmatrix} \frac{\partial^2 \varphi(u)}{\partial u^2} \varphi(v) & \frac{\partial \varphi(u)}{\partial u} \frac{\partial \varphi(v)}{\partial v} \\ \frac{\partial \varphi(u)}{\partial u} \frac{\partial \varphi(v)}{\partial v} & \frac{\partial^2 \varphi(v)}{\partial v^2} \varphi(u) \end{bmatrix}. \quad (67)$$

Hence, from (66) and (67) we have

$$\begin{aligned} \left\| \text{vec} \left( \frac{\partial \tilde{Q}_{\alpha, \beta}}{\partial \Sigma_x} \right) \right\|_{\infty} &\leq \frac{1}{2q^2\alpha\beta} \max \left\{ \mathbb{E}_{(u, v) \sim \mathcal{N}(0, \Sigma_x)} \left[ \left\| \frac{\partial^2 \varphi(u)}{\partial u^2} \varphi(v) \right\| \right], \mathbb{E}_{(u, v) \sim \mathcal{N}(0, \Sigma_x)} \left[ \left\| \frac{\partial \varphi(u)}{\partial u} \frac{\partial \varphi(v)}{\partial v} \right\| \right], \right. \\ &\quad \left. \mathbb{E}_{(u, v) \sim \mathcal{N}(0, \Sigma_x)} \left[ \left\| \frac{\partial^2 \varphi(v)}{\partial v^2} \varphi(u) \right\| \right] \right\}. \end{aligned} \quad (68)$$

Hence, by the assumption that  $\|\varphi\|_{\infty} \leq L, \|\varphi''\|_{\infty} \leq L$ , from (68) we obtain

$$\left\| \text{vec} \left( \frac{\partial \tilde{Q}_{\alpha, \beta}}{\partial \Sigma_x} \right) \right\|_{\infty} \leq \frac{L^2}{2q^2\alpha\beta}. \quad (69)$$

It follows that

$$|\tilde{Q}_{\alpha, \beta}(y) - \tilde{Q}_{\alpha, \beta}(x)| = \left| \int_x^y \frac{d\tilde{Q}_{\alpha, \beta}}{dt} dt \right| \quad (70)$$

$$= \left| \int_x^y \text{tr} \left( \left( \frac{\partial \tilde{Q}_{\alpha, \beta}}{\partial \Sigma_t} \right)^T \frac{\partial \Sigma_t}{dt} \right) dt \right| \quad (71)$$

$$\leq \int_x^y \left| \text{tr} \left( \left( \frac{\partial \tilde{Q}_{\alpha, \beta}}{\partial \Sigma_t} \right)^T \frac{\partial \Sigma_t}{dt} \right) \right| dt \quad (72)$$

$$= \int_x^y \left| \text{vec} \left( \frac{\partial \tilde{Q}_{\alpha, \beta}}{\partial \Sigma_t} \right)^T \text{vec} \left( \frac{\partial \Sigma_t}{dt} \right) \right| dt \quad (73)$$

$$\leq 4 \int_x^y \left\| \text{vec} \left( \frac{\partial \tilde{Q}_{\alpha, \beta}}{\partial \Sigma_t} \right) \right\|_{\infty} \left\| \text{vec} \left( \frac{\partial \Sigma_t}{dt} \right) \right\|_{\infty} dt \quad (74)$$

$$\leq \frac{4L^2}{2q^2\alpha\beta} \int_x^y \left\| \text{vec} \left( \frac{\partial \Sigma_t}{dt} \right) \right\|_{\infty} dt \quad (75)$$

$$= \frac{4L^2}{2q^2\alpha\beta} \alpha\beta |y - x| \quad (76)$$

$$= \frac{2L^2}{q^2} |y - x|. \quad (77)$$

## C Proof of Lemma 9

Observe that

$$\nu_{ii}^{(l)} = \frac{\sigma_w^2 \mathbf{K}_{ii}^{(l-1)} + d^{-1} \mathbf{x}_i^T \mathbf{x}_i}{\sigma_w^2 \mathbf{K}_{ii}^{(l-1)} + 1} \quad (78)$$

$$= 1. \quad (79)$$

From (15) and (19) in Definition 4 and (79), we have

$$\rho_{ii}^{(l)} = \sigma_w^2 \mathbf{K}_{ii}^{(l-1)} + 1. \quad (80)$$

In addition, from (16) and (18) in Definition 4 and (80), we also have

$$\rho_{ij}^{(l)} \nu_{ij}^{(l)} = \sigma_w^2 \mathbf{K}_{ij}^{(l-1)} + d^{-1} \mathbf{x}_i^T \mathbf{x}_j, \quad \forall i, j. \quad (81)$$

Replacing (19) in Definition 4 and (80) to (18) in Definition 4, we obtain for  $i \neq j$ ,

$$|\nu_{ij}^{(l)}| = \frac{|\sigma_w^2 \mathbf{K}_{ij}^{(l-1)} + d^{-1} \mathbf{x}_i^T \mathbf{x}_j|}{\sqrt{(\sigma_w^2 \mathbf{K}_{ii}^{(l-1)} + 1)(\sigma_w^2 \mathbf{K}_{jj}^{(l-1)} + 1)}} \quad (82)$$

$$= \frac{|2q^2 \sigma_w^2 \rho_{ij}^{(l-1)} Q_{ij}(\nu_{ij}^{(l-1)}) + d^{-1} \mathbf{x}_i^T \mathbf{x}_j|}{\sqrt{\rho_{ii}^{(l)} \rho_{jj}^{(l)}}} \quad (83)$$

$$= \frac{|Q_{ij}(\nu_{ij}^{(l-1)}) / \sqrt{Q_{ii}(1) Q_{jj}(1)} \sqrt{(2q^2 \sigma_w^2 \rho_{ii}^{(l-1)} Q_{ii}(1))(2q^2 \sigma_w^2 \rho_{jj}^{(l-1)} Q_{jj}(1))} + d^{-1} \mathbf{x}_i^T \mathbf{x}_j|}{\sqrt{\rho_{ii}^{(l)} \rho_{jj}^{(l)}}} \quad (84)$$

$$= \frac{|Q_{ij}(\nu_{ij}^{(l-1)}) / \sqrt{Q_{ii}(1) Q_{jj}(1)} \sqrt{(\rho_{ii}^{(l)} - 1)(\rho_{jj}^{(l)} - 1)} + d^{-1} \mathbf{x}_i^T \mathbf{x}_j|}{\sqrt{\rho_{ii}^{(l)} \rho_{jj}^{(l)}}} \quad (85)$$

$$\leq \frac{\sqrt{(\rho_{ii}^{(l)} - 1)(\rho_{jj}^{(l)} - 1)} + |d^{-1} \mathbf{x}_i^T \mathbf{x}_j|}{\sqrt{\rho_{ii}^{(l)} \rho_{jj}^{(l)}}} \quad (86)$$

$$\leq \frac{\sqrt{(\rho_{ii}^{(l)} - 1)(\rho_{jj}^{(l)} - 1)} + 1}{\sqrt{\rho_{ii}^{(l)} \rho_{jj}^{(l)}}} \quad (87)$$

$$\leq 1, \quad (88)$$

where (86) follows from Lemma 7, and (87) follows from  $d^{-1} |\mathbf{x}_i^T \mathbf{x}_j| \leq d^{-1} \|\mathbf{x}_i\|_2 \|\mathbf{x}_j\|_2 = 1$ .

## D Proof of Proposition 10

For all  $i, j \in [n] \times [n]$ , observe that

$$\begin{aligned} & |\mathbf{K}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l)}| \\ &= 2q^2 |\rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l+1)}) - \rho_{ij}^{(l)} Q_{ij}(\nu_{ij}^{(l)})| \end{aligned} \quad (89)$$

$$\leq 2q^2 |\rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l+1)}) - \rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l)})| + 2q^2 |\rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l)}) - \rho_{ij}^{(l)} Q_{ij}(\nu_{ij}^{(l)})|, \quad (90)$$

where (90) follows from the triangle inequality.

Now, we bound each term in (90). First, from Assumption 1 and Lemma 7, we have

$$2q^2 \sigma_w^2 Q_{ii}(1) \leq 8L^2 \sigma_w^2 < 1. \quad (91)$$

Therefore, from (15) we have

$$\rho_{ii}^{(l)} = \frac{1 - (2q^2\sigma_w^2 Q_{ii}(1))^l}{1 - 2q^2\sigma_w^2 Q_{ii}(1)}, \quad \forall i. \quad (92)$$

It follows that

$$|\rho_{ii}^{(l)} - \rho_{ii}^{(l+1)}| \leq O((2q^2\sigma_w^2 Q_{ii}(1))^l). \quad (93)$$

In addition, for  $i \neq j$ , we have

$$|\rho_{ij}^{(l)} - \rho_{ij}^{(l+1)}| = |\sqrt{\rho_{ii}^{(l+1)} \rho_{jj}^{(l+1)}} - \sqrt{\rho_{ii}^{(l)} \rho_{jj}^{(l)}}| \quad (94)$$

$$\leq \sqrt{\rho_{ii}^{(l+1)}} |\sqrt{\rho_{jj}^{(l+1)}} - \sqrt{\rho_{jj}^{(l)}}| + \sqrt{\rho_{jj}^{(l)}} |\sqrt{\rho_{ii}^{(l+1)}} - \sqrt{\rho_{ii}^{(l)}}| \quad (95)$$

$$\leq O((2q^2\sigma_w^2 Q_{ii}(1))^l). \quad (96)$$

From (91), (93), and (96), we obtain

$$|\rho_{ij}^{(l)} - \rho_{ij}^{(l+1)}| \leq O((8L^2\sigma_w^2)^l), \quad \forall i, j. \quad (97)$$

Now, we have

$$\begin{aligned} & |\rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l+1)}) - \rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l)})| \\ &= \left| \rho_{ij}^{(l+1)} \tilde{Q}_{\sqrt{2(\frac{\sigma_w^2}{m} \mathbf{G}_{ii} + 1)}, \sqrt{2(\frac{\sigma_w^2}{m} \mathbf{G}_{jj} + 1)}}(\nu_{ij}^{(l+1)}) \right. \\ &\quad \left. - \rho_{ij}^{(l+1)} \tilde{Q}_{\sqrt{2(\frac{\sigma_w^2}{m} \mathbf{G}_{ii} + 1)}, \sqrt{2(\frac{\sigma_w^2}{m} \mathbf{G}_{jj} + 1)}}(\nu_{ij}^{(l)}) \right| \end{aligned} \quad (98)$$

$$\leq \frac{2L^2}{q^2} |\rho_{ij}^{(l+1)} \nu_{ij}^{(l+1)} - \rho_{ij}^{(l+1)} \nu_{ij}^{(l)}| \quad (99)$$

$$\leq \frac{2L^2}{q^2} |\rho_{ij}^{(l+1)} \nu_{ij}^{(l+1)} - \rho_{ij}^{(l)} \nu_{ij}^{(l)}| + \frac{2L^2}{q^2} |\rho_{ij}^{(l)} - \rho_{ij}^{(l+1)}| |\nu_{ij}^{(l)}| \quad (100)$$

$$\leq \frac{2L^2}{q^2} |\rho_{ij}^{(l+1)} \nu_{ij}^{(l+1)} - \rho_{ij}^{(l)} \nu_{ij}^{(l)}| + \frac{2L^2}{q^2} |\rho_{ij}^{(l)} - \rho_{ij}^{(l+1)}| \quad (101)$$

$$= \frac{2L^2}{q^2} \sigma_w^2 |\mathbf{K}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l-1)}| + \frac{2L^2}{q^2} O((8L^2\sigma_w^2)^l), \quad (102)$$

where (99) follows from Lemma 7, (101) follows from Lemma 9, (102) follows from (25) in Lemma 9 and (97).

In addition, by using the fact that  $|Q_{\alpha,\beta}(x)| \leq \frac{4L^2}{q^2}$  for all  $\alpha \geq 1, \beta \geq 1$  in Lemma 7, we have

$$|\rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l)}) - \rho_{ij}^{(l)} Q_{ij}(\nu_{ij}^{(l)})| \leq \frac{4L^2}{q^2} |\rho_{ij}^{(l+1)} - \rho_{ij}^{(l)}| \quad (103)$$

$$= \frac{4L^2}{q^2} O((8L^2\sigma_w^2)^l), \quad (104)$$

where (104) follows from (97).

From (19), (102), and (104) we have

$$\begin{aligned} & |\mathbf{K}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l)}| \\ &= 2q^2 |\rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l+1)}) - \rho_{ij}^{(l)} Q_{ij}(\nu_{ij}^{(l)})| \end{aligned} \quad (105)$$

$$\leq 2q^2 |\rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l+1)}) - \rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l)})| + 2q^2 |\rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l)}) - \rho_{ij}^{(l)} Q_{ij}(\nu_{ij}^{(l)})| \quad (106)$$

$$\leq 2q^2 \left[ \frac{2L^2}{q^2} \sigma_w^2 |\mathbf{K}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l-1)}| + \frac{2L^2}{q^2} O((8L^2\sigma_w^2)^l) \right] + \frac{4L^2}{q^2} O((8L^2\sigma_w^2)^l). \quad (107)$$

By using induction, from (107) we have

$$|\mathbf{K}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l)}| = O((8L^2\sigma_w^2)^l). \quad (108)$$

Since  $\sigma_w^2 < 1/(8L^2)$ ,  $\{\mathbf{K}_{ij}^{(l)}\}_{l=1}^\infty$  can be easily shown to be a Cauchy sequence. From the completeness of  $\mathbb{R}$ , it holds that

$$\mathbf{K}_{ij}^{(l)} \rightarrow \mathbf{K}_{ij} \quad (109)$$

uniformly in  $i, j \in [n] \times [n]$  as  $l \rightarrow \infty$  for some matrix  $\mathbf{K}$ . By using the triangle inequality, we have

$$|\mathbf{K}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l)}| \geq |\mathbf{K}_{ij}^{(l)} - \mathbf{K}_{ij}| - |\mathbf{K}_{ij}^{(l+1)} - \mathbf{K}_{ij}|. \quad (110)$$

From (108) and (110), we obtain

$$|\mathbf{K}_{ij}^{(l)} - \mathbf{K}_{ij}| = O((8L^2\sigma_w^2)^l). \quad (111)$$

From (111), we obtain

$$\|\mathbf{K}^{(l)} - \mathbf{K}\|_F = O(n(8L^2\sigma_w^2)^l). \quad (112)$$

Now, by (19) and (109) we have

$$\mathbf{K}_{ij}^{(l)} = 2q^2\rho_{ij}^{(l)}Q_{ij}(\nu_{ij}^{(l)}) \quad (113)$$

and  $\mathbf{K}_{ij}^{(l)} \rightarrow \mathbf{K}_{ij}$ . On the other hand, by (91) we have  $2q^2\sigma_w^2Q_{ii}(1) < 1$ . It follows from (92) that

$$\rho_{ii}^{(l)} \rightarrow \frac{1}{1 - 2q^2\sigma_w^2Q_{ii}(1)} \quad (114)$$

as  $l \rightarrow \infty$ . Hence, it holds that  $\nu_{ij}^{(l)} \rightarrow \nu_{ij}$  uniformly in  $i, j \in [n] \times [n]$ .

Hence, by (30) in Lemma 9, we have

$$\nu_{ij} = \begin{cases} \frac{Q_{ij}(\nu_{ij})/\sqrt{Q_{ii}(1)Q_{jj}(1)}\sqrt{(\rho_{ii}-1)(\rho_{jj}-1)+d^{-1}\mathbf{x}_i^T\mathbf{x}_j}}{\sqrt{\rho_{ii}\rho_{jj}}}, & i \neq j, \\ 1, & i = j, \end{cases} \quad (115)$$

where

$$\rho_{ii} = \frac{1}{1 - 2q^2\sigma_w^2Q_{ii}(1)}. \quad (116)$$

## E Proof of Proposition 11

Assume that  $\mathbf{T}^{(l)} = [\mathbf{t}_1^{(l)}, \mathbf{t}_2^{(l)}, \dots, \mathbf{t}_n^{(l)}]$  where  $\mathbf{t}_i^{(l)} \in \mathbb{R}^m$  for all  $i \in [n]$ . By (1), we have

$$\mathbf{t}_i^{(l)} = \varphi(\mathbf{W}\mathbf{t}_i^{(l-1)} + \mathbf{U}\mathbf{x}_i), \quad \forall i \in [n]. \quad (117)$$

Hence, with probability at least  $1 - \exp(-\Omega(m))$ , we have

$$\|\mathbf{t}_i^{(l+1)} - \mathbf{t}_i^{(l)}\|_2 = \|\varphi(\mathbf{W}\mathbf{t}_i^{(l)} + \mathbf{U}\mathbf{x}_i) - \varphi(\mathbf{W}\mathbf{t}_i^{(l-1)} + \mathbf{U}\mathbf{x}_i)\|_2 \quad (118)$$

$$\leq L\|\mathbf{W}(\mathbf{t}_i^{(l)} - \mathbf{t}_i^{(l-1)})\|_2 \quad (119)$$

$$\leq L\|\mathbf{W}\|_2\|\mathbf{t}_i^{(l)} - \mathbf{t}_i^{(l-1)}\|_2 \quad (120)$$

$$\leq 2L\sqrt{2}\sigma_w\|\mathbf{t}_i^{(l)} - \mathbf{t}_i^{(l-1)}\| \quad (121)$$

where (119) is a consequence of the assumption that  $\varphi$  is  $L$ -bounded, and (121) follows from (Tao, 2012, Sect. (2.3)).

Therefore, for all  $l \geq 2$ , it holds that

$$\|\mathbf{t}_i^{(l)} - \mathbf{t}_i^{(l-1)}\|_2 \leq (2L\sqrt{2}\sigma_w)^l \|\mathbf{t}_i^{(1)} - \mathbf{t}_i^{(0)}\|_2 \quad (122)$$

$$= (2L\sqrt{2}\sigma_w)^l \|\mathbf{t}_i^{(1)}\|_2. \quad (123)$$

Now, let  $V \sim \mathcal{N}(0, 4)$  given  $\mathbf{x}_i$ . For each  $\mathbf{t}_i^{(l)}$ , we have

$$p_i := \mathbb{E} \left[ \frac{1}{m} (\mathbf{t}_i^{(1)})^T \mathbf{t}_i^{(1)} \right] = \mathbb{E} \left[ \frac{1}{m} \varphi(\mathbf{U}\mathbf{x}_i)^T \varphi(\mathbf{U}\mathbf{x}_i) \right] \quad (124)$$

$$= \mathbb{E}[\sigma(V^2)] \quad (125)$$

$$\leq 2(L^2 + L^2\mathbb{E}[V^2]) \quad (126)$$

$$= 10L^2, \quad (127)$$

where (125) follows from  $|\varphi(x) - \varphi(0)| \leq L|x|$  for all  $x \in \mathbb{R}$ .

Then, by using Bernstein's inequality, it holds with probability at least  $\geq 1 - 2\exp(-\Omega(mt^2))$  that

$$\left| \frac{1}{m} (\mathbf{t}_i^{(1)})^T \mathbf{t}_i^{(1)} - p_i \right| \leq t. \quad (128)$$

Hence, with probability at least  $1 - \exp(-\Omega(m)) - 2\exp(-\Omega(mt^2))$  it holds that

$$\|\mathbf{t}_i^{(l)} - \mathbf{t}_i^{(l-1)}\| \leq (2L\sqrt{2}\sigma_w)^l \sqrt{m(p_i + t)} \quad (129)$$

$$\leq (2L\sqrt{2}\sigma_w)^l \sqrt{m(10L^2 + t)}. \quad (130)$$

Then, for all  $r > s$ , with probability at least  $1 - \exp(-\Omega(m)) - 2\exp(-\Omega(mt^2))$ , we have

$$\|\mathbf{t}_i^{(r)} - \mathbf{t}_i^{(s)}\| \leq \sqrt{m(10L^2 + t)} (2L\sqrt{2}\sigma_w)^s \rightarrow 0 \quad (131)$$

as  $s \rightarrow \infty$  since  $2L\sqrt{2}\sigma_w < 1$ . Since  $\mathbb{R}$  is complete, hence we have

$$\|\mathbf{t}_i^{(l)} - \mathbf{t}_i\| \rightarrow 0 \quad (132)$$

for some vector  $\mathbf{t}_i$ .

It follows that

$$\|\mathbf{t}_i^{(l-1)} - \mathbf{t}_i\| - \|\mathbf{t}_i^{(l)} - \mathbf{t}_i\| \leq \|\mathbf{t}_i^{(l)} - \mathbf{t}_i^{(l-1)}\| \quad (133)$$

$$\leq \sqrt{m(10L^2 + t)} \|\mathbf{t}_i^{(1)}\| (2L\sqrt{2}\sigma_w)^l, \quad \forall l \geq 2. \quad (134)$$

From (134), with probability at least  $1 - \exp(-\Omega(m)) - 2\exp(-\Omega(mt^2))$  we have

$$\|\mathbf{t}_i^{(l)} - \mathbf{t}_i\| \leq \sqrt{m(10L^2 + t)} \|\mathbf{t}_i^{(1)}\| \sum_{k=l+1}^{\infty} (2L\sqrt{2}\sigma_w)^k \quad (135)$$

$$= \sqrt{m(10L^2 + t)} \frac{\|\mathbf{t}_i^{(1)}\| (2L\sqrt{2}\sigma_w)^{l+1}}{1 - 2L\sqrt{2}\sigma_w}. \quad (136)$$



Consequently, we have

$$\left| \mathbf{G}_{ij} - \mathbf{G}_{ij}^{(l)} \right| = |\mathbf{t}_i^T \mathbf{t}_j - (\mathbf{t}_i^{(l)})^T (\mathbf{t}_j^{(l)})| \quad (137)$$

$$\leq |\mathbf{t}_i^T \mathbf{t}_j - \mathbf{t}_i^T (\mathbf{t}_j^{(l)})| + |\mathbf{t}_i^T (\mathbf{t}_j^{(l)}) - (\mathbf{t}_i^{(l)})^T (\mathbf{t}_j^{(l)})| \quad (138)$$

$$\leq \|\mathbf{t}_i\| \|\mathbf{t}_j - \mathbf{t}_j^{(l)}\| + \|\mathbf{t}_j^{(l)}\| \|\mathbf{t}_i - \mathbf{t}_i^{(l)}\| \quad (139)$$

$$\begin{aligned} &\leq \sqrt{m(10L^2 + t)} \|\mathbf{t}_i\| \|\mathbf{t}_i^{(1)}\| \frac{(2L\sqrt{2}\sigma_w)^{l+1}}{1 - 2L\sqrt{2}\sigma_w} \\ &\quad + \sqrt{m(10L^2 + t)} \|\mathbf{t}_i^{(l)}\| \|\mathbf{t}_i^{(1)}\| \frac{(2L\sqrt{2}\sigma_w)^{l+1}}{1 - 2L\sqrt{2}\sigma_w}. \end{aligned} \quad (140)$$

Let  $t$  be an absolute constant. Finally, we obtain (31) from (140).

## F Proof of Proposition 12

Define

$$\hat{\mathbf{G}}_{ij}^{(l)} := \mathbb{E} \left[ \frac{1}{m} \mathbf{G}_{ij}^{(l)} \right]. \quad (141)$$

Then, by Lemma 8, we have

$$\hat{\mathbf{G}}_{ij}^{(l)} = \mathbb{E} \left[ \frac{1}{m} \varphi(\mathbf{M}\mathbf{h}_l)^T \varphi(\mathbf{M}\mathbf{h}_l') \right] \quad (142)$$

$$= \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, 2\mathbf{I})} [\varphi(\mathbf{w}^T \mathbf{h}_l) \varphi(\mathbf{w}^T \mathbf{h}_l')]. \quad (143)$$

Let

$$\hat{\mathbf{A}}_{ij}^{(l)} := \mathbf{h}_l^T \mathbf{h}_l', \quad \hat{\mathbf{A}}_{ii}^{(l)} := \|\mathbf{h}_l\|_2^2, \quad \hat{\mathbf{A}}_{jj}^{(l)} := \|\mathbf{h}_l'\|_2^2, \quad (144)$$

and define

$$\hat{\nu}_{ij}^{(l)} := \frac{\hat{\mathbf{A}}_{ij}^{(l)}}{\sqrt{\hat{\mathbf{A}}_{ii}^{(l)} \hat{\mathbf{A}}_{jj}^{(l)}}}. \quad (145)$$

Then, we have

$$\hat{\mathbf{G}}_{ij}^{(l)} = \mathbb{E}_{(u,v) \sim \mathcal{N}\left(0, 2 \begin{bmatrix} \|\mathbf{h}_l\|^2 & \mathbf{h}_l^T \mathbf{h}_l' \\ \mathbf{h}_l^T \mathbf{h}_l' & \|\mathbf{h}_l'\|^2 \end{bmatrix}\right)} [\varphi(u) \varphi(v)] \quad (146)$$

$$= \mathbb{E}_{(u,v) \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & \frac{\mathbf{h}_l^T \mathbf{h}_l'}{\|\mathbf{h}_l\| \|\mathbf{h}_l'\|} \\ \frac{\mathbf{h}_l^T \mathbf{h}_l'}{\|\mathbf{h}_l\| \|\mathbf{h}_l'\|} & 1 \end{bmatrix}\right)} [\varphi(\sqrt{2}\|\mathbf{h}_l\|u) \varphi(\sqrt{2}\|\mathbf{h}_l'\|v)] \quad (147)$$

$$= 2q^2 \|\mathbf{h}_l\| \|\mathbf{h}_l'\| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_l\|, \sqrt{2}\|\mathbf{h}_l'\|}(\hat{\nu}_{ij}^{(l)}) \quad (148)$$

$$= 2q^2 \sqrt{\hat{\mathbf{A}}_{ii}^{(l)} \hat{\mathbf{A}}_{jj}^{(l)}} \tilde{Q}_{\sqrt{2}\|\mathbf{h}_l\|, \sqrt{2}\|\mathbf{h}_l'\|}(\hat{\nu}_{ij}^{(l)}). \quad (149)$$

Now, we consider two cases:

- **Case 1:**  $i = j$ .

By Lemma 10, we have

$$\mathbf{G}_{ii}^{(l+1)} = \varphi(\mathbf{M}\mathbf{h}_{l+1})^T \varphi(\mathbf{M}\mathbf{h}_{l+1}), \quad (150)$$

where

$$\|\mathbf{h}_{l+1}\|^2 = \frac{\sigma_w^2}{m} \mathbf{G}_{ii}^{(l)} + 1. \quad (151)$$

Now, for a fixed  $\mathbf{h}_{l+1}$ , by Beinstein's inequality and (150), it holds with probability  $1 - \exp(-\Omega(m\varepsilon^2))$  that

$$\left| \frac{1}{m} \mathbf{G}_{ii}^{(l+1)} - \hat{\mathbf{G}}_{ii}^{(l+1)} \right| \leq \varepsilon/2. \quad (152)$$

On the other hand, by Preposition 11, with probability at least  $1 - n^2 \exp(-\Omega(m))$ , we have

$$\frac{1}{m} \|\mathbf{G} - \mathbf{G}^{(l+1)}\|_F = O\left(n(2L\sqrt{2}\sigma_w)^{l+1}\right). \quad (153)$$

Since  $2L\sqrt{2}\sigma_w < 1$ , it holds with probability at least  $1 - n^2 \exp(-\Omega(m))$  that

$$\left| \frac{1}{m} \mathbf{G}_{ii}^{(l+1)} - \mathbf{G}_{ii} \right| = O\left(n(2L\sqrt{2}\sigma_w)^{l+1}\right) = o(1). \quad (154)$$

From (154),  $\|\mathbf{h}_{l+1}\|^2 = O(1)$  with probability at least  $1 - n^2 \exp(-\Omega(m))$ . Since the  $\varepsilon$ -net size for  $\mathbf{h}_{l+1} \in \mathbb{R}^{2l+d+2}$  is at most  $\exp\{O(l \log \frac{1}{\varepsilon})\}$ , it holds with probability at least  $1 - n^2 \exp(-\Omega(m\varepsilon^2) + O(l \log \frac{1}{\varepsilon}))$ ,

$$\left| \frac{1}{m} \mathbf{G}_{ii}^{(l+1)} - \hat{\mathbf{G}}_{ii}^{(l+1)} \right| \leq \varepsilon/2. \quad (155)$$

Now, observe that

$$\hat{\mathbf{G}}_{ii}^{(l+1)} = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(0, 2\mathbf{I})} [\varphi^2(\mathbf{w}^T \mathbf{h}_{l+1})] \quad (156)$$

$$= 2q^2 \|\mathbf{h}_{l+1}\|^2 \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1). \quad (157)$$

On the other hand, we also have

$$\mathbf{K}_{ii}^{(l+1)} = 2q^2 \rho_{ii}^{(l+1)} Q_{ii}(1) \quad (158)$$

$$= 2q^2 (\sigma_w^2 \mathbf{K}_{ii}^{(l)} + 1) Q_{ii}(1). \quad (159)$$

It follows that

$$\begin{aligned} & \left| \hat{\mathbf{G}}_{ii}^{(l+1)} - \mathbf{K}_{ii}^{(l+1)} \right| \\ &= 2q^2 \left| \|\mathbf{h}_{l+1}\|^2 \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1) - (\sigma_w^2 \mathbf{K}_{ii}^{(l)} + 1) Q_{ii}(1) \right| \end{aligned} \quad (160)$$

$$= 2q^2 \left| \left( \frac{\sigma_w^2}{m} \mathbf{G}_{ii}^{(l)} + 1 \right) \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1) - (\sigma_w^2 \mathbf{K}_{ii}^{(l)} + 1) Q_{ii}(1) \right| \quad (161)$$

$$\begin{aligned} & \leq 2q^2 \left| \left( \frac{\sigma_w^2}{m} \mathbf{G}_{ii}^{(l)} + 1 \right) \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1) - (\sigma_w^2 \mathbf{K}_{ii}^{(l)} + 1) \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1) \right| \\ & \quad + 2q^2 (\sigma_w^2 \mathbf{K}_{ii}^{(l)} + 1) \left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1) - Q_{ii}(1) \right| \end{aligned} \quad (162)$$

$$\leq 8L^2 \sigma_w^2 \left| \frac{\mathbf{G}_{ii}^{(l)}}{m} - \mathbf{K}_{ii}^{(l)} \right| + 2q^2 (\sigma_w^2 \mathbf{K}_{ii}^{(l)} + 1) \left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1) - Q_{ii}(1) \right|, \quad (163)$$

where (163) follows from (23) in Lemma 7.

Now, let

$$\|\mathbf{h}\|^2 := \frac{\sigma_w^2}{m} \mathbf{G}_{ii} + 1. \quad (164)$$

Then, we have

$$|\|\mathbf{h}_{l+1}\|^2 - \|\mathbf{h}\|^2| = \frac{\sigma_w^2}{m} |\mathbf{G}_{ii}^{(l)} - \mathbf{G}_{ii}| \quad (165)$$

$$\leq \frac{1}{m} \|\mathbf{G}^{(l)} - \mathbf{G}\|_F \quad (166)$$

$$= O\left(n(2L\sqrt{2}\sigma_w)^l\right) \quad (167)$$

where (165) follows from (151) and (164), and (167) follows from (153). Since  $\mathbf{G}_{ii} = \|\mathbf{t}_i\|^2 \geq 0$ ,  $\mathbf{G}_{ii}^{(l)} = \|\mathbf{t}_i^{(l)}\|^2 \geq 0$ , from (151), (164), and (167), we obtain

$$|\|\mathbf{h}_{l+1}\| - \|\mathbf{h}\|| = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (168)$$

Now, for any  $a \in \mathbb{R}$  note that

$$\begin{aligned} & \left| \varphi^2(\sqrt{2}\|\mathbf{h}_{l+1}\|a) - \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right| \\ &= \left| \varphi(\sqrt{2}\|\mathbf{h}_{l+1}\|a) - \varphi(\sqrt{2}\|\mathbf{h}\|a) \right| \left| \sigma(\sqrt{2}\|\mathbf{h}_{l+1}\|a) + \sigma(\sqrt{2}\|\mathbf{h}\|a) \right|. \end{aligned} \quad (169)$$

On the other hand, we have

$$\left| \varphi(\sqrt{2}\|\mathbf{h}_{l+1}\|a) - \varphi(\sqrt{2}\|\mathbf{h}\|a) \right| \leq L\sqrt{2}|a| \|\mathbf{h}_{l+1}\| - \|\mathbf{h}\|, \quad (170)$$

$$\left| \varphi(\sqrt{2}\|\mathbf{h}_{l+1}\|a) + \varphi(\sqrt{2}\|\mathbf{h}\|a) \right| \leq 2|\varphi(0)| + L\sqrt{2}(\|\mathbf{h}_{l+1}\| + \|\mathbf{h}\|)|a| \quad (171)$$

where we use  $|\varphi(x)| - |\varphi(0)| \leq |\varphi(x) - \varphi(0)| \leq L|x|$  on (171).

From (169), (170), and (171), we obtain

$$\left| \varphi^2(\sqrt{2}\|\mathbf{h}_{l+1}\|a) - \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right| \leq 2L\sqrt{2}|\varphi(0)||a| \|\mathbf{h}_{l+1}\| - \|\mathbf{h}\| + 2L^2|a|^2 \|\mathbf{h}_{l+1}\|^2 - \|\mathbf{h}\|^2 \quad (172)$$

$$= |a| \left[ \varepsilon + O\left(nL(2L\sqrt{2}\sigma_w)^l\right) \right] + |a|^2 \left[ \varepsilon + O\left(nL^2(2L\sqrt{2}\sigma_w)^l\right) \right] \quad (173)$$

where (173) follows from (167) and (168).

From (173), we obtain

$$\begin{aligned} & \left| \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}_{l+1}\|a) \right] - \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right] \right| \\ & \leq \mathbb{E}_{a \sim \mathcal{N}(0,1)} [|a|] \left[ \varepsilon + O\left(nL(2L\sqrt{2}\sigma_w)^l\right) \right] + \mathbb{E}_{a \sim \mathcal{N}(0,1)} [|a|^2] \left[ \varepsilon + O\left(nL^2(2L\sqrt{2}\sigma_w)^l\right) \right] \end{aligned} \quad (174)$$

$$= O\left(\varepsilon + nL^2(2L\sqrt{2}\sigma_w)^l\right). \quad (175)$$

Similarly, we also have

$$\begin{aligned} \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right] &\leq \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \left( |\varphi(0)| + L\sqrt{2}\|\mathbf{h}\||a| \right)^2 \right] \\ &= O(1). \end{aligned} \quad (176)$$

It follows that

$$\begin{aligned} &\left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1) - Q_{ii}(1) \right| \\ &= \left| \frac{1}{2q^2\|\mathbf{h}_{l+1}\|^2} \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}_{l+1}\|a) \right] - \frac{1}{2q^2\|\mathbf{h}\|^2} \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right] \right| \end{aligned} \quad (178)$$

$$\begin{aligned} &\leq \left| \frac{1}{2q^2\|\mathbf{h}_{l+1}\|^2} \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}_{l+1}\|a) \right] - \frac{1}{2q^2\|\mathbf{h}_{l+1}\|^2} \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right] \right| \\ &+ \left| \frac{1}{2q^2\|\mathbf{h}_{l+1}\|^2} \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right] - \frac{1}{2q^2\|\mathbf{h}\|^2} \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right] \right| \end{aligned} \quad (179)$$

$$\begin{aligned} &\leq \frac{1}{2q^2\|\mathbf{h}_{l+1}\|^2} \left| \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}_{l+1}\|a) \right] - \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right] \right| \\ &+ \frac{1}{2q^2} \left| \frac{1}{\|\mathbf{h}_{l+1}\|^2} - \frac{1}{\|\mathbf{h}\|^2} \right| \mathbb{E}_{a \sim \mathcal{N}(0,1)} \left[ \varphi^2(\sqrt{2}\|\mathbf{h}\|a) \right]. \end{aligned} \quad (180)$$

By combining (167), (175), and (177), from (180), we obtain

$$\left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}_{l+1}\|}(1) - Q_{i,i}(1) \right| = O\left(\varepsilon + nL^2(2L\sqrt{2}\sigma_w)^l\right). \quad (181)$$

On the other hand, by Proposition 10, with probability at least  $1 - m \exp(-\Omega(m))$ , we have

$$\|\mathbf{K} - \mathbf{K}^{(l+1)}\|_F = O\left(n(8L^2\sigma_w^2)^{l+1}\right) = O\left(n(2L\sqrt{2}\sigma_w)^{l+1}\right). \quad (182)$$

It follows that

$$\|\mathbf{K}_{ii}^{(l+1)} - \mathbf{K}_{ii}\| = O\left(n(2L\sqrt{2}\sigma_w)^{l+1}\right). \quad (183)$$

From (181), (183), by setting

$$\varepsilon := O\left(nL^2(2L\sqrt{2}\sigma_w)^{l+1}\right) \quad (184)$$

from (163), we obtain

$$\left| \hat{\mathbf{G}}_{ii}^{(l+1)} - \mathbf{K}_{ii}^{(l+1)} \right| \leq 8L^2\sigma_w^2 \left| \frac{\mathbf{G}_{ii}^{(l)}}{m} - \mathbf{K}_{ii}^{(l)} \right| + 2\varepsilon. \quad (185)$$

It follows from (155) and (185) that with probability at least  $1 - n^2 \exp\{-\Omega(m\varepsilon^2) + O(l \log \frac{1}{\varepsilon})\}$ ,

$$\left| \frac{1}{m} \mathbf{G}_{ii}^{(l+1)} - \mathbf{K}_{ii}^{(l+1)} \right| \leq \left| \frac{1}{m} \mathbf{G}_{ii}^{(l+1)} - \hat{\mathbf{G}}_{ii}^{(l+1)} \right| + \left| \hat{\mathbf{G}}_{ii}^{(l+1)} - \mathbf{K}_{ii}^{(l+1)} \right| \quad (186)$$

$$\leq 8L^2\sigma_w^2 \left| \frac{1}{m} \mathbf{G}_{ii}^{(l)} - \mathbf{K}_{ii}^{(l)} \right| + 2\varepsilon, \quad (187)$$

which implies that with probability at least  $1 - n^2 l \exp\{-\Omega(m\varepsilon^2) + O(l \log \frac{1}{\varepsilon})\}$ , we have

$$\left| \mathbf{G}_{ii}^{(l)} - \mathbf{K}_{ii}^{(l)} \right| \leq \frac{1 - (8L^2\sigma_w^2)^l}{1 - 8L^2\sigma_w^2} 2\varepsilon. \quad (188)$$

Final note is that since  $\varepsilon = O(nL^2(2L\sqrt{2}\sigma_w)^{l+1})$ , it holds with probability at least  $1 - n^2l \exp\{-\Omega(8^l L^{2l} \sigma_w^{2l} mnL^2) + O(l^2)\} \geq 1 - n^2 \exp\{-\Omega(8^l L^{2l} \sigma_w^{2l} mnL^2) + O(l^2)\}$ , we have

$$\left| \mathbf{G}_{ii}^{(l)} - \mathbf{K}_{ii}^{(l)} \right| = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (189)$$

- **Case 2:**  $i \neq j$ .

For this case, let

$$\|\mathbf{h}\|^2 = \frac{\sigma_w^2}{m} \mathbf{G}_{ii} + 1, \quad (190)$$

$$\|\mathbf{h}'\|^2 = \frac{\sigma_w^2}{m} \mathbf{G}_{jj} + 1. \quad (191)$$

By Preposition 11, with probability at least  $1 - n^2 \exp(-\Omega(m))$ , we have

$$\frac{1}{m} \left\| \mathbf{G} - \mathbf{G}^{(l)} \right\|_F = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (192)$$

In addition, we also have

$$\|\mathbf{h}_{l+1}\|^2 = \frac{\sigma_w^2}{m} \mathbf{G}_{ii}^{(l)} + 1 \geq 1, \quad (193)$$

$$\|\mathbf{h}'_{l+1}\|^2 = \frac{\sigma_w^2}{m} \mathbf{G}_{jj}^{(l)} + 1 \geq 1. \quad (194)$$

Hence, we have

$$|\|\mathbf{h}_{l+1}\| - \|\mathbf{h}\|| = O\left(|\|\mathbf{h}_{l+1}\|^2 - \|\mathbf{h}\|^2|\right) \quad (195)$$

$$= \frac{\sigma_w^2}{m} \left\| \mathbf{G}_{ii}^{(l)} - \mathbf{G}_{ii} \right\| \quad (196)$$

$$\leq \frac{\sigma_w^2}{m} \left\| \mathbf{G}^{(l)} - \mathbf{G} \right\|_F \quad (197)$$

$$= O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (198)$$

Then, it holds that

$$\begin{aligned} & \left| \hat{\mathbf{G}}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l+1)} \right| \\ &= 2q^2 \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\hat{\nu}_{ij}^{(l)}) - \rho_{ij}^{(l+1)} Q_{ij}(\nu_{ij}^{(l+1)}) \right| \end{aligned} \quad (199)$$

$$\begin{aligned} &\leq 2q^2 \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\hat{\nu}_{ij}^{(l+1)}) - \rho_{ij}^{(l+1)} \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\nu_{ij}^{(l)}) \right| \\ &\quad + 2q^2 \rho_{ij}^{(l+1)} \left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\nu_{ij}^{(l)}) - Q_{ij}(\nu_{ij}^{(l+1)}) \right|. \end{aligned} \quad (200)$$

Now, for all  $|x| \leq 1$ , we have

$$\begin{aligned} \left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) - Q_{ij}(x) \right| &\leq \left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) - \tilde{Q}_{\sqrt{2}\|\mathbf{h}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) \right| \\ &\quad + \left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) - Q_{ij}(x) \right|. \end{aligned} \quad (201)$$

On the other hand, we have

$$\begin{aligned} & \left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) - Q_{ij}(x) \right| \\ &= \left| \frac{1}{2q^2\|\mathbf{h}\|\|\mathbf{h}'_{l+1}\|} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'_{l+1}\|b) \right| \end{aligned} \quad (202)$$

$$- \frac{1}{2q^2\|\mathbf{h}\|\|\mathbf{h}'\|} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'\|b) \quad (203)$$

$$\begin{aligned} & \leq \left| \frac{1}{2q^2\|\mathbf{h}\|\|\mathbf{h}'_{l+1}\|} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'_{l+1}\|b) \right. \\ & \quad \left. - \frac{1}{2q^2\|\mathbf{h}\|\|\mathbf{h}'_{l+1}\|} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'\|b) \right| \\ & \quad + \left| \frac{1}{2q^2\|\mathbf{h}\|\|\mathbf{h}'_{l+1}\|} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'\|b) \right. \\ & \quad \left. - \frac{1}{2q^2\|\mathbf{h}\|\|\mathbf{h}'\|} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'\|b) \right| \end{aligned} \quad (204)$$

$$\begin{aligned} & \leq \frac{1}{2q^2\|\mathbf{h}\|\|\mathbf{h}'_{l+1}\|} \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} \left| \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'_{l+1}\|b) \right. \\ & \quad \left. - \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'\|b) \right| \\ & \quad + \frac{1}{2q^2\|\mathbf{h}\|} \left| \frac{1}{\|\mathbf{h}'_{l+1}\|} - \frac{1}{\|\mathbf{h}'\|} \right| \mathbb{E}_{(a,b)^T \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & x \\ x & 1 \end{bmatrix}\right)} \left| \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'\|b) \right|. \end{aligned} \quad (205)$$

In addition, we have

$$|\varphi(\sqrt{2}\|\mathbf{h}\|a)| \leq |\varphi(0)| + L\sqrt{2}\|\mathbf{h}\||a| \quad (206)$$

$$|\varphi(\sqrt{2}\|\mathbf{h}'\|b)| \leq |\varphi(0)| + L\sqrt{2}\|\mathbf{h}'\||b|. \quad (207)$$

It follows that

$$\begin{aligned} & \left| \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'_{l+1}\|b) - \varphi(\sqrt{2}\|\mathbf{h}\|a) \varphi(\sqrt{2}\|\mathbf{h}'\|b) \right| \\ &= \left| \varphi(\sqrt{2}\|\mathbf{h}\|a) \right| \left| \varphi(\sqrt{2}\|\mathbf{h}'_{l+1}\|b) - \varphi(\sqrt{2}\|\mathbf{h}'\|b) \right| \end{aligned} \quad (208)$$

$$\leq \left( |\varphi(0)| + L\sqrt{2}\|\mathbf{h}\||a| \right) \left| \varphi(\sqrt{2}\|\mathbf{h}'_{l+1}\|b) - \varphi(\sqrt{2}\|\mathbf{h}'\|b) \right| \quad (209)$$

$$\leq L\sqrt{2} \left( |\varphi(0)| + L\sqrt{2}\mathbb{E}[\|\mathbf{h}\||a|] \right) \left| \|\mathbf{h}'_{l+1}\| - \|\mathbf{h}'\| \right| |b|. \quad (210)$$

On the other hand, by Bernstein's inequality, with probability at least  $1 - \exp(-\Omega(m)\varepsilon^2)$ , it holds that

$$\left| \|\mathbf{h}'\| - \mathbb{E}[\|\mathbf{h}'\|] \right| \leq \varepsilon. \quad (211)$$

From (198) and (211), we have

$$||\mathbf{h}'_{l+1}|| - ||\mathbf{h}'|| \leq ||\mathbf{h}'_{l+1}|| - ||\mathbf{h}|| + ||\mathbf{h}|| - ||\mathbf{h}'|| \quad (212)$$

$$\leq \varepsilon + O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (213)$$

Now, by setting

$$\varepsilon := O\left(n(2\sqrt{2}\sigma_w)^l\right), \quad (214)$$

from (213), we obtain

$$||\mathbf{h}'_{l+1}|| - \mathbb{E}[||\mathbf{h}'||] = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (215)$$

Similarly, we also have

$$||\mathbf{h}'_{l+1}|| - ||\mathbf{h}'|| = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (216)$$

From (205), (210), and (215), we obtain

$$\left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) - Q_{ij}(x) \right| = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (217)$$

Similarly, we can prove that

$$\left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) - \tilde{Q}_{\sqrt{2}\|\mathbf{h}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) \right| = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (218)$$

From (201), (217), and (218), we obtain

$$\left| \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(x) - Q_{ij}(x) \right| \leq O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (219)$$

Next, we aim to upper bound

$$2q^2 \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\hat{\nu}_{ij}^{(l+1)}) - \rho_{ij}^{(l+1)} \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\nu_{ij}^{(l+1)}) \right|.$$

Observe that with probability at least  $1 - n^2 \exp(-\Omega(m))$ , it holds for all  $l$  sufficiently large that

$$\begin{aligned} & \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\hat{\nu}_{ij}^{(l+1)}) - \rho_{ij}^{(l+1)} \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\nu_{ij}^{(l+1)}) \right| \\ & \leq \left| \left( \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} - \rho_{ij}^{(l+1)} \right) \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\hat{\nu}_{ij}^{(l+1)}) \right| \\ & \quad + \left| \rho_{ij}^{(l+1)} \left( \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\hat{\nu}_{ij}^{(l+1)}) - \tilde{Q}_{\sqrt{2}\|\mathbf{h}_{l+1}\|, \sqrt{2}\|\mathbf{h}'_{l+1}\|}(\nu_{ij}^{(l+1)}) \right) \right| \end{aligned} \quad (220)$$

$$\begin{aligned} & \leq \frac{4L^2}{q^2} \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} - \rho_{ij}^{(l+1)} \right| \\ & \quad + \rho_{ij}^{(l+1)} \frac{4L^2}{q^2} \max\{(\sqrt{2}\|\mathbf{h}_{l+1}\| + 1)^2, (\sqrt{2}\|\mathbf{h}'_{l+1}\| + 1)^2\} |\hat{\nu}_{ij}^{(l+1)} - \nu_{ij}^{(l+1)}| \end{aligned} \quad (221)$$

$$\begin{aligned} &\leq \frac{4L^2}{q^2} \max\{(\sqrt{2}\|\mathbf{h}_{l+1}\| + 1)^2, (\sqrt{2}\|\mathbf{h}'_{l+1}\| + 1)^2\} \left[ \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} - \rho_{ij}^{(l+1)} \right| \right. \\ &\quad \left. + \rho_{ij}^{(l+1)} |\hat{\nu}_{ij}^{(l+1)} - \nu_{ij}^{(l+1)}| \right] \end{aligned} \quad (222)$$

$$\begin{aligned} &\leq \frac{4L^2}{q^2} \max\left\{(\sqrt{2}\|\mathbf{h}\| + 1 + \varepsilon)^2, (\sqrt{2}\|\mathbf{h}\| + 1 + \varepsilon)^2\right\} \left[ \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} - \rho_{ij}^{(l+1)} \right| \right. \\ &\quad \left. + \rho_{ij}^{(l+1)} |\hat{\nu}_{ij}^{(l+1)} - \nu_{ij}^{(l+1)}| \right], \end{aligned} \quad (223)$$

where (221) follows from Lemma 7, and (223) follows from (215) and (216).

On the other hand, we have

$$\begin{aligned} &\left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} - \rho_{ij}^{(l+1)} \right| + \rho_{ij}^{(l+1)} |\hat{\nu}_{ij}^{(l+1)} - \nu_{ij}^{(l+1)}| \\ &= \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} - \rho_{ij}^{(l+1)} \right| \\ &\quad + \left| \left( \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} + \rho_{ij}^{(l+1)} - \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} \right) \hat{\nu}_{ij}^{(l+1)} - \rho_{ij}^{(l+1)} \nu_{ij}^{(l+1)} \right| \end{aligned} \quad (224)$$

$$\leq 2 \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} - \rho_{ij}^{(l+1)} \right| + \left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} \hat{\nu}_{ij}^{(l+1)} - \rho_{ij}^{(l+1)} \nu_{ij}^{(l+1)} \right|, \quad (225)$$

where (225) follows from  $|\hat{\nu}_{ij}^{(l)}| \leq 1$ .

On the other hand, since  $\rho_{ii}^{(l+1)} = \sqrt{\rho_{ii}^{(l+1)} \rho_{jj}^{(l+1)}}$ , we also have

$$\begin{aligned} &\left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} - \rho_{ij}^{(l+1)} \right| \\ &= 2q^2 \left| \sqrt{\left( \frac{\sigma_w^2}{m} \mathbf{G}_{ii}^{(l)} + 1 \right) \left( \frac{\sigma_w^2}{m} \mathbf{G}_{jj}^{(l)} + 1 \right)} - \sqrt{(\sigma_w^2 \mathbf{K}_{ii}^{(l)} + 1)(\sigma_w^2 \mathbf{K}_{jj}^{(l)} + 1)} \right| \end{aligned} \quad (226)$$

$$= O\left(n(2L\sqrt{2}\sigma_w)^l\right), \quad (227)$$

where (227) follows from (189).

Moreover, note that

$$\sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} \hat{\nu}_{ij}^{(l+1)} = \hat{\mathbf{A}}_{ij}^{(l+1)} \quad (228)$$

$$= \|\mathbf{h}_{l+1}\|^2 \quad (229)$$

$$= \frac{\sigma_w^2}{m} \mathbf{G}_{ij}^{(l)} + \frac{1}{d} \mathbf{x}_i^T \mathbf{x}_j \quad (230)$$

and

$$\rho_{ij}^{(l+1)} \nu_{ij}^{(l+1)} = \nu_{ij}^{(l+1)} \sqrt{\rho_{ii}^{(l+1)} \rho_{jj}^{(l+1)}} \quad (231)$$

$$= \nu_{ij}^{(l+1)} \sqrt{(\sigma_w^2 \mathbf{K}_{ii}^{(l)} + 1)(\sigma_w^2 \mathbf{K}_{jj}^{(l)} + 1)} \quad (232)$$

$$= \sigma_w^2 \mathbf{K}_{ij}^{(l)} + \frac{1}{d} \mathbf{x}_i^T \mathbf{x}_j. \quad (233)$$

Thus, it holds that

$$\left| \sqrt{\hat{\mathbf{A}}_{ii}^{(l+1)} \hat{\mathbf{A}}_{jj}^{(l+1)}} \hat{\nu}_{ij}^{(l+1)} - \rho_{ij}^{(l+1)} \nu_{ij}^{(l+1)} \right| = \sigma_w^2 \left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l)} \right|. \quad (234)$$



Thus, with probability at least  $1 - l \exp(-\Omega(m\varepsilon^2) + O(l \log 1/\varepsilon))$ , it holds that

$$\left| \hat{\mathbf{G}}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l+1)} \right| \leq \sigma_w^2 \left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l)} \right| + \varepsilon. \quad (235)$$

On the other hand, by Lemma 10, we have

$$\mathbf{G}_{ij}^{(l+1)} = \varphi(\mathbf{M}\mathbf{h}_{l+1})^T \varphi(\mathbf{M}\mathbf{h}'_{l+1}). \quad (236)$$

Hence, for a fixed vector pair  $\mathbf{h}_{l+1}, \mathbf{h}'_{l+1}$ , by Bernstein's inequality, with probability at least  $1 - \exp(-\Omega(m\varepsilon^2))$  it holds that

$$\left| \frac{1}{m} \mathbf{G}_{ij}^{(l+1)} - \hat{\mathbf{G}}_{ij}^{(l+1)} \right| \leq \varepsilon. \quad (237)$$

Then, by using  $\varepsilon$ -net arguments as in Case 1, with probability at least  $1 - l \exp(-\Omega(m\varepsilon^2) + O(l \log 1/\varepsilon))$ , we have

$$\left| \frac{1}{m} \mathbf{G}_{ij}^{(l+1)} - \hat{\mathbf{G}}_{ij}^{(l+1)} \right| \leq \varepsilon. \quad (238)$$

Consequently, we have

$$\left| \frac{1}{m} \mathbf{G}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l+1)} \right| \leq \left| \frac{1}{m} \mathbf{G}_{ij}^{(l+1)} - \hat{\mathbf{G}}_{ij}^{(l+1)} \right| + \left| \hat{\mathbf{G}}_{ij}^{(l+1)} - \mathbf{K}_{ij}^{(l+1)} \right| \quad (239)$$

$$\leq 2\varepsilon + \sigma_w^2 \left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l)} \right| \quad (240)$$

where (240) follows from (235) and (238).

By applying the induction argument, one can show that for  $l \geq 1$ , it holds with probability at least  $1 - l^2 \exp(-\Omega(m\varepsilon^2) + O(l \log 1/\varepsilon))$ , we have

$$\left| \frac{1}{m} \mathbf{G}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l)} \right| \leq \frac{2\varepsilon}{1 - \sigma_w^2}. \quad (241)$$

By the choice of  $\varepsilon$  in (214), it holds that with probability at least  $1 - n^2 \exp\{-\Omega(8^l L^{2l} \sigma_w^{2l} mn L^2) + O(l^2)\}$ , we have

$$\left| \mathbf{G}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l)} \right| = O\left(n(2L\sqrt{2}\sigma_w)^l\right). \quad (242)$$

## G Proof of Theorem 6

Since  $\mathbf{U}\mathbf{x}_i$  is a Gaussian vector with zero-mean and variance depending on  $\|\mathbf{x}_i\|^2$ . On the other hand, by the Assumption 2,  $\|\mathbf{x}_i\| = \sqrt{d}$ . Hence, from  $\mathbf{t}_i = \varphi(\mathbf{W}\mathbf{t}_i + \mathbf{U}\mathbf{x}_i)$ , it is easy to see that  $\mathbb{E}[G_{ii}] = \mathbb{E}[\|\mathbf{t}_i\|^2]$  does not depend on  $i \in [n]$ . This means that  $\mathbb{E}[\mathbf{G}_{ii}] = \mathbb{E}[\mathbf{G}_{jj}]$  for all  $i, j \in [n] \times [n]$ . Hence,  $Q_{ij}(x)$  has the form  $\bar{Q}_{\alpha, \alpha}(x)$  for some  $\alpha \geq 1$ .

Thanks to this fact, from Proposition 10 and the assumption on this theorem, for all  $(i, j) \in [n] \times [n]$ , it holds that

$$\mathbf{K}_{ij} = 2q^2 Q_{ij}(\nu_{ij}) \sqrt{\rho_{ii} \rho_{jj}} \quad (243)$$

$$= 2q^2 \sqrt{\rho_{ii} \rho_{jj}} \sum_{r=0}^{\infty} \mu_{r, \alpha}^2(\varphi) \nu_{ij}^r, \quad (244)$$

where

$$\nu_{ij} = \frac{Q_{ij}(\nu_{ij})/\sqrt{Q_{ii}(1)Q_{jj}(1)}\sqrt{(\rho_{ii}-1)(\rho_{jj}-1)} + d^{-1}\mathbf{x}_i^T\mathbf{x}_j}{\sqrt{\rho_{ii}\rho_{jj}}}. \quad (245)$$

Here,

$$\rho_{ii} = \frac{1}{1 - 2q^2\sigma_w^2 Q_{ii}(1)}. \quad (246)$$

Now, by Lemma 9, we have  $|\nu_{ij}| \leq 1$  for all  $(i, j) \in [n] \times [n]$ . Let  $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n]$  where  $\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n$  be unit vectors such that  $\nu_{ij} = \mathbf{h}_i^T \mathbf{h}_j$  for all  $(i, j) \in [n] \times [n]$ . It is easy to check that  $[(\mathbf{H}^T \mathbf{H})^{\odot r}]_{ij} = (\mathbf{h}_i^T \mathbf{h}_j)^r$  holds for all  $(i, j) \in [n] \times [n]$ . Let  $\tilde{\mathbf{K}}$  be a  $n \times n$  matrix such that

$$\tilde{\mathbf{K}}_{ij} = \mathbf{K}_{ij}/\sqrt{\rho_{ii}\rho_{jj}}, \quad \forall i, j \in [n] \times [n]. \quad (247)$$

Then,  $\tilde{\mathbf{K}}$  can be written as

$$\tilde{\mathbf{K}} = 2q^2 \sum_{r=0}^{\infty} \mu_{r,\alpha}^2(\varphi) (\mathbf{H}^T \mathbf{H})^{(\odot r)}. \quad (248)$$

Now, for any unit vector  $\mathbf{u} = [u_1, u_2, \dots, u_n]^T \in \mathbb{R}^n$ , it holds that

$$\mathbf{u}^T (\mathbf{H}^T \mathbf{H})^{(\odot r)} \mathbf{u} = \sum_{i,j} u_i u_j (\mathbf{h}_i^T \mathbf{h}_j)^r \quad (249)$$

$$= \sum_i u_i^2 + \sum_{i \neq j} u_i u_j \nu_{ij}^r \quad (250)$$

$$= 1 + \sum_{i \neq j} u_i u_j \nu_{ij}^r. \quad (251)$$

Next, we show that  $|\nu_{ij}| < 1$  if  $i \neq j$ . Indeed, assume that there exists  $i \neq j$  such that  $|\nu_{ij}| \geq 1$ . Then, from (30) in Lemma 9, we have

$$1 \leq |\nu_{ij}| \quad (252)$$

$$= \left| \frac{Q_{ij}(\nu_{ij})/\sqrt{Q_{ii}(1)Q_{jj}(1)}\sqrt{(\rho_{ii}-d^{-1}\|\mathbf{x}_i\|_2^2)(\rho_{jj}-d^{-1}\|\mathbf{x}_j\|_2^2)} + d^{-1}\mathbf{x}_i^T\mathbf{x}_j}{\sqrt{\rho_{ii}\rho_{jj}}} \right| \quad (253)$$

$$\leq \frac{\sqrt{(\rho_{ii}-d^{-1}\|\mathbf{x}_i\|_2^2)(\rho_{jj}-d^{-1}\|\mathbf{x}_j\|_2^2)} + |d^{-1}\mathbf{x}_i^T\mathbf{x}_j|}{\sqrt{\rho_{ii}\rho_{jj}}} \quad (254)$$

$$< \frac{\sqrt{(\rho_{ii}-d^{-1}\|\mathbf{x}_i\|_2^2)(\rho_{jj}-d^{-1}\|\mathbf{x}_j\|_2^2)} + 1}{\sqrt{\rho_{ii}\rho_{jj}}} \quad (255)$$

$$\leq 1, \quad (256)$$

where (254) follows from Lemma 7, and (255) follows by the fact that since  $\mathbf{x}_i \nparallel \mathbf{x}_j$ , from Cauchy-Schwarz inequality and Assumption 2, we have  $\mathbf{x}_i^T \mathbf{x}_j < \|\mathbf{x}_i\|_2 \|\mathbf{x}_j\|_2 = d$ . This is a contradiction. Hence, we have  $|\beta| < 1$  where

$$\beta := \max_{i \neq j} |\nu_{ij}|. \quad (257)$$

Now, by taking  $r > -\frac{\log n}{\log \beta}$ , we have

$$\left| \sum_{i \neq j} u_i u_j \nu_{ij}^r \right| \leq \sum_{i \neq j} |u_i| |u_j| \beta^r \quad (258)$$

$$\leq \left( \sum_i |\nu_i| \right)^2 \beta^r \quad (259)$$

$$\leq n\beta^r \quad (260)$$

$$< 1. \quad (261)$$

From (251) and (261), we obtain

$$\mathbf{u}^T (\mathbf{H}^T \mathbf{H})^{(\odot r)} \mathbf{u} > 0, \quad \forall \mathbf{u}, \quad (262)$$

so  $(\mathbf{H}^T \mathbf{H})^{(\odot r)}$  is positive definite. Following Theorem 6, it holds that  $\mu_{r,\alpha}^2(\varphi) > 0$  for infinitely many values of  $r$ . Hence,  $\tilde{\mathbf{K}}$  is positive definite.

Now, let  $\mathbf{\Gamma} = \{\sqrt{\rho_{ii}\rho_{jj}}\}_{i,j}$  be an  $n \times n$  matrix where the  $(i, j)$  element is  $\sqrt{\rho_{ii}\rho_{jj}}$ . Then, we have

$$\mathbf{K} = \tilde{\mathbf{K}} \odot \mathbf{\Gamma}. \quad (263)$$

Now, for any vector  $\mathbf{u} = [u_1, u_2, \dots, u_n]^T$ , we have

$$\mathbf{u}^T \mathbf{\Gamma} \mathbf{u} = \sum_{i,j} u_i u_j \sqrt{\rho_{ii}\rho_{jj}} \quad (264)$$

$$= \left( \sum_i u_i \sqrt{\rho_{ii}} \right)^2 \quad (265)$$

$$\geq 0. \quad (266)$$

Hence,  $\mathbf{\Gamma}$  is positive semi-definite. Now, by applying (Ling et al., 2022, Lemma 6), we have

$$\lambda_{\min}(\mathbf{K}) \geq \left( \min_i \rho_{ii} \right) \lambda_{\min}(\tilde{\mathbf{K}}) \quad (267)$$

$$\geq \lambda_{\min}(\tilde{\mathbf{K}}) > 0, \quad (268)$$

so  $\mathbf{K}$  is positive definite with the smallest eigenvalue  $\lambda_* > 0$ .

## H Proof of Theorem 2

The following proof follows the same steps as (Ling et al., 2022, Proof of Theorem 1). There are some small changes by the change of the activation function. First, we recall the two important auxiliary lemmas:

**Lemma 16.** (Horn & Johnson, 1985, Sect. 5.8) Let  $\mathbf{\Delta} = \mathbf{B} - \mathbf{A}$  where  $\mathbf{A}$  and  $\mathbf{B}$  are square complex matrices. Then, it holds that

$$\|\mathbf{B}^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\mathbf{A}^{-1}\mathbf{\Delta}\|}. \quad (269)$$

**Lemma 17.** (Weyl's inequality)(Ling et al., 2022, Lemma 5) Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$  with their singular values satisfying  $\sigma_1(\mathbf{A}) \geq \sigma_2(\mathbf{A}) \geq \dots \geq \sigma_r(\mathbf{A})$  and  $\sigma_1(\mathbf{B}) \geq \sigma_2(\mathbf{B}) \geq \dots \geq \sigma_r(\mathbf{B})$  and  $r = \min(m, n)$ . Then,

$$\max_{i \in [r]} |\sigma_i(\mathbf{A}) - \sigma_i(\mathbf{B})| \leq \|\mathbf{A} - \mathbf{B}\| \quad (270)$$

Based on these two lemmas, we can prove the following result:

**Lemma 18.** For each  $s \in [0, \tau]$ , suppose that  $\|\mathbf{W}(s)\|_2 \leq \bar{\rho}_w$ ,  $\|\mathbf{U}(s)\|_2 \leq \bar{\rho}_u$  and  $\|\mathbf{a}(s)\|_2 \leq \bar{\rho}_a$ . It holds that

$$\|\mathbf{T}(s)\|_{\mathcal{F}} \leq c_a \|\mathbf{X}\|_F + c_m \quad (271)$$

and

$$\|\nabla_{\mathbf{W}} \Phi(s)\|_F \leq c_u (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2, \quad (272)$$

$$\|\nabla_U \Phi(s)\|_F \leq c_u \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2, \quad (273)$$

$$\|\nabla_a \Phi(s)\|_F \leq (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2. \quad (274)$$

Furthermore, for each  $k, s \in [0, \tau]$ , it holds that

$$\begin{aligned} \|\mathbf{T}(k) - \mathbf{T}(s)\| &\leq \frac{L}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{W}(k) - \mathbf{W}(s)\|_2 \\ &\quad + \frac{L}{1 - L\bar{\rho}_w} \|\mathbf{U}(k) - \mathbf{U}(s)\|_2 \|\mathbf{X}\|_F \end{aligned} \quad (275)$$

and

$$\begin{aligned} \|\hat{\mathbf{y}}(k) - \hat{\mathbf{y}}(s)\|_2 &\leq \bar{\rho}_a \left[ \frac{L}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{W}(k) - \mathbf{W}(s)\|_2 \right. \\ &\quad \left. + \frac{L}{1 - L\bar{\rho}_w} \|\mathbf{U}(k) - \mathbf{U}(s)\|_2 \|\mathbf{X}\|_F \right] + (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{a}(k) - \mathbf{a}(s)\|_2. \end{aligned} \quad (276)$$

*Proof.* Observe that  $\mathbf{T}(s) = \varphi(\mathbf{W}(s)\mathbf{T}(s) + \mathbf{U}(s)\mathbf{X})$ . Using the fact that  $|\varphi(x) - \varphi(0)| \leq L|x|$  (Lipschitz condition of  $\varphi$ ), we have

$$\|\mathbf{T}(s) - \varphi(0)\|_F = \|\varphi(\mathbf{W}(s)\mathbf{T}(s) + \mathbf{U}(s)\mathbf{X}) - \varphi(0)\|_F \quad (277)$$

$$\leq L \|\mathbf{W}(s)\mathbf{T}(s) + \mathbf{U}(s)\mathbf{X}\|_F \quad (278)$$

$$\leq L \left( \|\mathbf{W}(s)\|_2 \|\mathbf{T}(s)\|_F + \|\mathbf{U}(s)\|_2 \|\mathbf{X}\|_F \right) \quad (279)$$

$$\leq L\bar{\rho}_w \|\mathbf{T}(s)\|_F + L\bar{\rho}_u \|\mathbf{X}\|_F. \quad (280)$$

From (280), we have

$$\|\mathbf{T}(s)\|_F \leq \|\varphi(0)\|_F + L\bar{\rho}_w \|\mathbf{T}(s)\|_F + L\bar{\rho}_u \|\mathbf{X}\|_F \quad (281)$$

$$= m^2 \varphi(0) + L\bar{\rho}_w \|\mathbf{T}(s)\|_F + L\bar{\rho}_u \|\mathbf{X}\|_F. \quad (282)$$

Since  $\bar{\rho}_w < 1/L$ , from (282), we obtain

$$\|\mathbf{T}(s)\|_F \leq c_a \|\mathbf{X}\|_F + c_m. \quad (283)$$

Now, we prove (272)-(274). By using Lemma 16 with  $\mathbf{A} = \mathbf{I}_{m,n}$ ,  $\mathbf{B} = \mathbf{J}(s)$ ,  $\mathbf{\Delta} = -\mathbf{D}(s)(\mathbf{I}_n \otimes \mathbf{W}(s))$ , we have

$$\|\mathbf{J}(s)^{-1}\|_2 \leq \frac{1}{1 - \|\mathbf{D}(\tau)(\mathbf{I}_n \otimes \mathbf{W}(s))\|_2} \quad (284)$$

$$\leq \frac{1}{1 - \|\mathbf{D}(s)\|_2 \|\mathbf{W}(s)\|_2}. \quad (285)$$

On the other hand since  $\|\varphi'\|_\infty \leq L$ , we have

$$\|\mathbf{D}(s)\|_2 \leq L. \quad (286)$$

Hence, from (285), we have

$$\|\mathbf{J}(s)^{-1}\|_2 \leq \frac{1}{1 - L\bar{\rho}_w}, \quad (287)$$

and thus it holds that

$$\|\mathbf{R}(s)\|_2 \leq \|\mathbf{a}(s)\|_2 \|\mathbf{J}(s)^{-1}\|_2 \|\mathbf{D}(s)\|_2 \quad (288)$$

$$\leq \frac{L\bar{\rho}_a}{1 - L\bar{\rho}_w}. \quad (289)$$

Then, we have

$$\|\nabla_{\mathbf{W}}\Phi(s)\|_F = \|\text{vec}(\nabla_{\mathbf{W}}\Phi(s))\|_2 \quad (290)$$

$$= \|(\mathbf{T}(s) \otimes \mathbf{I}_m)\mathbf{R}(s)^T(\hat{\mathbf{y}}(s) - \mathbf{y})\|_2 \quad (291)$$

$$\leq \|\mathbf{T}(s)\|_2 \|\mathbf{R}(s)\|_2 \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2 \quad (292)$$

$$\leq \frac{L\bar{\rho}_a}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2, \quad (293)$$

$$\|\nabla_{\mathbf{U}}\Phi(s)\|_F = \|\text{vec}(\nabla_{\mathbf{U}}\Phi(s))\|_2 \quad (294)$$

$$= \|(\mathbf{X} \otimes \mathbf{I}_m)\mathbf{R}(s)^T(\hat{\mathbf{y}}(s) - \mathbf{y})\|_2 \quad (295)$$

$$\leq \frac{L\bar{\rho}_a}{1 - L\bar{\rho}_w} \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2, \quad (296)$$

$$\|\nabla_{\mathbf{a}}\Phi(s)\|_F = \|\mathbf{T}(s)(\hat{\mathbf{y}}(s) - \mathbf{y})\| \quad (297)$$

$$\leq (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2. \quad (298)$$

Next, we prove (275). Observe that

$$\begin{aligned} \|\mathbf{T}(k) - \mathbf{T}(s)\|_F &= \|\varphi(\mathbf{W}(k)\mathbf{T}(k) + \mathbf{U}(k)\mathbf{X}) - \varphi(\mathbf{W}(s)\mathbf{T}(s) + \mathbf{U}(s)\mathbf{X})\|_F \end{aligned} \quad (299)$$

$$\leq L\|\mathbf{W}(k)\mathbf{T}(k) + \mathbf{U}(k)\mathbf{X} - \mathbf{W}(s)\mathbf{T}(s) - \mathbf{U}(s)\mathbf{X}\|_F \quad (300)$$

$$\begin{aligned} &\leq L(\|\mathbf{W}(k)\mathbf{T}(k) - \mathbf{W}(k)\mathbf{T}(s)\|_F + \|\mathbf{W}(k)\mathbf{T}(s) - \mathbf{W}(s)\mathbf{T}(s)\|_F \\ &\quad + \|\mathbf{U}(k)\mathbf{X} - \mathbf{U}(s)\mathbf{X}\|_F) \end{aligned} \quad (301)$$

$$\begin{aligned} &\leq L\|\mathbf{W}(k)\|_2 \|\mathbf{T}(k) - \mathbf{T}(s)\|_F + L\|\mathbf{W}(k) - \mathbf{W}(s)\|_2 \|\mathbf{T}(s)\|_F \\ &\quad + L\|\mathbf{U}(k) - \mathbf{U}(s)\|_2 \|\mathbf{X}\|_F \end{aligned} \quad (302)$$

$$\begin{aligned} &\leq L\bar{\rho}_w \|\mathbf{T}(k) - \mathbf{T}(s)\|_F + L(c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{W}(k) - \mathbf{W}(s)\|_2 \\ &\quad + L\|\mathbf{U}(k) - \mathbf{U}(s)\|_2 \|\mathbf{X}\|_F. \end{aligned} \quad (303)$$

From (303), we obtain

$$\begin{aligned} \|\mathbf{T}(k) - \mathbf{T}(s)\|_F &\leq \frac{L}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{W}(k) - \mathbf{W}(s)\|_2 \\ &\quad + \frac{L}{1 - L\bar{\rho}_w} \|\mathbf{U}(k) - \mathbf{U}(s)\|_2 \|\mathbf{X}\|_F. \end{aligned} \quad (304)$$

Finally, we prove (276). Observe that

$$\begin{aligned} \|\hat{\mathbf{y}}(k) - \hat{\mathbf{y}}(s)\|_F &= \|\mathbf{a}(k)\mathbf{T}(k) - \mathbf{a}(s)\mathbf{Z}(s)\|_F \end{aligned} \quad (305)$$

$$\leq \|\mathbf{a}(k)\mathbf{T}(k) - \mathbf{a}(k)\mathbf{T}(s)\|_F + \|\mathbf{a}(k)\mathbf{T}(s) - \mathbf{a}(s)\mathbf{T}(s)\|_F \quad (306)$$

$$\leq \|\mathbf{a}(k)\|_2 \|\mathbf{T}(k) - \mathbf{T}(s)\|_F + \|\mathbf{a}(k) - \mathbf{a}(s)\|_2 \|\mathbf{T}(s)\|_F \quad (307)$$

$$\begin{aligned} &\leq \bar{\rho}_a \left[ \frac{L}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{W}(k) - \mathbf{W}(s)\|_2 \right. \\ &\quad \left. + \frac{L}{1 - L\bar{\rho}_w} \|\mathbf{U}(k) - \mathbf{U}(s)\|_2 \|\mathbf{X}\|_F \right] + (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{a}(k) - \mathbf{a}(s)\|_2. \end{aligned} \quad (308)$$

□

Now, we return to prove Theorem 2. We prove by induction for every  $\tau > 0$ ,

$$\|\mathbf{W}(s)\| \leq \bar{\rho}_w, \|\mathbf{U}(s)\| \leq \bar{\rho}_u, \|\mathbf{a}(s)\|_2 \leq \bar{\rho}_a, s \in [0, \tau], \quad (309)$$

$$\lambda_s \geq \frac{\lambda_0}{2}, s \in [0, \tau], \quad (310)$$

$$\Phi(s) \leq \left(1 - \eta \frac{\lambda_0}{2}\right)^s \Phi(0), \quad s \in [0, \tau]. \quad (311)$$

For  $\tau = 0$ , it is clear that (309)-(311) hold. Assume that (309)-(311) holds up to  $\tau$  iterations. Then, by using triangle inequality, we have

$$\|\mathbf{W}(\tau + 1) - \mathbf{W}(0)\|_F \leq \sum_{s=0}^{\tau} \|\mathbf{W}(s+1) - \mathbf{W}(s)\|_F \quad (312)$$

$$= \sum_{s=0}^{\tau} \eta \|\nabla_{\mathbf{W}} \Phi(s)\|_F \quad (313)$$

$$\leq \eta \sum_{s=0}^{\tau} c_u (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2 \quad (314)$$

$$= \eta c_u (c_a \|\mathbf{X}\|_F + c_m) \sum_{s=0}^{\tau} \left(1 - \eta \frac{\lambda_0}{2}\right)^{s/2} \|\hat{\mathbf{y}}(0) - \hat{\mathbf{y}}\|_2 \quad (315)$$

where (314) follows from Lemma 18. Let  $u := \sqrt{1 - \eta \lambda_0 / 2}$ . Then  $\|\mathbf{W}(\tau + 1) - \mathbf{W}(0)\|_F$  can be bounded with

$$\begin{aligned} & \frac{2}{\lambda_0} (1 - u^2) \frac{1 - u^{\tau+1}}{1 - u} c_u (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(0) - \mathbf{y}\| \\ & \leq \frac{4}{\lambda_0} c_u (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(0) - \mathbf{y}\| \end{aligned} \quad (316)$$

$$\leq \delta. \quad (317)$$

Then, we have

$$\|\mathbf{W}(\tau + 1)\| \leq \|\mathbf{W}(0)\|_2 + \delta = \bar{\rho}_w < 1. \quad (318)$$

Using the similar technique, one can show that

$$\|\mathbf{U}(\tau + 1) - \mathbf{U}(0)\|_F \leq \sum_{s=0}^{\tau} \|\mathbf{U}(s+1) - \mathbf{U}(s)\|_2 \quad (319)$$

$$= \sum_{s=0}^{\tau} \eta \|\nabla_{\mathbf{U}} \Phi(s)\|_F \quad (320)$$

$$\leq \sum_{s=0}^{\tau} \eta c_u \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2 \quad (321)$$

$$\leq \eta c_u \|\mathbf{X}\|_F \sum_{s=0}^{\tau} \left(1 - \eta \frac{\lambda_0}{2}\right)^{s/2} \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2 \quad (322)$$

$$\leq \frac{4}{\lambda_0} c_u \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2 \quad (323)$$

$$\leq \delta. \quad (324)$$

$$\|\mathbf{a}(\tau+1) - \mathbf{a}(0)\|_F \leq \sum_{s=0}^{\tau} \|\mathbf{a}(s+1) - \mathbf{a}(s)\|_F \quad (325)$$

$$= \sum_{s=0}^{\tau} \eta \|\nabla_{\mathbf{a}} \Phi(s)\|_F \quad (326)$$

$$\leq \eta(c_a \|\mathbf{X}\|_F + c_m) \sum_{s=0}^{\tau} \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2 \quad (327)$$

$$\leq \eta(c_a \|\mathbf{X}\|_F + c_m) \sum_{s=0}^{\tau} \left(1 - \eta \frac{\lambda_0}{2}\right)^{s/2} \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2 \quad (328)$$

$$\leq \frac{4}{\lambda_0} (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2 \quad (329)$$

$$\leq \delta. \quad (330)$$

Finally, using (275), we have

$$\begin{aligned} \|\mathbf{T}(\tau+1) - \mathbf{T}(0)\| &\leq \frac{L}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{W}(\tau+1) - \mathbf{W}(0)\|_2 \\ &\quad + \frac{L}{1 - L\bar{\rho}_w} \|\mathbf{U}(\tau+1) - \mathbf{U}(0)\|_2 \|\mathbf{X}\|_F \end{aligned} \quad (331)$$

$$\begin{aligned} &\leq \frac{L}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \frac{4}{\lambda_0} c_u (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2 \\ &\quad + \frac{L}{1 - L\bar{\rho}_w} \frac{4}{\lambda_0} c_u \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2 \|\mathbf{X}\|_F \end{aligned} \quad (332)$$

$$= \frac{4L}{(1 - L\bar{\rho}_w)\lambda_0} \left[ (c_a \|\mathbf{X}\|_F + c_m)^2 + c_u \|\mathbf{X}\|_F^2 \right] \|\hat{\mathbf{y}}(0) - \mathbf{y}\|_2 \quad (333)$$

$$\leq \frac{2 - \sqrt{2}}{2} \sqrt{\lambda_0} \quad (334)$$

by (11).

By Wely's inequality, it implies that the least singular value of  $\mathbf{T}(\tau+1)$  satisfies  $\sigma_{\min}(\mathbf{T}(\tau+1)) \geq \sqrt{\frac{\lambda_0}{2}}$ . Thus, it holds  $\lambda_{\tau+1} \geq \frac{\lambda_0}{2}$ .

Now, we define  $\mathbf{g} := \mathbf{a}(\tau+1)^T \mathbf{T}(\tau)$  and note that

$$\begin{aligned} &\Phi(\tau+1) - \Phi(\tau) \\ &= \frac{1}{2} \|\hat{\mathbf{y}}(\tau+1) - \hat{\mathbf{y}}(\tau)\|_2^2 + (\hat{\mathbf{y}}(\tau+1) - \mathbf{g})^T (\hat{\mathbf{y}}(\tau) - \mathbf{y}) + (\mathbf{g} - \hat{\mathbf{y}}(\tau))^T (\hat{\mathbf{y}}(\tau) - \mathbf{y}). \end{aligned} \quad (335)$$

We bound each term of the RHS of this equation individually. First, using (276), we have

$$\begin{aligned} &\|\hat{\mathbf{y}}(\tau+1) - \hat{\mathbf{y}}(\tau)\|_2 \\ &\leq \bar{\rho}_a \left[ \frac{L}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{W}(\tau+1) - \mathbf{W}(\tau)\|_2 \right. \\ &\quad \left. + \frac{L}{1 - L\bar{\rho}_w} \|\mathbf{U}(\tau+1) - \mathbf{U}(\tau)\|_2 \|\mathbf{X}\|_F \right] + (c_a \|\mathbf{X}\|_F + c_m) \|\mathbf{a}(k) - \mathbf{a}(s)\|_2 \end{aligned} \quad (336)$$

$$\begin{aligned} &= \bar{\rho}_a \left[ \frac{L}{1 - L\bar{\rho}_w} (c_a \|\mathbf{X}\|_F + c_m) \eta c_u (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \right. \\ &\quad \left. + \frac{L}{1 - L\bar{\rho}_w} \eta c_u \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \|\mathbf{X}\|_F \right] \\ &\quad + (c_a \|\mathbf{X}\|_F + c_m) \eta (c_a \|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \end{aligned} \quad (337)$$

$$= \eta C_1 \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2, \quad (338)$$

where  $C_1 := c_u^2(c_a\|\mathbf{X}\|_F + c_m)^2 + c_u^2\|\mathbf{X}\|_F^2 + (c_a\|\mathbf{X}\|_F + c_m)^2$ .

On the other hand, we have

$$\begin{aligned} & (\hat{\mathbf{y}}(\tau+1) - \mathbf{g})^T(\hat{\mathbf{y}}(\tau) - \mathbf{y}) \\ &= \mathbf{a}(\tau+1)^T(\mathbf{T}(\tau+1) - \mathbf{T}(\tau))(\hat{\mathbf{y}}(\tau) - \mathbf{y}) \end{aligned} \quad (339)$$

$$\leq \|\mathbf{a}(\tau+1)\|_2 \|\mathbf{T}(\tau+1) - \mathbf{T}(\tau)\|_2 \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \quad (340)$$

$$\begin{aligned} & \leq \|\mathbf{a}(\tau+1)\|_2 \left[ \frac{L}{1-L\bar{\rho}_w} (c_a\|\mathbf{X}\|_F + c_m) \|\mathbf{W}(\tau+1) - \mathbf{W}(\tau)\|_2 \right. \\ & \quad \left. + \frac{L}{1-L\bar{\rho}_w} \|\mathbf{U}(k) - \mathbf{U}(s)\|_2 \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \right] \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \end{aligned} \quad (341)$$

$$\begin{aligned} & \leq \bar{\rho}_a \left[ \frac{L}{1-L\bar{\rho}_w} (c_a\|\mathbf{X}\|_F + c_m) \|\eta c_u (c_a\|\mathbf{X}\|_F + c_m) \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \right. \\ & \quad \left. + \frac{L}{1-L\bar{\rho}_w} \eta c_u \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(s) - \mathbf{y}\|_2 \|\mathbf{X}\|_F \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \right] \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2 \end{aligned} \quad (342)$$

$$= \eta C_2 \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2^2, \quad (343)$$

where

$$C_2 := c_u^2(c_a\|\mathbf{X}\|_F + c_m)^2 + c_u^2\|\mathbf{X}\|_F^2. \quad (344)$$

Furthermore, we also have

$$\begin{aligned} & (\mathbf{g} - \hat{\mathbf{y}}(\tau))^T(\hat{\mathbf{y}}(\tau) - \mathbf{y}) \\ &= (\mathbf{a}(\tau+1) - \mathbf{a}(\tau))^T \mathbf{T}(\tau)(\hat{\mathbf{y}}(\tau) - \mathbf{y}) \end{aligned} \quad (345)$$

$$= -(\eta \nabla_{\mathbf{a}} \Phi(\tau))^T \mathbf{T}(\tau)(\hat{\mathbf{y}}(\tau) - \mathbf{y}) \quad (346)$$

$$= -(\hat{\mathbf{y}}(\tau) - \mathbf{y})^T \mathbf{Z}(\tau)^T \mathbf{Z}(\tau)(\hat{\mathbf{y}}(\tau) - \mathbf{y}) \quad (347)$$

$$\leq -\eta \frac{\lambda_0}{2} \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2^2 \quad (348)$$

where we use induction  $\lambda_\tau \geq \frac{\lambda_0}{2}$ .

From (335)-(348), we obtain

$$\begin{aligned} & \Phi(\tau+1) - \Phi(\tau) \\ & \leq \frac{1}{2} \eta^2 C_1^2 \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2^2 + \eta C_2 \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2^2 - \eta \frac{\lambda_0}{2} \|\hat{\mathbf{y}}(\tau) - \mathbf{y}\|_2^2 \end{aligned} \quad (349)$$

$$= 2\Phi(\tau) \left[ \frac{1}{2} \eta^2 C_1^2 + \eta C_2 - \eta \frac{\lambda_0}{2} \right] \quad (350)$$

$$= \Phi(\tau) \left[ \eta^2 C_1^2 + 2\eta C_2 - \eta \lambda_0 \right], \quad (351)$$

which leads to

$$\Phi(\tau+1) \leq \Phi(\tau) \left[ 1 - \eta(\lambda_0 - \eta C_1^2 - 2C_2) \right] \Phi(\tau) \quad (352)$$

$$\leq (1 - \eta(\lambda_0 - 4C_2)) \Phi(\tau) \quad (353)$$

$$\leq \left( 1 - \eta \frac{\lambda_0}{2} \right) \Phi(\tau). \quad (354)$$