

---

# Theory-Inspired Path-Regularized Differential Network Architecture Search (Supplementary File)

---

Pan Zhou Caiming Xiong Richard Socher Steven C.H. Hoi  
Salesforce Research  
{pzhou, cxiong, rsocher, shoi}@salesforce.com

This supplementary document contains the technical proofs of convergence results and some additional experimental results of the NeurIPS’20 submission entitled “Theory-Inspired Path-Regularized Differential Network Architecture Search”. It is structured as follows. In Appendix A, we provides more experimental results and details, including the robustness investigation of PR-DARTS to regularization parameters, effects of group-structured sparse regularization to gate activate probability, and training algorithms and details of PR-DARTS. Appendix B summarizes the notations throughout this document and also provides the existing auxiliary theories and lemmas for subsequent analysis. Then Appendix C gives the proofs of the main results in Sec. 3, namely Theorem 1, by first introducing auxiliary theories and lemmas for subsequent analysis whose proofs are deferred to Appendix E. Next, in Appendix D we presents the results in Sec. 4, including Thoerems 2, 3 and 4. Finally, Appendix E provides the proofs for auxiliary theories and lemmas in Appendix C.

## A More Experimental Results and Details

Due to space limitation, we defer more experimental results and details to this appendix. Here we first investigate robustness of PR-DARTS to regularization parameters. Then we present effects of group-structured sparse regularization to gate activate probability, and also show the reduction cell of PR-DARTS on CIFAR10. Next, we introduce the training algorithm of PR-DARTS, and finally present more setting details of optimizers for searching architectures and retraining from scratch.

### A.1 Robustness to Regularization Parameters

Fig. 3 reports the effects of regularization parameters  $\lambda_1 \sim \lambda_3$  to the performance of PR-DARTS. Due to the high training cost, we fix two regularization parameters and then investigate the third one. From Fig. 3, one can observe that for each  $\lambda$  ( $\lambda_1$  or  $\lambda_2$  or  $\lambda_3$ ), when tuning it in a relatively large range, e.g.  $\lambda_1 \in [10^{-2}, 1]$ ,  $\lambda_2 \in [10^{-4.5}, 10^{-2.5}]$  and  $\lambda_3 \in [10^{-4}, 10^{-1.5}]$ , PR-DARTS has relatively stable performance on CIFAR10. This testifies the robustness of PR-DARTS to regularization parameters.

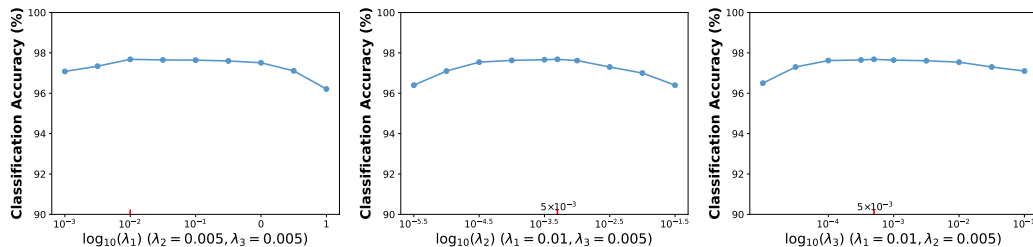


Figure 3: Effects of regularization parameters  $\lambda_1 \sim \lambda_3$  to the performance of PR-DARTS.

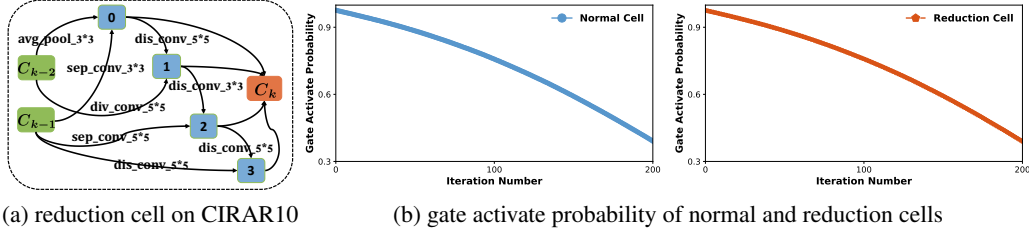


Figure 4: Visualization of search results. (a) denotes the selected reduction cell on CIRAR10. The normal cell is displayed in Fig. 1 in the manuscript. (b) shows the gate activate probability of normal cell and reduction cell in PR-DARTS.

## A.2 Effects of Group-Structured Sparse Regularization to Gate Activate Probability

Here we first display the selected reduction cell on CIRAR10 in Fig. 4 (a). The normal cell selected on CIFAR10 is displayed in Fig. 1 in the manuscript.

Next, we also report the average gate activate probability in the normal and reduction cells in Fig. 4 (b). At the beginning of the search, we initialize the activation probability of each gate to be one. This is because (1) as shown in Theorem 3, the activation probability of the gate  $g_{s,t}^{(l)}$  is  $\mathbb{P}(g_{s,t}^{(l)} \neq 0) = \Theta(\beta_{s,t}^{(l)} - \tau \ln \frac{-a}{b})$ ; (2) we set  $a = -0.1, b = 1.1, \beta_{s,t}^{(l)} = 0.5$  and initialize  $\tau = 10$  which leads to  $\mathbb{P}(g_{s,t}^{(l)} \neq 0) = \Theta(\beta_{s,t}^{(l)} - \tau \ln \frac{-a}{b}) \approx 1$ . In this way, all gates will be well explored. With along more iterations, the group structured sparsity regularization encourages competition and cooperation among all operations to improve the performance, and also prunes redundancy and unnecessary connections in the cells as well. To measure the overall sparsity of the normal cell, we compute its overall average activation probability  $\frac{1}{|\mathcal{G}|} \sum_{g_{s,t}^{(l)} \in \mathcal{G}} \mathbb{P}(g_{s,t}^{(l)} \neq 0)$ , where the gate set  $\mathcal{G}$  collects all the operation gate in the normal cell. Similarly, we can compute the average activation probability of gates in the reduction cell. As shown in Fig. 4 (b), for both normal and reduction cells, their average gate activate probability becomes smaller with along more iterations. This indicates the activation probability of the gates on redundancy and unnecessary connections becomes smaller, which means that sparsity regularizer gradually and automatically prunes redundancy and unnecessary connections which reduces the information loss of pruning at the end of search. Moreover, this sparsity regularizer defined on the whole cell can encourage global competition and cooperation of all operations in the cell, which differs from DARTS that only introduces local competition among the operations between two nodes. Actually, sparse cell also can reduce the computation cost and boost the search efficiency.

## A.3 Algorithm Framework of PR-DARTS

In this subsection, we introduce the training algorithm of PR-DARTS in details. Same as DARTS, we alternatively update the network parameter  $\mathbf{W}$  and the architecture parameter  $\beta$  via gradient descent which is detailed in Algorithm 1. For notation in Algorithm 1,  $F_{\mathcal{B}_{\text{train}}}(\mathbf{W}, \beta) = \frac{1}{|\mathcal{B}_{\text{train}}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}_{\text{train}}} f(\mathbf{W}, \beta; (\mathbf{x}, \mathbf{y}))$  denotes the training loss on mini-batch  $\mathcal{B}_{\text{train}}$ . Similarly, the loss  $F_{\mathcal{B}_{\text{val}}}(\mathbf{W}, \beta)$  denotes the validation loss on mini-batch  $\mathcal{B}_{\text{val}}$ . When we compute the gradient  $\nabla_{\beta} F_{\mathcal{B}_{\text{train}}}(\mathbf{W}, \beta)$ , we ignore the second-order Hessian to accelerate the computation which is the same as first-order DARTS.

---

### Algorithm 1 Searching Algorithm for PR-DARTS

---

**Input:** training dataset  $\mathcal{D}_{\text{train}}$  and validation dataset  $\mathcal{D}_{\text{val}}$ , mini-batch size  $b$ , learning rate  $\eta$ .  
**while** not convergence **do**  
    sample mini-batch  $\mathcal{B}_{\text{train}}$  from  $\mathcal{D}_{\text{train}}$  to update  $\mathbf{W}$  by gradient descent  $\mathbf{W} = \mathbf{W} - \eta \nabla_{\mathbf{W}} F_{\mathcal{B}_{\text{train}}}(\mathbf{W}, \beta)$ .  
    sample mini-batch  $\mathcal{B}_{\text{val}}$  from  $\mathcal{D}_{\text{val}}$  to update  $\beta$  by gradient descent  $\beta = \beta - \eta \nabla_{\beta} F_{\mathcal{B}_{\text{val}}}(\mathbf{W}, \beta)$ .  
**end while**  
**Output:**  $\beta$

---

## A.4 Algorithm Parameter Settings

**CIFAR10 and CIAFR100.** In the search phase, following DARTS, we use momentum SGD to optimize network parameter  $\mathbf{W}$ , with an initial learning rate 0.025 (annealed down to zero via cosine

decay [1]), a momentum of 0.9, and a weight decay of  $3 \times 10^{-4}$ . Architecture parameter  $\beta$  is updated by ADAM [2] with a learning rate of  $3 \times 10^{-4}$  and a weight decay of  $10^{-3}$ . For evaluation on CIFAR10 and CIFAR100, we use momentum SGD with an initial learning 0.025 (cosine decayed to zero), a momentum of 0.9, a weight decay of  $3 \times 10^{-4}$ , and gradient norm clipping parameter 5.0.

**ImageNet.** We evaluate the transfer ability of the cells selected on CIFAR10 by testing them on ImageNet. Following DARTS, we use momentum SGD with an initial learning 0.025 (cosine decayed to zero), a momentum of 0.9, a weight decay of  $3 \times 10^{-4}$ , and gradient norm clipping parameter 5.0.

## B Notation and Preliminarily

### B.1 Notations

In this document, we use  $\mathbf{X}_i^{(l)}(k)$  to denote the output  $\mathbf{X}_i^{(l)}$  of the  $i$ -th sample in the  $l$ -th layer at the  $k$ -th iteration. For brevity, we usually ignore the notation  $(k)$  and  $i$  and use  $\mathbf{X}^{(l)}$  to denote the output  $\mathbf{X}^{(l)}$  of any sample  $\mathbf{X}_i$  ( $\forall i = 1, \dots, n$ ) in the  $l$ -th layer at any iteration. We use  $\Omega = \{\mathbf{W}^{(0)}, \mathbf{W}_0^{(1)}, \mathbf{W}_0^{(2)}, \mathbf{W}_1^{(2)}, \dots, \mathbf{W}_0^{(l)}, \dots, \mathbf{W}_{l-1}^{(l)}, \dots, \mathbf{W}_0^{h-1}, \dots, \mathbf{W}_{h-2}^{(h-1)}, \mathbf{W}_0, \dots, \mathbf{W}_{h-1}\}$  to denote the set of all  $\frac{h(h+3)}{2}$  learnable matrix parameters, including the convolution parameters  $\mathbf{W}_s^{(l)}$  and the linear mapping parameters  $\mathbf{W}_s$ . Let  $\Omega_i$  denote the  $i$ -th matrix parameters in  $\Omega$ , e.g.  $\Omega_1 = \mathbf{W}^{(0)}$ . For notation simplicity, here we assume the input size is  $m \times p$  to avoid using  $\bar{m} \times \bar{p}$ . The operation  $\text{vec}(\mathbf{X})$  vectorizes the matrix  $\mathbf{X}$ .

Then we define the loss

$$F(\Omega) = \frac{1}{2n} \|\mathbf{y} - \mathbf{u}(k)\|_2^2 = \frac{1}{2n} \sum_{i=1}^n (y_i - u_i)^2 = \frac{1}{n} \sum_{i=1}^n \ell_i,$$

where  $\mathbf{u}(k) = [u_1(k); u_2(k); \dots, u_n(k)] \in \mathbb{R}^n$  denotes the prediction at the  $k$ -th iteration,  $\mathbf{y} = [y_1; y_2; \dots, y_n] \in \mathbb{R}^n$  is the labels for the  $n$  samples  $\{\mathbf{X}_i\}_{i=1}^n$ , and  $\ell_i = (y_i - u_i)^2$  denotes the individual loss of the  $i$ -th sample  $\mathbf{X}_i$ .

Then for brevity,  $\ell(\Omega)$  and  $\ell_i(\Omega)$  respectively denote the losses when feeding the input  $(\mathbf{X}, \mathbf{y})$  and  $(\mathbf{X}_i, y_i)$ . Then we denote the gradient of  $\ell(\Omega)$  with respect to all learnable parameters  $\Omega$  as

$$\nabla_{\Omega} \ell(\Omega) = \left[ \text{vec} \left( \frac{\partial \ell}{\partial \mathbf{W}^{(0)}} \right); \left\{ \text{vec} \left( \frac{\partial \ell}{\partial \mathbf{W}_s^{(l)}} \right) \right\}_{0 \leq l \leq h-1, 0 \leq s \leq l-1}; \left\{ \text{vec} \left( \frac{\partial \ell}{\partial \mathbf{W}_s} \right) \right\}_{0 \leq s \leq h-1} \right],$$

where the  $\text{vec}(\mathbf{X})$  operation vectorizes the matrix  $\mathbf{X}$  into vector. Here we also let  $\nabla_{\Omega_i} \ell(\Omega)$  denotes the gradient of  $\ell(\Omega)$  with the  $i$ -th matrix parameter, e.g.  $\nabla_{\Omega_1} \ell(\Omega) = \text{vec} \left( \frac{\partial \ell}{\partial \mathbf{W}^{(0)}} \right)$ . Therefore,  $\nabla_{\Omega} F(\Omega) = \frac{1}{n} \sum_{i=1}^n \nabla_{\Omega} \ell_i(\Omega)$  where  $\ell_i(\Omega)$  is the loss given input  $(\mathbf{X}_i, y_i)$ . In this way, we can define the Gram matrix  $\mathbf{G}(k) \in \mathbb{R}^{n \times n}$  at the  $k$ -th iteration in which its  $(i, j)$ -th entry is defined as

$$\mathbf{G}_{ij}(k) = \langle \nabla_{\Omega} \ell_i(\Omega(k)), \nabla_{\Omega} \ell_j(\Omega(k)) \rangle,$$

where  $\nabla_{\Omega} \ell_i(\Omega(k))$  denote the gradient of the loss  $\ell_i$  on the  $i$ -th sample  $(\mathbf{X}_i, y_i)$  with respect to all parameter  $\Omega$  at the  $k$ -th iteration. We often ignore the notation  $k$  and use  $\mathbf{G}$  to denote the Gram matrix that does not depend on iteration number  $k$ .

According to the definitions, we have

$$\begin{aligned} \mathbf{G}_{ij}(k) &= \langle \nabla_{\Omega} \ell_i(\Omega(k)), \nabla_{\Omega} \ell_j(\Omega(k)) \rangle = \sum_{t=1}^{\frac{h(h+3)}{2}} \langle \nabla_{\Omega_t} \ell_i(\Omega(k)), \nabla_{\Omega_t} \ell_j(\Omega(k)) \rangle \\ &= \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}^{(0)}(k)}, \frac{\partial \ell_j}{\partial \mathbf{W}^{(0)}(k)} \right\rangle + \sum_{l=1}^{h-1} \sum_{s=0}^{l-1} \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}_s^{(l)}(k)}, \frac{\partial \ell_j}{\partial \mathbf{W}_s^{(l)}(k)} \right\rangle + \sum_{s=0}^{h-1} \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}_s(k)}, \frac{\partial \ell_j}{\partial \mathbf{W}_s(k)} \right\rangle \end{aligned}$$

For brevity, we let

$$\bar{\mathbf{G}}_{ij}^0(k) = \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}^{(0)}(k)}, \frac{\partial \ell_j}{\partial \mathbf{W}^{(0)}(k)} \right\rangle, \mathbf{G}_{ij}^{ls}(k) = \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}_s^{(l)}(k)}, \frac{\partial \ell_j}{\partial \mathbf{W}_s^{(l)}(k)} \right\rangle, \mathbf{G}_{ij}^s(k) = \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}_s(k)}, \frac{\partial \ell_j}{\partial \mathbf{W}_s(k)} \right\rangle.$$

Therefore, we have

$$\mathbf{G}_{ij}(k) = \bar{\mathbf{G}}_{ij}^0(k) + \sum_{l=1}^{h-1} \sum_{s=0}^{l-1} \mathbf{G}_{ij}^{ls}(k) + \sum_{s=0}^{h-1} \mathbf{G}_{ij}^s(k), \quad \mathbf{G}(k) = \bar{\mathbf{G}}^0(k) + \sum_{l=1}^{h-1} \sum_{s=0}^{l-1} \mathbf{G}^{ls}(k) + \sum_{s=0}^{h-1} \mathbf{G}^s(k).$$

Finally, since we need to compute the gradient. Here we define an operation for computing the gradient for convolution operation. For back-propagate, we define the inverse operation of  $\Phi(\mathbf{X})$  as  $\Psi(\frac{1}{\tau}\Phi(\mathbf{X})) = \mathbf{X} \in \mathbb{R}^{m \times p}$ . For the  $(i, j)$ -th entry in  $\Psi(\mathbf{X})$ , it equals to the sum of all  $\mathbf{X}_{i,j}$  in  $\Phi(\mathbf{X})$ .

## B.2 Auxiliary Lemmas

**Lemma 1.** [3][Chebyshev's inequality] For any variable  $x$ , we have

$$\mathbb{P}(|x - \mathbb{E}[x]| \geq a) \leq \frac{\text{Var}(x)}{a^2},$$

where  $a$  is a positive constant,  $\text{Var}(x)$  denotes the variance of  $x$ .

**Lemma 2.** [4] Given a set of matrices  $\{\mathbf{A}_i, \mathbf{B}_i\}$  with proper sizes, if  $\|\mathbf{A}_i\|_2 \leq a_i$  and  $\|\mathbf{B}_i\|_2 \leq a_i$  and  $\|\mathbf{A}_i - \mathbf{B}_i\|_F \leq b_i a_i$ , we have

$$\left\| \prod_{i=1}^n \mathbf{A}_i - \prod_{i=1}^n \mathbf{B}_i \right\|_F \leq \left( \sum_{i=1}^n b_i \right) \prod_{i=1}^n a_i.$$

**Lemma 3.** [5][Cauchy Interlace Theorem] Let  $\mathbf{A}$  be a Hermitian matrix of order  $n$  and let  $\mathbf{B}$  be a principal submatrix of  $\mathbf{A}$  of order  $n-1$ . If  $\lambda_n \leq \lambda_{n-1} \leq \dots \leq \lambda_1$  lists the eigenvalues of  $\mathbf{A}$  and  $\mu_n \leq \mu_{n-1} \leq \dots \leq \mu_2$  the eigenvalues of  $\mathbf{B}$ , then  $\lambda_n \leq \mu_n \leq \lambda_{n-1} \leq \mu_{n-1} \leq \dots \leq \lambda_2 \leq \mu_2 \leq \lambda_1$ .

**Lemma 4.** [6][Chi-Square Variable Bound] Let  $x$  be chi-square variable with  $n$  degree of freedom. Then for any  $t > 0$ , it holds

$$\mathbb{P}(x - n \geq 2\sqrt{nt} + 2t) \leq \exp(-t), \quad \text{and} \quad \mathbb{P}(x - n \leq -2\sqrt{nt}) \leq \exp(-t).$$

**Lemma 5.** [4] Suppose  $\sigma$  is analytic and not a polynomial function. Consider data  $\{\mathbf{X}_{i=1}^n\}_{i=1}^n$  are not parallel, namely  $\text{vec}(\mathbf{X}_i) \notin \text{span}(\text{vec}(\mathbf{X}_j))$  for all  $i \neq j$ , Then the smallest eigenvalue the matrix  $\mathbf{G}$  which is defined as

$$\mathbf{G}(\mathbf{X})_{ij} = \mathbb{E}_{\mathbf{W} \sim \mathcal{N}(0, \mathbf{I})} \sigma(\langle \mathbf{W}, \mathbf{X}_i \rangle) \sigma(\langle \mathbf{W}, \mathbf{X}_j \rangle)$$

is larger than zero, namely  $\lambda_{\min}(\mathbf{G}) > 0$ .

**Lemma 6.** [4] Suppose  $\sigma$  is analytic and not a polynomial function. Consider data  $\{\mathbf{X}_{i=1}^n\}_{i=1}^n$  are not parallel, namely  $\text{vec}(\mathbf{X}_i) \notin \text{span}(\text{vec}(\mathbf{X}_j))$  for all  $i \neq j$ , Then the smallest eigenvalue the matrix  $\mathbf{G}$  which is defined as

$$\mathbf{G}(\mathbf{X})_{ij} = \mathbb{E}_{\mathbf{W} \sim \mathcal{N}(0, \mathbf{I})} \sigma'(\langle \mathbf{W}, \mathbf{X}_i \rangle) \sigma'(\langle \mathbf{W}, \mathbf{X}_j \rangle)$$

is larger than zero, namely  $\lambda_{\min}(\mathbf{G}) > 0$ .

**Lemma 7.** [4] Suppose the activation function  $\sigma(\cdot)$  satisfies Assumption 1. Suppose there exists  $c > 0$  such that

$$\mathbf{A} = \begin{bmatrix} a_1^2 & \rho a_1 b_1 \\ \rho a_1 b_1 & b_1^2 \end{bmatrix} \succ 0, \quad \mathbf{B} = \begin{bmatrix} a_2^2 & \rho a_2 b_2 \\ \rho a_2 b_2 & b_2^2 \end{bmatrix} \succ 0,$$

where the parameter satisfies  $1/c \leq x \leq c$  in which  $x$  could be  $a_1, a_2, b_1, b_2$ . Let  $g(\mathbf{A}) = \mathbb{E}_{(u,v) \sim \mathcal{N}(0, \mathbf{A})} \sigma(u) \sigma(v)$ . Then we have

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

where  $C$  is a constant that only depends on  $c$  and the Lipschitz and smooth parameter of  $\sigma(\cdot)$ .

## C Proofs of Results in Sec. 3

### C.1 Proof of Theorem 1

Suppose Assumptions 1, 2 and 3 hold. To prove our main results, namely the results in Theorem 1, we have two steps. In the first step, from Lemma 21, we have that if  $m$  and  $\eta$  satisfy

$$m \geq \frac{c'_m c^2 \rho k_c^2 c_{w0} \mu^2}{\lambda^2 n}, \quad \eta \leq \frac{c'_\eta \lambda}{\sqrt{m} \mu^4 h^3 k_c^2 c^4},$$

where  $c'_m$  and  $c'_\eta$  are two constants,  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c c_{w0}})^h$ ,  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . Then with probability at least  $1 - \delta/2$  we have

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq \left(1 - \frac{\eta\lambda_{\min}(\mathbf{G}(0))}{4}\right) \|\mathbf{y} - \mathbf{u}(k-1)\|_2^2.$$

where  $k$  denotes the iteration number,  $\lambda_{\min}(\mathbf{G}(0))$  denotes the smallest eigenvalue of the Gram matrix  $\mathbf{G}(0)$  at the initialization. For this part, we prove it in Appendix C.3.

In the second step, we will prove that the smallest eigenvalue of can be lower bounded. Specifically, we prove this results in Lemma 24: if  $m \geq \frac{c_4\mu^2 p^2 n^2 \log(n/\delta)}{\lambda^2}$ , it holds that with probability at least  $1 - \delta/2$ , the smallest eigenvalue the matrix  $\mathbf{G}$  satisfies

$$\lambda_{\min}(\mathbf{G}(0)) \geq \frac{3c_\sigma}{4} \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) \lambda_{\min}(\mathbf{K}).$$

where  $\lambda = 3c_\sigma \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) \lambda_{\min}(\mathbf{K})$ ,  $c_\sigma$  is a constant that only depends on  $\sigma$  and the input data,  $\lambda_{\min}(\mathbf{K}) = \min_{i,j} \lambda_{\min}(\mathbf{K}_{ij})$  is larger than zero in which  $\lambda_{\min}(\mathbf{K}_{ij})$  is the the smallest eigenvalue of  $\mathbf{K}_{ij} = \begin{bmatrix} \mathbf{X}_i^\top \mathbf{X}_j & \mathbf{X}_i^\top \mathbf{X}_j \\ \mathbf{X}_j^\top \mathbf{X}_i & \mathbf{X}_j^\top \mathbf{X}_j \end{bmatrix}$ . Appendix C.4 provides the proof for this result.

Finally, we combine these results in the above two steps and can obtain that if  $m \geq \frac{c_m\mu^2}{\lambda^2} [\rho p^2 n^2 \log(n/\delta) + c^2 k_c^2 c_{w0}^2/n]$  and  $\eta \leq \frac{c_\eta\lambda}{\sqrt{m}\mu^4 h^3 k_c^2 c^4}$ , where  $c_{w0}, c_m, c_\eta$  are constants, with probability at least  $(1 - \delta/2)^2 > 1 - \delta$ , we have

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq (1 - \eta\lambda/4) \|\mathbf{y} - \mathbf{u}(k-1)\|_2^2 \quad (\forall k \geq 1),$$

where  $\lambda = \frac{3c_\sigma}{4} \lambda_{\min}(\mathbf{K}) \sum_{s=0}^{h-2} (\alpha_{s,3}^{(h-1)})^2 \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2$ , the positive constant  $c_\sigma$  only depends on  $\sigma$  and input data. On the other hand, we have

$$F_{\text{train}}(\mathbf{W}(k+1), \beta) = \frac{1}{2n} \|\mathbf{y} - \mathbf{u}(k+1)\|_2^2,$$

then we can obtain the desired results in Theorem 1. Please refer to the proof details in Appendix C.3 and C.4 for the above two steps respectively.

Note that our proof framework is similar to [4]. But there are essential differences. The main difference is that here our network architecture is much complex (e.g. each layer connects all the previous layers) and each edge in our network also involves more operations, including zero operation, skip operation and convolution operation, which requires bounding many terms in this work differently and more elaborately.

For the following proofs, Appendix C.2 provides the auxiliary lemmas for the proofs for Step 1 and Step 2. Then Appendix C.3 and C.4 respectively present the proof details in Step 1 and Step 2.

## C.2 Auxiliary Lemmas

**Lemma 8.** *The gradient of the loss  $\ell = \frac{1}{2}(u - y)^2$  with parameter and temporary output can be written as follows:*

$$\frac{\partial \ell}{\partial \mathbf{X}^{(l)}} = (u - y)\mathbf{W}_l + \sum_{s=l+1}^{h-1} \left( \alpha_{l,2}^{(s)} \frac{\partial \ell}{\partial \mathbf{X}^{(s)}} + \alpha_{l,3}^{(s)} \tau \Psi \left( (\mathbf{W}_l^{(s)})^\top \left( \sigma' \left( \mathbf{W}_l^{(s)} \Phi(\mathbf{X}^{(l)}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(s)}} \right) \right) \right),$$

( $0 \leq l \leq h-1, 0 \leq s \leq l-1$ ),

$$\frac{\partial \ell}{\partial \mathbf{X}} = \tau \Psi \left( (\mathbf{W}^{(0)})^\top \left( \sigma' \left( \mathbf{W}^{(0)} \Phi(\mathbf{X}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(0)}} \right) \right) \in \mathbb{R}^{m \times p},$$

$$\frac{\partial \ell}{\partial \mathbf{W}_s^{(l)}} = \alpha_{s,3}^{(l)} \tau \Phi(\mathbf{X}^{(s)}) \left( \sigma' \left( \mathbf{W}_s^{(l)} \Phi(\mathbf{X}^{(s)}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(l)}} \right)^\top \in \mathbb{R}^{m \times p} \quad (0 \leq l \leq h-1, 0 \leq s \leq l-1),$$

$$\frac{\partial \ell}{\partial \mathbf{W}^{(0)}} = \tau \Phi(\mathbf{X}) \left( \sigma' \left( \mathbf{W}^{(0)} \Phi(\mathbf{X}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(0)}} \right)^\top \in \mathbb{R}^{m \times p},$$

$$\frac{\partial \ell}{\partial \mathbf{W}_s} = (u - y)\mathbf{X}^{(s)} \in \mathbb{R}^{m \times p},$$

where  $\odot$  denotes the dot product,  $\frac{\partial \ell}{\partial \mathbf{X}^{(l)}} \in \mathbb{R}^{m \times p}$ .

See its proof in Appendix E.1.

**Lemma 9.** *The gradient of the network output  $u$  with respect to the output and convolution parameter can be written as follows:*

$$\begin{aligned}\frac{\partial u}{\partial \mathbf{X}^{(l)}} &= \mathbf{W}_l + \sum_{s=l+1}^{h-1} \left( \alpha_{l,2}^{(s)} \frac{\partial u}{\partial \mathbf{X}^{(s)}} + \alpha_{l,3}^{(s)} \tau \Psi \left( (\mathbf{W}_l^{(s)})^\top \left( \sigma' \left( \mathbf{W}_l^{(s)} \Phi(\mathbf{X}^{(l)}) \right) \odot \frac{\partial u}{\partial \mathbf{X}^{(s)}} \right) \right) \right), \\ &\quad (0 \leq l \leq h-1, 0 \leq s \leq l-1), \\ \frac{\partial u}{\partial \mathbf{X}} &= \tau \Psi \left( (\mathbf{W}^{(0)})^\top \left( \sigma' \left( \mathbf{W}^{(0)} \Phi(\mathbf{X}) \right) \odot \frac{\partial u}{\partial \mathbf{X}^{(0)}} \right) \right) \in \mathbb{R}^{m \times p}, \\ \frac{\partial u}{\partial \mathbf{W}_s^{(l)}} &= \alpha_{s,3}^{(l)} \tau \Phi(\mathbf{X}^{(s)}) \left( \sigma' \left( \mathbf{W}_s^{(l)} \Phi(\mathbf{X}^{(s)}) \right) \odot \frac{\partial u}{\partial \mathbf{X}^{(l)}} \right)^\top \in \mathbb{R}^{m \times p} \quad (0 \leq l \leq h-1, 0 \leq s \leq l-1), \\ \frac{\partial u}{\partial \mathbf{W}^{(0)}} &= \tau \Phi(\mathbf{X}) \left( \sigma' \left( \mathbf{W}^{(0)} \Phi(\mathbf{X}) \right) \odot \frac{\partial u}{\partial \mathbf{X}^{(0)}} \right)^\top \in \mathbb{R}^{m \times p}, \\ \frac{\partial u}{\partial \mathbf{W}_s} &= \mathbf{X}^{(s)} \in \mathbb{R}^{m \times p}, \quad (0 \leq s \leq h-1),\end{aligned}$$

where  $\odot$  denotes the dot product and  $\frac{\partial u}{\partial \mathbf{X}^{(l)}} \in \mathbb{R}^{m \times p}$ .

See its proof in Appendix E.2.

**Lemma 10.** *Suppose Assumptions 1, 2 and 3 hold. Given a constant  $\delta \in (0, 1)$ , assume  $m \geq \frac{16c_1np^2}{c_2\delta}$ , where  $c_1 = \sigma^4(0) + 4|\sigma^3(0)|\mu\sqrt{2/\pi} + 8|\sigma(0)|\mu^3\sqrt{2/\pi} + 32\mu^4$  and  $c = \mathbb{E}_{\omega \sim \mathcal{N}(0, \frac{1}{\sqrt{p}})} [\sigma^2(\omega)]$ . Suppose  $\mathbf{W}_s^{(l)}(0) \leq \sqrt{m}c_{w0} \forall 0 \leq l \leq h, 0 \leq s \leq l-1$ . Then with probability at least  $1 - \delta/4$ , we have*

$$\frac{1}{c_{x0}} \leq \|\mathbf{X}^{(l)}(0)\|_F \leq c_{x0}.$$

where  $c_{x0} \geq 1$  is a constant.

See its proof in Appendix E.3.

**Lemma 11.** *Suppose Assumptions 1, 2 and 3 hold. Assume  $\|\mathbf{W}_s^{(l)}(0)\|_2 \leq \sqrt{m}c_{w0}$ ,  $\|\mathbf{W}_s^{(l)}(k) - \mathbf{W}_s^{(l)}(0)\|_F \leq \sqrt{m}r$ . Then for  $\forall l$ , we have*

$$\begin{aligned}\|\mathbf{X}^{(l)}(k) - \mathbf{X}^{(l)}(0)\|_F &\leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right)^l \mu \sqrt{k_c} r, \\ \left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k)) - \mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F &\leq \frac{1}{\alpha_3} \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right)^l \sqrt{k_c m} r,\end{aligned}$$

where  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ , and  $c_{x0} \geq 1$  is given in Lemma 10.

See its proof in Appendix E.4.

**Lemma 12.** *Suppose Assumptions 1, 2 and 3 hold. Assume  $\frac{1}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_F = c_y$  and  $\|\mathbf{W}_h(t)\|_F \leq c_u$ ,  $\|\mathbf{W}_l^{(s)}(t) - \mathbf{W}_l^{(s)}(0)\|_F \leq \sqrt{m}r$ , and  $\|\mathbf{W}_l^{(s)}(0)\|_F \leq \sqrt{m}c_{w0}$ . Then for  $\forall l$ , we have*

$$\frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)} \right\|_F \leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right)^l c_y c_u,$$

where  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ .

See its proof in Appendix E.5.

**Lemma 13.** *Suppose Assumptions 1, 2 and 3 hold. Assume  $\|\mathbf{y} - \mathbf{u}(t)\|_2^2 \leq (1 - \frac{\eta\lambda}{2})^t \|\mathbf{y} - \mathbf{u}(0)\|_2^2$  holds for  $t = 1, \dots, k$ . Then by setting*

$$\tilde{r} = \frac{8c_{x0} \|\mathbf{y} - \mathbf{u}(0)\|_2}{\lambda \sqrt{mn}} \max \left( 1, 2 \left( 1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0} \right)^l \alpha_{s,3}^{(l)} \mu \sqrt{k_c} c_{w0} \right) \leq c_{w0},$$

we have that for any  $s = 1, \dots, k+1$ ,

$$\begin{aligned} \|\mathbf{W}^{(0)}(t) - \mathbf{W}^{(0)}(0)\|_F &\leq \sqrt{m}\tilde{r}, \quad \|\mathbf{W}_s^{(l)}(t) - \mathbf{W}_s^{(l)}(0)\|_F \leq \sqrt{m}\tilde{r}, \quad \|\mathbf{W}_s(t) - \mathbf{W}_s(0)\|_F \leq \sqrt{m}\tilde{r}, \\ \|\mathbf{W}^{(0)}(t+1) - \mathbf{W}^{(0)}(t)\|_F &= \eta \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}^{(0)}(t)} \right\|_F \leq \frac{4c\eta\mu c_{x0}c_{w0}\sqrt{k_c}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2, \\ \|\mathbf{W}_s^{(l)}(t+1) - \mathbf{W}_s^{(l)}(t)\|_F &= \eta \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}_s^{(l)}(t)} \right\|_F \leq \frac{4c\eta\alpha_{s,3}^{(l)}\mu c_{x0}c_{w0}\sqrt{k_c}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2, \\ \|\mathbf{W}_s(t+1) - \mathbf{W}_s(t)\|_F &= \eta \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}_s(t)} \right\|_F \leq \frac{2\eta c_{x0}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2, \end{aligned}$$

where  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ .

See its proof in Appendix E.6.

**Lemma 14.** Suppose Assumptions 1, 2 and 3 hold. Then we have

$$\begin{aligned} &\left\| \mathbf{X}^{(l)}(k+1) - \mathbf{X}^{(l)}(k) \right\|_F \\ &\leq \left(1 + \alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu\right)^l \left(1 + \frac{2(\alpha_3)^2 c_{x0}}{(\alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu)\sqrt{n}}\right) \frac{4c\tau\eta\mu^2 c_{x0}c_{w0}k_c}{\sqrt{n}} \|\mathbf{u}(k) - \mathbf{y}\|_F, \end{aligned}$$

where  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ .

See its proof in Appendix E.7.

**Lemma 15.** Suppose Assumptions 1, 2 and 3 hold. Then we have

$$\left\| \mathbf{W}^{(0)}(k) \right\|_F \leq 2\sqrt{m}c_{w0}, \quad \left\| \mathbf{W}_s^{(l)}(k) \right\|_F \leq 2\sqrt{m}c_{w0}, \quad \|\mathbf{W}_s(k)\|_F \leq 2\sqrt{m}c_{w0}.$$

If  $\tilde{r}$  in Lemma 13 satisfies  $\tilde{r} \leq \frac{c_{x0}}{(1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l \mu\sqrt{k_c}}$  which can be achieved by using large  $m$ , then we have

$$\left\| \mathbf{X}_i^{(l)}(k) \right\|_F \leq 2c_{x0},$$

where  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ .

See its proof in Appendix E.8.

**Lemma 16.** Suppose Assumptions 1, 2 and 3 hold. Then we have

$$\|\mathbf{X}_i^{(0)}(k) - \mathbf{X}_i^{(0)}(0)\|_F \leq \mu\sqrt{k_c}\tilde{r}, \quad \|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0)\|_F \leq c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c}\tilde{r},$$

where  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . Here  $\tilde{r}$  is given in Lemma 13.

See its proof in Appendix E.9.

**Lemma 17.** Suppose Assumptions 1, 2 and 3 hold.

$$|u_i(k) - u_i(0)| \leq 2\sqrt{m}h \left( c_{x0} + c_{w0}c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c} \right) \tilde{r},$$

where  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . Here  $\tilde{r}$  is given in Lemma 13. Besides, we have

$$\left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(0)} \right\|_F \leq c_1 c \alpha_3 c_{w0}^2 c_{x0} \rho k_c m \tilde{r},$$

where  $c_1$  is a constant.

See its proof in Appendix E.10.

**Lemma 18.** Suppose Assumption 2 holds. Then with probability at least  $1 - \delta/4$ , it holds

$$\begin{cases} \|\mathbf{W}^0\|_F \leq \sqrt{m}c_{w0}, \\ \|\mathbf{W}_s^{(l)}(0)\|_F \leq \sqrt{m}c_{w0} \quad (\forall 0 \leq l \leq h-1, 0 \leq s \leq l-1), \\ \|\mathbf{W}_s(0)\|_F \leq \sqrt{m}c_{w0} \quad (\forall 0 \leq s \leq h-1). \end{cases}$$

See its proof in Appendix E.11.

### C.3 Step 1 Linear Convergence of $\|\mathbf{y} - \mathbf{u}(k)\|_2^2$

Here we first present our results and then provides their proofs.

**Lemma 19.** *Suppose Assumptions 1, 2 and 3 hold. If  $m$  and  $\eta$  satisfy*

$$\begin{cases} m \geq \frac{c_1 \rho k_c^2 c_{w0}^2 \|\mathbf{y} - \mathbf{u}(0)\|_2^2}{\lambda^2 n} (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^{2h}, \\ \eta \leq \frac{c_2 \lambda}{\sqrt{m} \mu^4 c_{w0}^2 c_{x0}^2 h^3 k_c^2 (1 + \alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu)^{4h}}, \end{cases}$$

where  $c_1$  and  $c_2$  are two constants and  $\lambda$  is smallest eigenvalue of the Gram matrix  $\mathbf{G}(t)$  ( $t = 1, \dots, k-1$ ), then with probability at least  $1 - \delta/2$  we have

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq \left(1 - \frac{\eta \lambda}{2}\right) \|\mathbf{y} - \mathbf{u}(k-1)\|_2^2 \leq \left(1 - \frac{\eta \lambda}{2}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|_2^2.$$

See its proof in Appendix C.3.1.

**Lemma 20.** *Suppose Assumptions 1, 2 and 3 hold. If  $m$  satisfy*

$$m \geq \frac{c_3 \alpha_3^2 \mu^2 k_c c_{x0}^2 c^2}{\lambda^2 n},$$

where  $c_3$  is a constant,  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^h$ ,  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ , then we have

$$\|\mathbf{G}(k) - \mathbf{G}(0)\|_2 \leq \frac{\eta \lambda_{\min}(\mathbf{G}(0))}{2},$$

where  $\lambda_{\min}(\mathbf{G}(0))$  is the smallest eigenvalue of  $\mathbf{G}(0)$ .

See its proof in Appendix C.3.2.

**Lemma 21.** *Suppose Assumptions 1, 2 and 3 hold. If  $m$  and  $\eta$  satisfy*

$$\begin{cases} m \geq \frac{c'_m c^2 \rho k_c^2 c_{w0}^2 \mu^2}{\lambda^2 n}, \\ \eta \leq \frac{c'_\eta \lambda}{\sqrt{m} \mu^4 h^3 k_c^2 c^4}, \end{cases}$$

where  $c_m$  and  $c_\eta$  are two constants,  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^h$ ,  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . Then with probability at least  $1 - \delta$  we have

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq \left(1 - \frac{\eta \lambda_{\min}(\mathbf{G}(0))}{4}\right) \|\mathbf{y} - \mathbf{u}(k-1)\|_2^2 \leq \left(1 - \frac{\eta \lambda_{\min}(\mathbf{G}(0))}{4}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|_2^2.$$

See its proof in Appendix C.3.3.

#### C.3.1 Proof of Lemma 19

*Proof.* Here we use mathematical induction to prove the result. For  $k = 0$ , the results in Theorem 19 holds. Then we assume for  $j = 1, \dots, k$ , it holds

$$\|\mathbf{y} - \mathbf{u}(j)\|_2^2 \leq \left(1 - \frac{\eta \lambda}{2}\right) \|\mathbf{y} - \mathbf{u}(j-1)\|_2^2 \leq \left(1 - \frac{\eta \lambda}{2}\right)^j \|\mathbf{y} - \mathbf{u}(0)\|_2^2 \quad (j = 1, \dots, k).$$

Then we need to prove  $j = k+1$  still holds. Our proof has four steps. In the first step, we establish the relation between  $\|\mathbf{y} - \mathbf{u}(j)\|_2^2 \leq \|\mathbf{y} - \mathbf{u}(j)\|_2^2 + H_1 + H_2$ . Then in the second, third and fourth steps, we bound the terms  $H_1, H_2, H_3$  respectively. Finally, we combine results to obtain the desired result.

**Step 1. Establishing relation between  $\|\mathbf{y} - \mathbf{u}(j)\|_2^2 \leq \|\mathbf{y} - \mathbf{u}(j)\|_2^2 + H_1 + H_2 + H_3$ .**

According to the definition, we can obtain

$$\begin{aligned} \|\mathbf{y} - \mathbf{u}(k+1)\|_2^2 &= \|\mathbf{y} - \mathbf{u}(k) + \mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2 \\ &= \|\mathbf{y} - \mathbf{u}(k)\|_2^2 + 2\langle \mathbf{y} - \mathbf{u}(k), \mathbf{u}(k) - \mathbf{u}(k+1) \rangle + \|\mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2. \end{aligned}$$



Then for brevity,  $\ell(\boldsymbol{\Omega})$  and  $\ell_i(\boldsymbol{\Omega})$  respectively denote the losses when feeding the input  $(\mathbf{X}, \mathbf{y})$  and  $(\mathbf{X}_i, y_i)$ . Then as introduced in Sec. B, we denote the gradient of  $\ell(\boldsymbol{\Omega})$  with respect to all learnable parameters  $\boldsymbol{\Omega}$  as

$$\nabla_{\boldsymbol{\Omega}} \ell(\boldsymbol{\Omega}) = \left[ \text{vec} \left( \frac{\partial \ell}{\partial \mathbf{W}^{(0)}} \right); \left\{ \text{vec} \left( \frac{\partial \ell}{\partial \mathbf{W}_s^{(l)}} \right) \right\}_{0 \leq l \leq h-1, 0 \leq s \leq l-1}; \left\{ \text{vec} \left( \frac{\partial \ell}{\partial \mathbf{W}_s} \right) \right\}_{0 \leq s \leq h-1} \right].$$

Based on the above definitions, when we use gradient descent algorithm to update the variables with learning rate  $\eta$ , we have

$$\begin{aligned} u_i(k+1) - u_i(k) &= u_i(\boldsymbol{\Omega}(k) - \eta \nabla_{\boldsymbol{\Omega}} F(\boldsymbol{\Omega}(k))) - u_i(\boldsymbol{\Omega}(k)) \\ &= - \int_{t=0}^{\eta} \langle \nabla_{\boldsymbol{\Omega}} F(\boldsymbol{\Omega}(k)), \nabla_{\boldsymbol{\Omega}} u_i(\boldsymbol{\Omega}(k) - s \nabla_{\boldsymbol{\Omega}} F(\boldsymbol{\Omega}(k))) \rangle dt = \boldsymbol{\Delta}_1^i(k) + \boldsymbol{\Delta}_2^i(k), \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\Delta}_1^i(k) &= - \int_{t=0}^{\eta} \langle \nabla_{\boldsymbol{\Omega}} F(\boldsymbol{\Omega}(k)), \nabla_{\boldsymbol{\Omega}} u_i(\boldsymbol{\Omega}(k)) \rangle dt \\ \boldsymbol{\Delta}_2^i(k) &= \int_{t=0}^{\eta} \langle \nabla_{\boldsymbol{\Omega}} F(\boldsymbol{\Omega}(k)), \nabla_{\boldsymbol{\Omega}} u_i(\boldsymbol{\Omega}(k)) - \nabla_{\boldsymbol{\Omega}} u_i(\boldsymbol{\Omega}(k) - t \nabla_{\boldsymbol{\Omega}} F(\boldsymbol{\Omega}(k))) \rangle dt. \end{aligned}$$

Then we define two important notations:

$$\boldsymbol{\Delta}_1(k) = [\boldsymbol{\Delta}_1^1(k); \boldsymbol{\Delta}_1^2(k); \dots; \boldsymbol{\Delta}_1^n(k)] \in \mathbb{R}^n, \quad \boldsymbol{\Delta}_2(k) = [\boldsymbol{\Delta}_2^1(k); \boldsymbol{\Delta}_2^2(k); \dots; \boldsymbol{\Delta}_2^n(k)] \in \mathbb{R}^n.$$

In this way, we have  $\mathbf{u}(k+1) - \mathbf{u}(k) = \boldsymbol{\Delta}_1(k) + \boldsymbol{\Delta}_2(k)$ . Now we consider

$$\begin{aligned} \boldsymbol{\Delta}_1^i(k) &= - \int_{s=0}^{\eta} \langle \nabla_{\boldsymbol{\Omega}} F(\boldsymbol{\Omega}(k)), \nabla_{\boldsymbol{\Omega}} u_i(\boldsymbol{\Omega}(k)) \rangle \\ &= - \eta \langle \nabla_{\boldsymbol{\Omega}} F(\boldsymbol{\Omega}(k)), \nabla_{\boldsymbol{\Omega}} u_i(\boldsymbol{\Omega}(k)) \rangle \\ &= - \frac{\eta}{n} \sum_{j=1}^n (y_j - u_j) \langle \nabla_{\boldsymbol{\Omega}} u_j(\boldsymbol{\Omega}(k)), \nabla_{\boldsymbol{\Omega}} u_i(\boldsymbol{\Omega}(k)) \rangle \\ &= - \frac{\eta}{n} \sum_{j=1}^n (y_j - u_j) \sum_{t=1}^{(h+1)(\frac{h}{2}+1)} \langle \nabla_{\boldsymbol{\Omega}_t} u_j(\boldsymbol{\Omega}(k)), \nabla_{\boldsymbol{\Omega}_t} u_i(\boldsymbol{\Omega}(k)) \rangle. \end{aligned}$$

Let  $\mathbf{G}_{ij}^t(k) = \langle \nabla_{\boldsymbol{\Omega}_t} u_j(\boldsymbol{\Omega}(k)), \nabla_{\boldsymbol{\Omega}_t} u_i(\boldsymbol{\Omega}(k)) \rangle$ . In this way, we have  $\mathbf{G}(k) = \sum_{t=1}^{(h+1)(\frac{h}{2}+1)} \mathbf{G}^t$ . Then  $\boldsymbol{\Delta}_1(k)$  can be formulated as follows:

$$\boldsymbol{\Delta}_1(k) = -\eta \mathbf{G}(k)(\mathbf{u}(k) - \mathbf{y}).$$

In this way, we can compute

$$\begin{aligned} 2\langle \mathbf{y} - \mathbf{u}(k), \mathbf{u}(k) - \mathbf{u}(k+1) \rangle &= -2\langle \mathbf{y} - \mathbf{u}(k), \boldsymbol{\Delta}_1(k) + \boldsymbol{\Delta}_2(k) \rangle \\ &= -2\eta(\mathbf{u}(k) - \mathbf{y})^\top \mathbf{G}(k)(\mathbf{u}(k) - \mathbf{y}) - 2\langle \mathbf{y} - \mathbf{u}(k), \boldsymbol{\Delta}_2(k) \rangle \end{aligned}$$

Therefore, we can decompose  $\|\mathbf{y} - \mathbf{u}(k+1)\|_2^2$  into

$$\begin{aligned} &\|\mathbf{y} - \mathbf{u}(k+1)\|_2^2 \\ &= \|\mathbf{y} - \mathbf{u}(k)\|_2^2 + 2\langle \mathbf{y} - \mathbf{u}(k), \mathbf{u}(k) - \mathbf{u}(k+1) \rangle + \|\mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2 \\ &= \|\mathbf{y} - \mathbf{u}(k)\|_2^2 - 2\eta(\mathbf{u}(k) - \mathbf{y})^\top \mathbf{G}(k)(\mathbf{u}(k) - \mathbf{y}) - 2\langle \mathbf{y} - \mathbf{u}(k), \boldsymbol{\Delta}_2(k) \rangle + \|\mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2 \\ &\leq \|\mathbf{y} - \mathbf{u}(k)\|_2^2 - 2\eta(\mathbf{u}(k) - \mathbf{y})^\top \mathbf{G}(k)(\mathbf{u}(k) - \mathbf{y}) + 2\|\mathbf{y} - \mathbf{u}(k)\|_2 \|\boldsymbol{\Delta}_2(k)\|_2 + \|\mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2. \end{aligned} \tag{9}$$

Let  $H_1 = -2\eta(\mathbf{u}(k) - \mathbf{y})^\top \mathbf{G}(k)(\mathbf{u}(k) - \mathbf{y})$ ,  $H_2 = 2\|\mathbf{y} - \mathbf{u}(k)\|_2 \|\boldsymbol{\Delta}_2(k)\|_2$  and  $H_3 = \|\mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2$ . The remaining task is to upper bound  $H_1 \sim H_3$ .

**Step 2. Bound of  $H_1$ .**

To bound  $H_1$ , we can easily to bound it as follows:

$$H_1 = -2\eta(\mathbf{u}(k) - \mathbf{y})^\top \mathbf{G}(k)(\mathbf{u}(k) - \mathbf{y}) \leq -2\eta\lambda\|\mathbf{u}(k) - \mathbf{y}\|_2^2,$$

where  $\lambda = \min_k \lambda_{\min}(\mathbf{G}(k))$ .

**Step 3. Bound of  $H_2$ .**

In this step, we aim to bound  $H_2 = 2\|\mathbf{y} - \mathbf{u}(k)\|_2\|\Delta_2(k)\|_2$  by bounding  $\|\Delta_2^i(k)\|_2$ . According to the definition, we have

$$\begin{aligned} \Delta_2^i(k) &= \int_{t=0}^{\eta} \langle \nabla_{\Omega} F(\Omega(k)), \nabla_{\Omega} u_i(\Omega(k)) - \nabla_{\Omega} u_i(\Omega(k) - s\nabla_{\Omega} F(\Omega(k))) \rangle dt \\ &\leq \eta \max_{t \in [0, \eta]} \|\nabla_{\Omega} F(\Omega(k))\|_F \|\nabla_{\Omega} u_i(\Omega(k)) - \nabla_{\Omega} u_i(\Omega(k) - t\nabla_{\Omega} F(\Omega(k)))\|_F. \end{aligned}$$

In this way, we need to bound  $\max_{t \in [0, \eta]} \|\nabla_{\Omega} u_i(\Omega(k)) - \nabla_{\Omega} u_i(\Omega(k) - t\nabla_{\Omega} F(\Omega(k)))\|_F$  and  $\|\nabla_{\Omega} F(\Omega(k))\|_F$ .

**Step 3.1 Bound of  $\|\nabla_{\Omega} F(\Omega(k))\|_F$  in  $H_2$ .** According to the definition, we have

$$\begin{aligned} \|\nabla_{\Omega} F(\Omega(k))\|_F &\leq \sum_{t=1}^{(h+1)(h/2+1)} \|\nabla_{\Omega_t} F(\Omega(k))\|_F \\ &= \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}^{(0)}(k)} \right\|_F + \sum_{l=0}^{h-1} \sum_{s=0}^{l-1} \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}_s^{(l)}(k)} \right\|_F + \sum_{s=0}^{h-1} \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}_s(k)} \right\|_F \\ &\stackrel{\textcircled{1}}{\leq} \left( h + 2c\mu c_{w0} \sqrt{k_c} \left( 1 + \sum_{l=0}^{h-1} \sum_{s=0}^{l-1} \alpha_{s,3}^{(l)} \right) \right) \frac{2c_{x0}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2, \end{aligned}$$

where  $\textcircled{1}$  holds by using Lemma 13 with  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l$ ,  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$  since Lemma 13 proves

$$\begin{aligned} \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}^{(0)}(t)} \right\|_F &\leq \frac{4c\mu c_{x0} c_{w0} \sqrt{k_c}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2, \quad \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}_s^{(l)}(t)} \right\|_F \leq \frac{4c\alpha_{s,3}^{(l)} \mu c_{x0} c_{w0} \sqrt{k_c}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2, \\ \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}_s(t)} \right\|_F &\leq \frac{2c_{x0}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2, \end{aligned}$$

**Step 3.2 Bound of  $\|\nabla_{\Omega} u_i(\Omega(k)) - \nabla_{\Omega} u_i(\Omega(k) - t\nabla_{\Omega} F(\Omega(k)))\|_F$  in  $H_2$ .**

For brevity, let  $\Omega(k, t) = \Omega(k) - t\nabla_{\Omega} F(\Omega(k))$ . In this way, we can bound

$$\begin{aligned} \|\nabla_{\Omega} u_i(\Omega(k)) - \nabla_{\Omega} u_i(\Omega(k, t))\|_F &\leq \sum_{o=1}^{(h+1)(h/2+1)} \|\nabla_{\Omega_o} u_i(\Omega(k)) - \nabla_{\Omega_o} u_i(\Omega(k, s))\|_F \\ &= \left\| \frac{\partial u_i}{\partial \mathbf{W}^{(0)}(k)} - \frac{\partial u_i}{\partial \mathbf{W}^{(0)}(k, t)} \right\|_F + \sum_{l=0}^{h-1} \sum_{s=0}^{l-1} \left\| \frac{\partial u_i}{\partial \mathbf{W}_s^{(l)}(k)} - \frac{\partial u_i}{\partial \mathbf{W}_s^{(l)}(k, t)} \right\|_F + \sum_{s=0}^{h-1} \left\| \frac{\partial u_i}{\partial \mathbf{W}_s(k)} - \frac{\partial u_i}{\partial \mathbf{W}_s(k, t)} \right\|_F. \end{aligned}$$

In the following, we will bound each term. We first look at  $\left\| \frac{\partial u_i}{\partial \mathbf{W}_s(k)} - \frac{\partial u_i}{\partial \mathbf{W}_s(k, t)} \right\|_F$ . By using Lemma 8, we have  $\frac{\partial u_i}{\partial \mathbf{W}_s(k)} = \mathbf{X}_i^{(l)}(k)$ . Therefore, we can obtain

$$\begin{aligned} \left\| \frac{\partial u_i}{\partial \mathbf{W}_s(k)} - \frac{\partial u_i}{\partial \mathbf{W}_s(k, t)} \right\|_F &= \left\| \mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(k, t) \right\|_F = t \left\| \frac{\partial F(\Omega)}{\partial \mathbf{X}_i^{(l)}(k)} \right\|_F \\ &\leq t \frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell_i}{\partial \mathbf{X}_i^{(l)}(k)} \right\|_F \stackrel{\textcircled{1}}{\leq} \eta \left( 1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0} \right)^l c_y c_u, \end{aligned} \quad (10)$$

where  $\textcircled{1}$  holds since in Lemma 13, we have show

$$\max \left( \|\mathbf{W}^{(0)}(t) - \mathbf{W}^{(0)}(0)\|_F, \|\mathbf{W}_s^{(l)}(t) - \mathbf{W}_s^{(l)}(0)\|_F, \|\mathbf{W}_s(t) - \mathbf{W}_s(0)\|_F \right) \leq \sqrt{m\tilde{r}} \leq \sqrt{m}c_{w0}, \quad (11)$$

which allows us to use Lemma 12 which shows

$$\frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell_i}{\partial \mathbf{X}_i^{(l)}(k)} \right\|_F \leq \left( 1 + \alpha_2 + \alpha_3\mu\sqrt{k_c}(\tilde{r} + c_{w0}) \right)^l c_y c_u \leq \left( 1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0} \right)^l c_y c_u, \quad (12)$$

where parameters  $\frac{1}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2 = c_y$  and  $\|\mathbf{W}_h(t)\|_F \leq c_u$ ,  $\alpha_2 = \max_{s,t} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,t} \alpha_{s,3}^{(l)}$ . Moreover, from Lemma 13, we have  $\|\mathbf{W}_h(t)\|_F \leq \|\mathbf{W}_h(t) - \mathbf{W}_h(0)\|_F + \|\mathbf{W}_h(0)\|_F \leq 2\sqrt{m}c_{w0}$ . In this way, we have

$$\begin{aligned} \sum_{s=1}^h \left\| \frac{\partial u_i}{\partial \mathbf{W}_s(k)} - \frac{\partial u_i}{\partial \mathbf{W}_s(k,t)} \right\|_F &\leq \eta h \left(1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0}\right)^l \sqrt{m} c_{w0} \frac{1}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2 \\ &\leq \eta h \left(1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0}\right)^l \sqrt{m} c_{w0} \frac{1}{\sqrt{n}} \left(1 - \frac{\eta \lambda}{2}\right)^{t/2} \|\mathbf{u}(0) - \mathbf{y}\|_2 = \eta c_1, \end{aligned}$$

where  $c_1 = h \left(1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0}\right)^l \sqrt{m} c_{w0} \frac{1}{\sqrt{n}} \left(1 - \frac{\eta \lambda}{2}\right)^{t/2} \|\mathbf{u}(0) - \mathbf{y}\|_F$  is a constant.

Then we consider  $\left\| \frac{\partial u_i}{\partial \mathbf{W}_s^{(l)}(k)} - \frac{\partial u_i}{\partial \mathbf{W}_s^{(l)}(k,t)} \right\|_F$  as follows:

$$\begin{aligned} \left\| \frac{\partial u_i}{\partial \mathbf{W}_s^{(l)}(k)} - \frac{\partial u_i}{\partial \mathbf{W}_s^{(l)}(k,t)} \right\|_F &= \alpha_{s,3}^{(l)} \tau \left\| \left[ \Phi(\mathbf{X}_i^{(s)}(k)) \left( \sigma' \left( \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}_i^{(s)}(k)) \right) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k)} \right)^\top \right. \right. \\ &\quad \left. \left. - \Phi(\mathbf{X}_i^{(s)}(k,t)) \left( \sigma' \left( \mathbf{W}_s^{(l)}(k,t) \Phi(\mathbf{X}_i^{(s)}(k,t)) \right) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k,t)} \right)^\top \right] \right\|_F \\ &\stackrel{\textcircled{1}}{\leq} \alpha_{s,3}^{(l)} \tau \frac{a_1 a_2 (b_1 + b_2)}{\max(a_1, a_2)}, \end{aligned}$$

where  $\textcircled{1}$  uses Lemma 2. For parameters  $a_1, a_2, b_1, b_2$  satisfies

$$\begin{aligned} a_1 &= \max \left( \left\| \Phi(\mathbf{X}_i^{(s)}(k)) \right\|_2, \left\| \Phi(\mathbf{X}_i^{(s)}(k,t)) \right\|_2 \right) \leq \sqrt{k_c} \max \left( \left\| \mathbf{X}_i^{(s)}(k) \right\|_2, \left\| \mathbf{X}_i^{(s)}(k,t) \right\|_2 \right), \\ a_2 &= \max \left( \left\| \sigma' \left( \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}_i^{(s)}(k)) \right) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k)} \right\|_2, \left\| \sigma' \left( \mathbf{W}_s^{(l)}(k,t) \Phi(\mathbf{X}_i^{(s)}(k,t)) \right) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k,t)} \right\|_2 \right), \\ b_1 &= \left\| \Phi(\mathbf{X}_i^{(s)}(k)) - \Phi(\mathbf{X}_i^{(s)}(k,t)) \right\|_2 \leq \sqrt{k_c} \left\| \mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(k,t) \right\|_2, \\ b_2 &= \left\| \sigma' \left( \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}_i^{(s)}(k)) \right) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k)} - \sigma' \left( \mathbf{W}_s^{(l)}(k,t) \Phi(\mathbf{X}_i^{(s)}(k,t)) \right) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k,t)} \right\|_2. \end{aligned}$$

In Lemma 10, we show that when Eqn. (10) holds which is proven in Lemma 13, then  $\|\mathbf{X}_i^{(l)}(0)\|_F \leq c_{x0}$ . Under Eqn. (10), Lemma 11 shows

$$\|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0)\|_F \leq \left(1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0}\right)^l \mu \sqrt{k_c} \tilde{r} \stackrel{\textcircled{1}}{\leq} c_{x0}, \quad (13)$$

where  $\textcircled{1}$  holds since in Lemma 13, we set  $m = \mathcal{O} \left( \frac{\rho k_c^2 c_{w0}^2 \|\mathbf{y} - \mathbf{u}(0)\|_2^2}{\lambda^2 n} \left(1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0}\right)^{2h} \right)$  such that

$$\begin{aligned} \tilde{r} &= \frac{8c_{x0} \|\mathbf{y} - \mathbf{u}(0)\|_2}{\lambda \sqrt{mn}} \max \left( 1, 2 \left(1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0}\right)^l \alpha_{s,3}^{(l)} \mu \sqrt{k_c} c_{w0} \right) \\ &\leq \frac{c_{x0}}{\left(1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0}\right)^l \mu \sqrt{k_c}}. \end{aligned}$$

By using Lemma 11 and Lemma 10, we have

$$\|\mathbf{X}^{(s)}(t)\| \leq \|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0)\|_F + \|\mathbf{X}_i^{(l)}(0)\|_F \leq 2c_{x0}. \quad (14)$$

Then by using Eqn. (12) we upper bound  $\left\| \mathbf{X}_i^{(s)}(k,t) \right\|_2$  as follows:

$$\begin{aligned} \left\| \mathbf{X}_i^{(s)}(k,t) \right\|_2 &\leq \left\| \mathbf{X}_i^{(s)}(k) - t \frac{\partial F(\boldsymbol{\Omega})}{\partial \mathbf{X}_i^{(s)}(k)} \right\|_2 \leq \left\| \mathbf{X}_i^{(s)}(k) \right\|_2 + t \frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell_i}{\partial \mathbf{X}_i^{(s)}(k)} \right\|_2 \\ &\leq 2c_{x0} + \eta \left(1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0}\right)^l \sqrt{m} c_{w0} \frac{1}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_F \leq c_2, \end{aligned}$$

where  $c_2 = 2c_{x0} + \eta(1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l \sqrt{m}c_{w0} \frac{1}{\sqrt{n}}(1 - \frac{\eta\lambda}{2})^{t/2} \|\mathbf{u}(0) - \mathbf{y}\|_F$  is a constant. In this way, we can upper bound

$$a_1 \leq \sqrt{k_c} \max(2c_{w0}, c_2), \quad b_1 \stackrel{\textcircled{1}}{\leq} \frac{\sqrt{k_c}c_1\eta}{h},$$

where  $\textcircled{1}$  uses the results in Eqn. (10). Now we try to bound  $a_2$  and  $b_2$  as follows:

$$\begin{aligned} a_2 &= \max \left( \left\| \sigma'(\mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}_i^{(s)}(k))) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k)} \right\|_2, \left\| \sigma'(\mathbf{W}_s^{(l)}(k,t)\Phi(\mathbf{X}_i^{(s)}(k,t))) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k,t)} \right\|_2 \right) \\ &\leq \mu \max \left( \left\| \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k)} \right\|_2, \left\| \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k,t)} \right\|_2 \right) \stackrel{\textcircled{1}}{\leq} \mu(1+L)c_1^2\eta^2, \end{aligned}$$

where  $\textcircled{1}$  uses  $\left\| \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k,t)} \right\|_2 \leq \left\| \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k,t)} \right\|_F \leq \left\| \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k)} \right\|_F + L\|\mathbf{X}_i^{(l)}(k,t) - \mathbf{X}_i^{(l)}(k)\|_F^2 \stackrel{\textcircled{2}}{\leq} (1+L)c_1^2\eta^2$  where  $L$  is the Lipschitz constant of  $\frac{\partial u_i}{\partial \mathbf{X}^{(l)}}$ . In  $\textcircled{2}$  we use the results in Eqn. (14). Since  $\sigma$  is  $\rho$ -smooth and  $u$  is  $h$ -layered, by computing, we know  $L$  is at the order of  $\mathcal{O}(\beta^h)$  and is a constant. For  $b_2$  we can bound it as follows:

$$b_2 \leq \mu \left\| \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k)} - \frac{\partial u_i}{\partial \mathbf{X}_i^{(l)}(k,t)} \right\|_2 \leq 2\mu(1+L)c_1^2\eta^2.$$

Therefore, we can bound

$$\sum_{l=1}^h \sum_{s=0}^{l-1} \left\| \frac{\partial u_i}{\partial \mathbf{W}_s^{(l)}(k)} - \frac{\partial u_i}{\partial \mathbf{W}_s^{(l)}(k,t)} \right\|_F \leq \tau \frac{a_1 a_2 (b_1 + b_2)}{\max(a_1, a_2)} \sum_{l=1}^h \sum_{s=0}^{l-1} \alpha_{s,3}^{(l)} = c_3\eta,$$

where  $\alpha_3 = \max \alpha_{s,3}^{(l)}$  and  $c_3 = \frac{\tau\sqrt{k_c} \max(2c_{w0}, c_2)\mu(1+L)c_1^2\eta^2}{\max(\sqrt{k_c} \max(2c_{w0}, c_2), \mu(1+L)c_1^2\eta^2, \frac{\sqrt{k_c}c_1}{h} + 2\mu(1+L)c_1^2\eta^2)}$  is a constant. By using the same method, we can bound

$$\begin{aligned} &\left\| \frac{\partial u_i}{\partial \mathbf{W}^{(0)}(k)} - \frac{\partial u_i}{\partial \mathbf{W}^{(0)}(k,t)} \right\|_F \\ &= \tau \left\| \Phi(\mathbf{X}_i) \left( \sigma'(\mathbf{W}^{(0)}(k)\Phi(\mathbf{X}_i)) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(0)}(k)} \right)^\top - \Phi(\mathbf{X}_i) \left( \sigma'(\mathbf{W}^{(0)}(k,t)\Phi(\mathbf{X}_i)) \odot \frac{\partial u_i}{\partial \mathbf{X}_i^{(0)}(k,t)} \right)^\top \right\|_F \\ &\stackrel{\textcircled{1}}{\leq} \tau\sqrt{k_c} \left\| \frac{\partial u_i}{\partial \mathbf{X}_i^{(0)}(k)} - \frac{\partial u_i}{\partial \mathbf{X}_i^{(0)}(k,t)} \right\|_F \leq 2\mu(1+L)c_1^2\eta^2 = c_4\eta, \end{aligned}$$

where  $\textcircled{1}$  uses  $\|\Phi(\mathbf{X}_i)\|_F \leq \sqrt{k_c}\|\mathbf{X}_i\|_F \leq \sqrt{k_c}$  and  $\sigma$  is  $\mu$ -Lipschitz, and  $c_4 = 2\mu(1+L)c_1^2\eta$ . By combing the above results, we can further conclude

$$\|\nabla_{\Omega} u_i(\Omega(k)) - \nabla_{\Omega} u_i(\Omega(k,t))\|_F \leq (c_1 + c_3 + c_4)\eta = c_5\eta,$$

which further gives

$$\begin{aligned} \Delta_2^i(k) &\leq \eta \max_{t \in [0, \eta]} \|\nabla_{\Omega} F(\Omega(k))\|_F \|\nabla_{\Omega} u_i(\Omega(k)) - \nabla_{\Omega} u_i(\Omega(k,t)) - t\nabla_{\Omega} F(\Omega(k))\|_F \\ &\leq \eta^2 c_5 \left( h + 2c\mu c_{w0} \sqrt{k_c} \left( 1 + \sum_{l=1}^h \sum_{s=0}^{l-1} \alpha_{s,3}^{(l)} \right) \right) \frac{2c_{x0}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_F = \hat{c}\eta^2 \|\mathbf{u}(t) - \mathbf{y}\|_F, \end{aligned}$$

where  $\hat{c} = c_5 \left( h + 2c\mu c_{w0} \sqrt{k_c} \left( 1 + \sum_{l=1}^h \sum_{s=0}^{l-1} \alpha_{s,3}^{(l)} \right) \right) \frac{2c_{x0}}{\sqrt{n}}$ . Therefore we have

**Step 3.3 Upper bound  $H_2 = 2\|\mathbf{y} - \mathbf{u}(k)\|_2 \|\Delta_2(k)\|_2$ .** By combining the above results, we can bound

$$H_2 = 2\|\mathbf{y} - \mathbf{u}(k)\|_2 \|\Delta_2(k)\|_2 \leq \hat{c}\eta^2 \|\mathbf{u}(t) - \mathbf{y}\|_2^2,$$

where  $\hat{c} = \mathcal{O} \left( \frac{\mu c_{x0} c_{w0}^2 \sqrt{k_c} m h^3 (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0})^h}{n} \right)$ .

**Step 4. Upper bound**  $H_3 = \|\mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2$ .

$$\begin{aligned} \|\mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2 &= \sum_{i=1}^n \left( \sum_{s=0}^{h-1} \left( \langle \mathbf{W}_s(k), \mathbf{X}_i^{(l)}(k) \rangle - \langle \mathbf{W}_s(k+1), \mathbf{X}_i^{(l)}(k+1) \rangle \right) \right)^2 \\ &\leq \sqrt{h} \sum_{i=1}^n \sum_{s=0}^{h-1} \left( \langle \mathbf{W}_s(k), \mathbf{X}_i^{(l)}(k) \rangle - \langle \mathbf{W}_s(k+1), \mathbf{X}_i^{(l)}(k+1) \rangle \right)^2. \end{aligned}$$

Now we consider each term:

$$\begin{aligned} &\left( \langle \mathbf{W}_s(k), \mathbf{X}_i^{(l)}(k) \rangle - \langle \mathbf{W}_s(k+1), \mathbf{X}_i^{(l)}(k+1) \rangle \right)^2 \\ &= \left( \langle \mathbf{W}_s(k) - \mathbf{W}_s(k+1), \mathbf{X}_i^{(l)}(k+1) \rangle + \langle \mathbf{W}_s(k), \mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(k+1) \rangle \right)^2 \\ &\leq 2\|\mathbf{W}_s(k) - \mathbf{W}_s(k+1)\|_F^2 \|\mathbf{X}_i^{(l)}(k+1)\|_F^2 + 2\|\mathbf{W}_s(k)\|_F^2 \|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(k+1)\|_F^2 \\ &\stackrel{\textcircled{1}}{\leq} 8c_{x0}^2 \|\mathbf{W}_s(k) - \mathbf{W}_s(k+1)\|_F^2 + 8mc_{w0}^2 \|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(k+1)\|_F^2 \\ &\stackrel{\textcircled{2}}{\leq} \frac{32\eta^2 c_{x0}^2}{n} \left[ c_{x0}^2 + 4c^2 \mu^4 c_{w0}^4 k_c^2 \left( 1 + \alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu \right)^{2l} \left( 1 + \frac{2(\alpha_3)^2 c_{x0}}{(\alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu) \sqrt{n}} \right)^2 \right] \\ &\quad \cdot \|\mathbf{u}(k) - \mathbf{y}\|_2^2, \end{aligned}$$

where  $\textcircled{1}$  uses  $\|\mathbf{X}_i^{(l)}(k+1)\|_F^2 \leq 4c_{x0}^2$  in Eqn. (14), and the results in Eqn. (11) that  $\|\mathbf{W}_s(k)\|_F \leq \|\mathbf{W}_s(k) - \mathbf{W}_s(0)\|_F + \|\mathbf{W}_s(0)\|_F \leq 2\sqrt{m}c_{w0}$ ;  $\textcircled{2}$  holds since (1) in Lemma 13 we have  $\|\mathbf{W}_s(t+1) - \mathbf{W}_s(t)\|_F = \eta \left\| \frac{\partial F(\Omega)}{\partial \mathbf{W}_s(t)} \right\|_F \leq \frac{2\eta c_{x0}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2$  where  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ , and (2) in Lemma 14 we have

$$\begin{aligned} &\left\| \mathbf{X}^{(l)}(k+1) - \mathbf{X}^{(l)}(k) \right\|_F \\ &\leq \left( 1 + \alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu \right)^l \left( 1 + \frac{2(\alpha_3)^2 c_{x0}}{(\alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu) \sqrt{n}} \right) \frac{4c\tau\eta\mu^2 c_{x0} c_{w0} k_c}{\sqrt{n}} \|\mathbf{u}(k) - \mathbf{y}\|_2. \end{aligned}$$

In this way, we can conclude

$$\|\mathbf{u}(k) - \mathbf{u}(k+1)\|_2^2 \leq \eta^2 \tilde{c} \|\mathbf{u}(k) - \mathbf{y}\|_2^2,$$

where  $\tilde{c} = 32c_{x0}^2 h^{1.5} \left[ c_{x0}^2 + 4c^2 \mu^4 c_{w0}^4 k_c^2 \left( 1 + \alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu \right)^{2l} \left( 1 + \frac{2(\alpha_3)^2 c_{x0}}{(\alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu) \sqrt{n}} \right)^2 \right] = \mathcal{O} \left( \mu^4 c_{w0}^4 c_{x0}^2 h^{1.5} k_c^2 \left( 1 + \alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu \right)^{4l} \right)$ .

**Step 5. Upper bound**  $\|\mathbf{y} - \mathbf{u}(k+1)\|_2^2$ .

In this way, by using Eqn. (9) we can finally obtain

$$\begin{aligned} \|\mathbf{y} - \mathbf{u}(k+1)\|_2^2 &\leq \|\mathbf{y} - \mathbf{u}(k)\|_2^2 + H_1 + H_2 + H_3 \\ &\stackrel{\textcircled{1}}{\leq} \|\mathbf{y} - \mathbf{u}(k)\|_2^2 - 2\eta\lambda \|\mathbf{u}(k) - \mathbf{y}\|_2^2 + 2\hat{c}\eta^2 \|\mathbf{u}(t) - \mathbf{y}\|_2^2 + \eta^2 \tilde{c} \|\mathbf{u}(k) - \mathbf{y}\|_2^2 \\ &= (1 - \eta\lambda + (2\hat{c} + \tilde{c})\eta^2) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \\ &\stackrel{\textcircled{2}}{\leq} \left( 1 - \frac{\eta\lambda}{2} \right) \|\mathbf{y} - \mathbf{u}(k)\|_2^2 \end{aligned}$$

where  $\textcircled{1}$  holds by using  $H_1 \leq -2\eta\lambda \|\mathbf{u}(k) - \mathbf{y}\|_2^2$ ,  $H_2 \leq 2\hat{c}\eta^2 \|\mathbf{u}(t) - \mathbf{y}\|_2^2$  and  $H_3 \leq \eta^2 \tilde{c} \|\mathbf{u}(k) - \mathbf{y}\|_2^2$ ;

$\textcircled{2}$  holds by setting  $\eta \leq \frac{\lambda}{2(2\hat{c} + \tilde{c})} = \mathcal{O} \left( \frac{\lambda}{\sqrt{m}\mu^4 c_{w0}^4 c_{x0}^2 h^3 k_c^2 (1 + \alpha_2 + 2\sqrt{k_c c_{w0}} \alpha_3 \mu)^{4l}} \right)$ . The proof is completed.  $\square$

### C.3.2 Proof of Lemma 20

*Proof.* According to the definitions in Sec. B, we can write

$$\|\mathbf{G}(k) - \mathbf{G}(0)\|_2 \leq \|\bar{\mathbf{G}}^0(k) - \bar{\mathbf{G}}^0(0)\|_2 + \sum_{l=0}^{h-1} \sum_{s=0}^{l-1} \left\| \mathbf{G}^{ls}(k) - \mathbf{G}^{ls}(0) \right\|_2 + \sum_{s=0}^{h-1} \|\mathbf{G}^s(k) - \mathbf{G}^s(0)\|_2.$$

In this way, we only need to upper bound  $\|\bar{\mathbf{G}}^0(k) - \bar{\mathbf{G}}^0(0)\|_2$ ,  $\|\mathbf{G}^{ls}(k) - \mathbf{G}^{ls}(0)\|_2$  and  $\|\mathbf{G}^s(k) - \mathbf{G}^s(0)\|_2$ .

**Step 1. Bound of  $\|\mathbf{G}^s(k) - \mathbf{G}^s(0)\|_2$  ( $s = 0, \dots, h-1$ ).**

For analysis, we first recall existing results. Lemma 13 shows

$$\max\left(\|\mathbf{W}^{(0)}(t) - \mathbf{W}^{(0)}(0)\|_F, \|\mathbf{W}_s^{(l)}(t) - \mathbf{W}_s^{(l)}(0)\|_F, \|\mathbf{W}_s(t) - \mathbf{W}_s(0)\|_F\right) \leq \sqrt{m}\tilde{r} \leq \sqrt{m}c_{w0}, \quad (15)$$

where  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . Based on this result, Lemma 15 shows

$$\left\|\mathbf{W}^{(0)}(k)\right\|_F \leq 2\sqrt{m}c_{w0}, \quad \left\|\mathbf{W}_s^{(l)}(k)\right\|_F \leq 2\sqrt{m}c_{w0}, \quad \|\mathbf{W}_s(k)\|_F \leq 2\sqrt{m}c_{w0}, \quad \left\|\mathbf{X}_i^{(l)}(k)\right\|_F \leq 2c_{x0}. \quad (16)$$

Moreover, Lemma 16 shows

$$\|\mathbf{X}_i^{(0)}(k) - \mathbf{X}_i^{(0)}(0)\|_F \leq \mu\sqrt{k_c}\tilde{r}, \quad \|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0)\|_F \leq c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c}\tilde{r}.$$

To bound  $H_s$ , we only need to bound each entry in  $(\mathbf{G}^s(k) - \mathbf{G}^s(0))$ :

$$\begin{aligned} |\mathbf{G}^s(k) - \mathbf{G}^s(0)| &= \left| \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}_s(k)}, \frac{\partial \ell_j}{\partial \mathbf{W}_s(k)} \right\rangle - \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}_s(0)}, \frac{\partial \ell_j}{\partial \mathbf{W}_s(0)} \right\rangle \right| \\ &= \left| \left\langle \mathbf{X}_i^{(s)}(k), \mathbf{X}_j^{(s)}(k) \right\rangle - \left\langle \mathbf{X}_i^{(s)}(0), \mathbf{X}_j^{(s)}(0) \right\rangle \right| \\ &\leq \left| \left\langle \mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0), \mathbf{X}_j^{(s)}(k) \right\rangle \right| + \left| \left\langle \mathbf{X}_i^{(s)}(0), \mathbf{X}_j^{(s)}(k) - \mathbf{X}_j^{(s)}(0) \right\rangle \right| \\ &\leq \left\| \mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0) \right\|_F \left\| \mathbf{X}_j^{(s)}(k) \right\|_F + \left\| \mathbf{X}_i^{(s)}(0) \right\|_F \left\| \mathbf{X}_j^{(s)}(k) - \mathbf{X}_j^{(s)}(0) \right\|_F \\ &\stackrel{\textcircled{1}}{\leq} 4c_{x0}c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c}\tilde{r}, \end{aligned}$$

So we can further bound

$$\|\mathbf{G}^s(k) - \mathbf{G}^s(0)\|_2 \leq \sqrt{n} \|\mathbf{G}^s(k) - \mathbf{G}^s(0)\|_\infty \leq 4c_{x0}c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c}\tilde{r}, \quad (1 \leq s \leq h).$$

**Step 2. Bound of  $\|\mathbf{G}^{ls}(k) - \mathbf{G}^{ls}(0)\|_2$  ( $0 \leq l \leq h-1, 0 \leq s \leq l-1$ ).**

**We first consider  $l = h-1$ , namely bound of  $\|\mathbf{G}^{hs}(k) - \mathbf{G}^{hs}(0)\|_2$  ( $0 \leq s \leq h-2$ ).** For notation simplicity, we use  $h$  to denote  $h-1$ . In this way, according to Lemma 8, we have

$$\frac{\partial u}{\partial \mathbf{W}_s^{(h)}} = \alpha_{s,3}^{(h)} \tau \Phi(\mathbf{X}^{(s)}) \left( \sigma' \left( \mathbf{W}_s^{(h)} \Phi(\mathbf{X}^{(s)}) \right) \odot \mathbf{W}_h \right)^\top \quad (1 \leq s \leq h-1).$$

Let  $\mathbf{H}_i = \Phi(\mathbf{X}_i^{(s)})$ ,  $\mathbf{H}_{i,:t} = [\mathbf{H}_i]_{:,t}$ ,  $\mathbf{H}_{i,tr} = [\mathbf{H}_i]_{t,r}$ , and  $\mathbf{Z}_{i,tr} = (\mathbf{W}_{s,:r}^{(h)})^\top \mathbf{H}_{i,:t}$ . In this way, for  $1 \leq s \leq h-1$  we can write  $\mathbf{G}_{ij}^{hs}$  as

$$\begin{aligned} \mathbf{G}_{ij}^{hs} &= (\alpha_{s,3}^{(h)} \tau)^2 \sum_{r=1}^m \left[ \sum_{t=1}^p \mathbf{W}_{h,tr} \mathbf{H}_{i,:t} \left( \sigma' \left( (\mathbf{W}_{s,:r}^{(h)})^\top \mathbf{H}_{i,:t} \right) \right) \right]^\top \left[ \sum_{q=1}^p \mathbf{W}_{h,qr} \mathbf{H}_{j,:q} \left( \sigma' \left( (\mathbf{W}_{s,:r}^{(h)})^\top \mathbf{H}_{j,:q} \right) \right) \right] \\ &= (\alpha_{s,3}^{(h)} \tau)^2 \sum_{t=1}^p \sum_{q=1}^p \mathbf{H}_{i,:t}^\top \mathbf{H}_{j,:q} \sum_{r=1}^m \mathbf{W}_{h,tr} \mathbf{W}_{h,qr} \sigma'(\mathbf{Z}_{i,tr}) \sigma'(\mathbf{Z}_{j,qr}). \end{aligned}$$

Then we can obtain

$$\begin{aligned} &|\mathbf{G}_{ij}^{hs}(k) - \mathbf{G}_{ij}^{hs}(0)| \\ &= (\alpha_{s,3}^{(h)} \tau)^2 \left| \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(k))^\top \mathbf{H}_{j,:q}(k) \sum_{r=1}^m \mathbf{W}_{h,tr}(k) \mathbf{W}_{h,qr}(k) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) \right. \\ &\quad \left. - \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m \mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0) \sigma'(\mathbf{Z}_{i,tr}(0)) \sigma'(\mathbf{Z}_{j,qr}(0)) \right|. \end{aligned}$$

For brevity, we define  $\mathbf{A}_1$ ,  $\mathbf{A}_2$  and  $\mathbf{A}_3$  as follows:

$$\begin{aligned}\mathbf{A}_1 &= \left| \sum_{t=1}^p \sum_{q=1}^p \left( (\mathbf{H}_{i,:t}(k))^\top \mathbf{H}_{j,:q}(k) - (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \right) \sum_{r=1}^m \mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) \right|, \\ \mathbf{A}_2 &= \left| \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m \mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0) (\sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) - \sigma'(\mathbf{Z}_{i,tr}(0)) \sigma'(\mathbf{Z}_{j,qr}(0))) \right|, \\ \mathbf{A}_3 &= \left| \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m (\mathbf{W}_{h,tr}(k) \mathbf{W}_{h,qr}(k) - \mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0)) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) \right|.\end{aligned}$$

Then we have

$$|\mathbf{G}_{ij}^{hs}(k) - \mathbf{G}_{ij}^{hs}(0)| = (\alpha_{s,3}^{(h)} \tau)^2 (\mathbf{A}_1 + \mathbf{A}_2 + \mathbf{A}_3).$$

The remaining work is to upper bound  $\mathbf{A}_1$ ,  $\mathbf{A}_2$  and  $\mathbf{A}_3$ . We first look at  $\mathbf{A}_1$ :

$$\begin{aligned}\mathbf{A}_1 &= \left| \sum_{t=1}^p \sum_{q=1}^p \left( (\mathbf{H}_{i,:t}(k))^\top \mathbf{H}_{j,:q}(k) - (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \right) \sum_{r=1}^m \mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) \right| \\ &\leq m\mu^2 c_{u0}^2 \left| \sum_{t=1}^p \sum_{q=1}^p \left( (\mathbf{H}_{i,:t}(k))^\top \mathbf{H}_{j,:q}(k) - (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \right) \right| \\ &\stackrel{\textcircled{1}}{\leq} m\mu^2 c_{u0}^2 \sum_{t=1}^p \sum_{q=1}^p \left[ \left| (\mathbf{H}_{i,:t}(k) - \mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(k) \right| + \left| (\mathbf{H}_{i,:t}(0))^\top (\mathbf{H}_{j,:q}(k) - \mathbf{H}_{j,:q}(0)) \right| \right] \\ &\leq m\mu^2 c_{u0}^2 \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{H}_{i,:t}(k) - \mathbf{H}_{i,:t}(0)\|_2^2} \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{H}_{j,:q}(k)\|_2^2} \\ &\quad + m\mu^2 c_{u0}^2 \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{H}_{j,:q}(k) - \mathbf{H}_{j,:q}(0)\|_2^2} \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{H}_{i,:t}(0)\|_2^2} \\ &\leq mp\mu^2 c_{u0}^2 (\|\mathbf{H}_i(k) - \mathbf{H}_i(0)\|_F \|\mathbf{H}_j(k)\|_F + \|\mathbf{H}_j(k) - \mathbf{H}_j(0)\|_F \|\mathbf{H}_i(k)\|_F) \\ &\leq mp\mu^2 c_{u0}^2 (\|\mathbf{H}_i(k) - \mathbf{H}_i(0)\|_F \|\mathbf{H}_j(k)\|_F + \|\mathbf{H}_j(k) - \mathbf{H}_j(0)\|_F \|\mathbf{H}_i(k)\|_F)\end{aligned}$$

where  $\textcircled{1}$  holds since the activation function  $\sigma(\cdot)$  is  $\mu$ -Lipschitz and  $\rho$ -smooth and the assumption  $\|\mathbf{W}_s\|_\infty \leq c_{u0}$ . To bound  $\|\mathbf{H}_i(k) - \mathbf{H}_i(0)\|_F \|\mathbf{H}_j(k)\|_F$ , we first recall our existing results. Lemma 16 that

$$\|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0)\|_F \leq c(1 + 2\alpha_3 c_{x0}) \mu \sqrt{k_c} \tilde{r},$$

where  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . Here  $\tilde{r}$  is given in Lemma 13. Based on this result, Lemma 15 shows that (16) holds. So we have

$$\begin{aligned}\|\mathbf{H}_i(k) - \mathbf{H}_i(0)\|_F &\leq \|\Phi(\mathbf{X}_i^{(s)}(k)) - \Phi(\mathbf{X}_i^{(s)}(0))\|_F \leq \sqrt{k_c} \|\mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0)\|_F \\ &\leq c(1 + 2\alpha_3 c_{x0}) \mu k_c \tilde{r}, \\ \|\mathbf{H}_j(k)\|_F &= \|\Phi(\mathbf{X}_j^{(s)}(k))\|_F \leq \sqrt{k_c} \|\mathbf{X}_j^{(s)}(k)\|_F \leq 2\sqrt{k_c} c_{w0},\end{aligned}\tag{17}$$

which indicates

$$(\|\mathbf{H}_i(k) - \mathbf{H}_i(0)\|_F \|\mathbf{H}_j(k)\|_F + \|\mathbf{H}_j(k) - \mathbf{H}_j(0)\|_F \|\mathbf{H}_i(k)\|_F) \leq 4cc_{w0}(1 + 2\alpha_3 c_{x0}) \mu k_c^{1.5} \tilde{r}.$$

Therefore, we can upper bound

$$\mathbf{A}_1 \leq 4cmp\mu^3 k_c^{1.5} c_{u0}^2 c_{w0} (1 + 2\alpha_3 c_{x0}) \tilde{r}.$$

Then we consider to bound  $\mathbf{A}_2$ . To begin with, we have

$$\begin{aligned}& \left| \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) - \sigma'(\mathbf{Z}_{i,tr}(0)) \sigma'(\mathbf{Z}_{j,qr}(0)) \right| \\ &\leq \left| (\sigma'(\mathbf{Z}_{i,tr}(k)) - \sigma'(\mathbf{Z}_{i,tr}(0))) \sigma'(\mathbf{Z}_{j,qr}(k)) \right| + \left| \sigma'(\mathbf{Z}_{i,tr}(0)) (\sigma'(\mathbf{Z}_{j,qr}(k)) - \sigma'(\mathbf{Z}_{j,qr}(0))) \right| \\ &\stackrel{\textcircled{1}}{\leq} \mu \left| \sigma'(\mathbf{Z}_{i,tr}(k)) - \sigma'(\mathbf{Z}_{i,tr}(0)) \right| + \mu \left| \sigma'(\mathbf{Z}_{j,qr}(k)) - \sigma'(\mathbf{Z}_{j,qr}(0)) \right| \\ &\stackrel{\textcircled{2}}{\leq} \mu\rho \left| \mathbf{Z}_{i,tr}(k) - \mathbf{Z}_{i,tr}(0) \right| + \mu\rho \left| \mathbf{Z}_{j,qr}(k) - \mathbf{Z}_{j,qr}(0) \right|,\end{aligned}$$

where ① holds since the activation function  $\sigma(\cdot)$  is  $\mu$ -Lipschitz; ② holds since the activation function  $\sigma(\cdot)$  is  $\rho$ -smooth. Therefore, we can upper bound

$$\begin{aligned}
\mathbf{A}_2 &\leq \sum_{t=1}^p \sum_{q=1}^p \left| \mathbf{H}_{i,:t}(0)^\top \mathbf{H}_{j,:q}(0) \right| \sum_{r=1}^m |\mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0)| \\
&\quad \cdot \left| (\sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) - \sigma'(\mathbf{Z}_{i,tr}(0)) \sigma'(\mathbf{Z}_{j,qr}(0))) \right| \\
&\leq \mu \rho \sum_{t=1}^p \sum_{q=1}^p \left| \mathbf{H}_{i,:t}(0)^\top \mathbf{H}_{j,:q}(0) \right| \sum_{r=1}^m |\mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0)| \left[ |\mathbf{Z}_{i,tr}(k) - \mathbf{Z}_{i,tr}(0)| + |\mathbf{Z}_{j,qr}(k) - \mathbf{Z}_{j,qr}(0)| \right] \\
&\leq \mu \rho \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{H}_{i,:t}(0)\|_2^2 \|\mathbf{H}_{j,:q}(0)\|_2^2} \left[ \sqrt{\sum_{t=1}^p \sum_{q=1}^p \left( \sum_{r=1}^m |\mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0)| |\mathbf{Z}_{i,tr}(k) - \mathbf{Z}_{i,tr}(0)| \right)^2} \right. \\
&\quad \left. + \sqrt{\sum_{t=1}^p \sum_{q=1}^p \left( \sum_{r=1}^m |\mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0)| |\mathbf{Z}_{j,qr}(k) - \mathbf{Z}_{j,qr}(0)| \right)^2} \right] \\
&\leq \mu \rho c_{u0} \sqrt{m} \|\mathbf{H}_i(0)\|_F \|\mathbf{H}_j(0)\|_F \cdot \\
&\quad \left[ \sqrt{\sum_{t=1}^p \sum_{q=1}^p \sum_{r=1}^m |\mathbf{Z}_{i,tr}(k) - \mathbf{Z}_{i,tr}(0)|^2} + \sqrt{\sum_{t=1}^p \sum_{q=1}^p \sum_{r=1}^m |\mathbf{Z}_{j,tr}(k) - \mathbf{Z}_{j,tr}(0)|^2} \right] \\
&\leq \mu \rho c_{u0} \sqrt{mp} \|\mathbf{H}_i(0)\|_F \|\mathbf{H}_j(0)\|_F \left[ \|\mathbf{Z}_i(k) - \mathbf{Z}_i(0)\|_F + \|\mathbf{Z}_j(k) - \mathbf{Z}_j(0)\|_F \right].
\end{aligned}$$

From Eqn. (17), we have  $\|\mathbf{H}_j(k)\|_F \leq 2\sqrt{k}c_{w0}$ . Lemma 13 shows that Eqn. (15) holds. Based on this result and the fact that  $\tilde{r} \leq c_{w0}$ , Lemma 11 shows

$$\left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k)) - \mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F \leq \frac{c}{\alpha_3} \sqrt{k} c_{w0} \tilde{r}.$$

Therefore we can bound

$$\mathbf{A}_2 \leq \frac{8cmk_c^{1.5} c_{w0}^2 \mu \rho c_{u0} \sqrt{p\tilde{r}}}{\alpha_3}.$$

Now we bound  $\mathbf{A}_3$  as follows:

$$\begin{aligned}
\mathbf{A}_3 &= \left| \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m (\mathbf{W}_{h,tr}(k) \mathbf{W}_{h,qr}(k) - \mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0)) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) \right| \\
&\leq \mu^2 \left| \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m (\mathbf{W}_{h,tr}(k) \mathbf{W}_{h,qr}(k) - \mathbf{W}_{h,tr}(0) \mathbf{W}_{h,qr}(0)) \right| \\
&\leq \mu^2 \sum_{t=1}^p \sum_{q=1}^p \left| (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \right| \sum_{r=1}^m (|\mathbf{W}_{h,tr}(k) - \mathbf{W}_{h,tr}(0)| \|\mathbf{W}_{h,qr}(k)\| + |\mathbf{W}_{h,tr}(0)| \|\mathbf{W}_{h,qr}(k) - \mathbf{W}_{h,qr}(0)\|) \\
&\leq \mu^2 \sum_{t=1}^p \sum_{q=1}^p \left| (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \right| (\|\mathbf{W}_{h,t}(k) - \mathbf{W}_{h,t}(0)\|_2 \|\mathbf{W}_{h,q}(k)\|_2 + \|\mathbf{W}_{h,t}(0)\|_2 \|\mathbf{W}_{h,qr}(k) - \mathbf{W}_{h,qr}(0)\|_2) \\
&\leq \mu^2 \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{H}_{i,:t}(0)\|_2^2 \|\mathbf{H}_{j,:q}(0)\|_2^2} \left[ \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{W}_{h,t}(k) - \mathbf{W}_{h,t}(0)\|_2 \|\mathbf{W}_{h,q}(k)\|_2} \right. \\
&\quad \left. + \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{W}_{h,t}(k) - \mathbf{W}_{h,t}(0)\|_2 \|\mathbf{W}_{h,q}(k)\|_2} \right] \\
&\leq \mu^2 \|\mathbf{H}_i(0)\|_F \|\mathbf{H}_j(0)\|_F \left[ \|\mathbf{W}_h(k) - \mathbf{W}_h(0)\|_F \|\mathbf{W}_h(k)\|_F + \|\mathbf{W}_h(k) - \mathbf{W}_h(0)\|_F \|\mathbf{W}_h(k)\|_F \right] \\
&\stackrel{\textcircled{1}}{\leq} 8k_c \mu^2 c_{w0}^3 m \tilde{r},
\end{aligned}$$

where ① holds by using Eqn.s (15), (16), (17).

By combining the above results, we have that for  $s = 0, \dots, h-1$

$$\begin{aligned}
\|\mathbf{G}^{hs}(k) - \mathbf{G}^{hs}(0)\|_2 &\leq \sqrt{n} \|\mathbf{G}_{ij}^{hs}(k) - \mathbf{G}_{ij}^{hs}(0)\|_\infty \\
&\leq 4(\alpha_{s,3}^{(h)})^2 k_c \mu c_{w0} n^{0.5} \tilde{r} \left( cp \mu^2 k_c^{0.5} c_{u0}^2 (1 + 2\alpha_3 c_{x0}) + \frac{2ck_c^{0.5} c_{w0} \rho c_{u0} \sqrt{p}}{\alpha_3} + 2\mu c_{w0}^2 \right).
\end{aligned}$$



Then we consider  $1 \leq l < h$ , namely bound of  $H_{ls}$  ( $0 \leq s \leq h-1$ ). For brevity, let  $\mathbf{B}_i(k) = \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(k)}$ . Here we use the same strategy as above. Let

$$\begin{aligned} \mathbf{A}_1 &= \sum_{t=1}^p \sum_{q=1}^p \left( (\mathbf{H}_{i,:t}(k))^\top \mathbf{H}_{j,:q}(k) - (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \right) \sum_{r=1}^m \mathbf{B}_{i,tr}(0) \mathbf{B}_{j,qr}(0) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)), \\ \mathbf{A}_2 &= \sum_{t=1}^p \sum_{q=1}^p \mathbf{H}_{i,:t}(0)^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m \mathbf{B}_{i,tr}(0) \mathbf{B}_{j,qr}(0) \left( \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) - \sigma'(\mathbf{Z}_{i,tr}(0)) \sigma'(\mathbf{Z}_{j,qr}(0)) \right), \\ \mathbf{A}_{3,ij} &= \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m \left( \mathbf{B}_{i,tr}(k) \mathbf{B}_{j,qr}(k) - \mathbf{B}_{i,tr}(0) \mathbf{B}_{j,qr}(0) \right) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)). \end{aligned}$$

By assuming  $\|\mathbf{B}_i(k)\|_\infty \leq c_{u0}$ , we can use the same method to bound  $\mathbf{A}_1$  and  $\mathbf{A}_2$  as follows:

$$|\mathbf{A}_1| \leq 4c_m p \mu^3 k_c^{1.5} c_{u0}^2 c_{w0} (1 + 2\alpha_3 c_{x0}) \tilde{r}, \quad |\mathbf{A}_2| \leq \frac{8c_m k_c^{1.5} c_{w0}^2 \mu \rho c_{u0} \sqrt{p\tilde{r}}}{\alpha_3}.$$

Then we need to carefully bound  $\mathbf{A}_3$ :

$$\begin{aligned} |\mathbf{A}_{3,ij}| &= \left| \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m \left( \mathbf{B}_{i,tr}(k) \mathbf{B}_{j,qr}(k) - \mathbf{B}_{i,tr}(0) \mathbf{B}_{j,qr}(0) \right) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) \right| \\ &\leq \mu^2 \left| \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m \left( \mathbf{B}_{i,tr}(k) \mathbf{B}_{j,qr}(k) - \mathbf{B}_{i,tr}(0) \mathbf{B}_{j,qr}(0) \right) \right| \\ &\leq \mu^2 \sum_{t=1}^p \sum_{q=1}^p \left| (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \left[ \sum_{r=1}^m \left( \|\mathbf{B}_{i,tr}(k) - \mathbf{B}_{i,tr}(0)\| \|\mathbf{B}_{j,qr}(k)\| + \|\mathbf{B}_{i,tr}(0)\| \|\mathbf{B}_{j,qr}(k) - \mathbf{B}_{j,qr}(0)\| \right) \right] \right| \\ &\leq \mu^2 \sum_{t=1}^p \sum_{q=1}^p \left| (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \right| \left( \|\mathbf{B}_{i,t}:(k) - \mathbf{B}_{i,t}:(0)\|_2 \|\mathbf{B}_{j,q}:(k)\|_2 + \|\mathbf{B}_{i,t}:(0)\|_2 \|\mathbf{B}_{j,q}:(k) - \mathbf{B}_{j,q}:(0)\|_2 \right) \\ &\leq \mu^2 \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{H}_{i,:t}(0)\|_2^2 \|\mathbf{H}_{j,:q}(0)\|_2^2} \left[ \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{B}_{i,t}:(k) - \mathbf{B}_{i,t}:(0)\|_2^2 \|\mathbf{B}_{j,q}:(k)\|_2^2} \right. \\ &\quad \left. + \sqrt{\sum_{t=1}^p \sum_{q=1}^p \|\mathbf{B}_{i,t}:(0)\|_2^2 \|\mathbf{B}_{j,q}:(k) - \mathbf{B}_{j,q}:(0)\|_2^2} \right] \\ &\leq \mu^2 \|\mathbf{H}_i(0)\|_F \|\mathbf{H}_j(0)\|_F \left[ \|\mathbf{B}_i(k) - \mathbf{B}_i(0)\|_F \|\mathbf{B}_j(k)\|_F + \|\mathbf{B}_j(k) - \mathbf{B}_j(0)\|_F \|\mathbf{B}_i(0)\|_F \right] \\ &\stackrel{\textcircled{1}}{\leq} 4\mu^2 c_{w0}^2 \left[ \|\mathbf{B}_i(k) - \mathbf{B}_i(0)\|_F \|\mathbf{B}_j(k)\|_F + \|\mathbf{B}_j(k) - \mathbf{B}_j(0)\|_F \|\mathbf{B}_i(0)\|_F \right], \end{aligned}$$

where  $\textcircled{1}$  holds by using Eqn.s (15), (16), (17). Then when for  $c_y = \frac{1}{\sqrt{n}} \|\mathbf{u}^t - \mathbf{y}\|_2$  and  $c_u = \|\mathbf{W}_t\|_F$ , Lemma 12 shows

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)} \right\|_F &\leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right)^l c_y c_u \\ &\stackrel{\textcircled{1}}{\leq} 2c \sqrt{m} c_{w0} \left( 1 - \frac{\eta \lambda}{2} \right)^{t/2} \|\mathbf{u}^0 - \mathbf{y}\|_2, \end{aligned}$$

where  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0})^l$ ,  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ .  $\textcircled{1}$  holds since  $c_u = \|\mathbf{W}_t\|_F \leq \|\mathbf{W}_t - \mathbf{W}_0\|_F + \|\mathbf{W}_0\|_F \leq \sqrt{m}(\tilde{r} + c_{w0}) \leq 2\sqrt{m}c_{w0}$  and  $\|\mathbf{u}^t - \mathbf{y}\|_2 \leq (1 - \frac{\eta \lambda}{2})^{t/2} \|\mathbf{u}^0 - \mathbf{y}\|_2$  in Theorem 19. Lemma 17 proves

$$\left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(0)} \right\|_F \leq c_1 c \alpha_3 c_{w0}^2 c_{x0} \rho k_c m \tilde{r},$$

where  $c_1$  is a constant. The remaining work is to bound

$$\|\mathbf{B}_i(k) - \mathbf{B}_i(0)\|_F \|\mathbf{B}_j(k)\|_F \leq c_1 c \alpha_3 c_{w0}^2 c_{x0} \rho k_c m \tilde{r} \|\mathbf{B}_j(k)\|_F.$$

In this way, we have

$$\begin{aligned}\|\mathbf{A}_3\|_1 &\leq \sum_{j=1}^n \sum_{i=1}^n \|\mathbf{A}_{3,ij}\| \leq 4\mu^2 c_{w0}^2 c_1 c \alpha_3 c_{w0}^2 c_{x0} \rho k_c m \tilde{r} \sum_{j=1}^n \sum_{i=1}^n (\|\mathbf{B}_j(k)\|_F + \|\mathbf{B}_i(k)\|_F) \\ &\leq 8c_1 n \mu^2 c^2 \alpha_3 c_{w0}^5 c_{x0} \rho k_c m^{1.5} \tilde{r} \left(1 - \frac{\eta\lambda}{2}\right)^{t/2} \|\mathbf{u}^0 - \mathbf{y}\|_2.\end{aligned}$$

Then combining all above results gives

$$\begin{aligned}\|\mathbf{G}^{hs}(k) - \mathbf{G}^{hs}(0)\|_2 &= (\alpha_{s,3}^{(h)} \tau)^2 \|\mathbf{A}_1 + \mathbf{A}_2 + \mathbf{A}_3\|_2 \leq (\alpha_{s,3}^{(h)} \tau)^2 (\|\mathbf{A}_1\|_2 + \|\mathbf{A}_2\|_2 + \|\mathbf{A}_3\|_2) \\ &\leq (\alpha_{s,3}^{(h)} \tau)^2 \sqrt{n} (\|\mathbf{A}_1\|_\infty + \|\mathbf{A}_2\|_\infty + \|\mathbf{A}_3\|_1) \\ &\leq 4(\alpha_{s,3}^{(h)})^2 k_c \mu c_{w0} n^{0.5} \tilde{r} \left( c \rho \mu^2 k_c^{0.5} c_{u0}^2 (1 + 2\alpha_3 c_{x0}) + \frac{2ck_c^{0.5} c_{w0} \rho c_{u0} \sqrt{p}}{\alpha_3} \right) \\ &\quad + 8(\alpha_{s,3}^{(h)})^2 n c_1 \mu^2 c^2 \alpha_3 c_{w0}^5 c_{x0} \rho k_c m^{0.5} \tilde{r} \left(1 - \frac{\eta\lambda}{2}\right)^{t/2} \|\mathbf{u}^0 - \mathbf{y}\|_2.\end{aligned}$$

In this way, we only need to upper bound  $\|\mathbf{G}^0(k) - \mathbf{G}^0(0)\|_2$ ,  $\|\mathbf{G}^{ls}(k) - \mathbf{G}^{ls}(0)\|_2$  and  $\|\mathbf{G}^s(k) - \mathbf{G}^s(0)\|_2$ .

**Step 3. Bound of  $\|\bar{\mathbf{G}}^0(k) - \bar{\mathbf{G}}^0(0)\|_2$ .**

Here we use the same method when we bound  $\|\mathbf{G}^{ls}(k) - \mathbf{G}^{ls}(0)\|_2$  to bound  $\|\mathbf{G}^0(k) - \mathbf{G}^0(0)\|_2$ . Let  $\mathbf{H}_i = \Phi(\mathbf{X}_i)$ ,  $\mathbf{H}_{i,:t} = [\mathbf{H}_i]_{:,t}$ ,  $\mathbf{H}_{i,tr} = [\mathbf{H}_i]_{t,r}$ ,  $\mathbf{Z}_{i,tr} = (\mathbf{W}_{s,:r}^{(0)})^\top \mathbf{H}_{i,:t}$  and  $\mathbf{B}_i(k) = \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(k)}$ . In this way, for  $1 \leq s \leq h-1$  we can write  $\mathbf{G}_{ij}^{hs}$  as Then we define

$$\begin{aligned}\mathbf{A}_1 &= \sum_{t=1}^p \sum_{q=1}^p \left( (\mathbf{H}_{i,:t}(k))^\top \mathbf{H}_{j,:q}(k) - (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \right) \sum_{r=1}^m \mathbf{B}_{i,tr}(0) \mathbf{B}_{j,qr}(0) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)), \\ \mathbf{A}_2 &= \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m \mathbf{B}_{i,tr}(0) \mathbf{B}_{j,qr}(0) \left( \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)) - \sigma'(\mathbf{Z}_{i,tr}(0)) \sigma'(\mathbf{Z}_{j,qr}(0)) \right), \\ \mathbf{A}_{3,ij} &= \sum_{t=1}^p \sum_{q=1}^p (\mathbf{H}_{i,:t}(0))^\top \mathbf{H}_{j,:q}(0) \sum_{r=1}^m (\mathbf{B}_{i,tr}(k) \mathbf{B}_{j,qr}(k) - \mathbf{B}_{i,tr}(0) \mathbf{B}_{j,qr}(0)) \sigma'(\mathbf{Z}_{i,tr}(k)) \sigma'(\mathbf{Z}_{j,qr}(k)).\end{aligned}$$

Then by using the same method, we can prove

$$\begin{aligned}\|\bar{\mathbf{G}}^0(k) - \bar{\mathbf{G}}^0(0)\|_2 &= \tau^2 \|\mathbf{A}_1 + \mathbf{A}_2 + \mathbf{A}_3\|_2 \leq (\alpha_{s,3}^{(h)} \tau)^2 (\|\mathbf{A}_1\|_2 + \|\mathbf{A}_2\|_2 + \|\mathbf{A}_3\|_2) \\ &\leq \tau^2 \sqrt{n} (\|\mathbf{A}_1\|_\infty + \|\mathbf{A}_2\|_\infty + \|\mathbf{A}_3\|_1) \\ &\leq 4k_c \mu c_{w0} n^{0.5} \tilde{r} \left( c \rho \mu^2 k_c^{0.5} c_{u0}^2 (1 + 2\alpha_3 c_{x0}) + \frac{2ck_c^{0.5} c_{w0} \rho c_{u0} \sqrt{p}}{\alpha_3} \right) \\ &\quad + 8c_1 n \mu^2 c^2 \alpha_3 c_{w0}^5 c_{x0} \rho k_c m^{0.5} \tilde{r} \left(1 - \frac{\eta\lambda}{2}\right)^{k/2} \|\mathbf{u}^0 - \mathbf{y}\|_2.\end{aligned}$$

**Step 4. Bound of  $\|\mathbf{G}(k) - \mathbf{G}(0)\|_2$ .**

By combining the above results and ignoring all constants for brevity, we can bound

$$\begin{aligned}\|\mathbf{G}(k) - \mathbf{G}(0)\|_2 &\leq \|\bar{\mathbf{G}}^0(k) - \bar{\mathbf{G}}^0(0)\|_2 + \sum_{l=0}^{h-1} \sum_{s=0}^{l-1} \|\mathbf{G}^{ls}(k) - \mathbf{G}^{ls}(0)\|_2 + \sum_{s=0}^{h-1} \|\mathbf{G}^s(k) - \mathbf{G}^s(0)\|_2 \\ &\leq c_2 c h \mu k_c^{0.5} c_{x0} \tilde{r} n^{0.5} (\rho h \mu^2 k_c c_{u0}^2 c_{w0} + \alpha_3 c \rho h \mu k_c^{0.5} c_{w0}^5 n^{0.5})\end{aligned}$$

where  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^h$  and  $c_2$  is a constant. Considering

$$\tilde{r} = \frac{8c_{x0} \|\mathbf{y} - \mathbf{u}(0)\|_2}{\lambda \sqrt{mn}} \max \left( 1, 2 \left( 1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}} \right)^h \alpha_3 \mu \sqrt{k_c c_{w0}} \right) \leq c_{w0},$$

to achieve

$$\|\mathbf{G}(k) - \mathbf{G}(0)\|_2 \leq \frac{\lambda}{2},$$

$m$  should be at the order of

$$m \geq \frac{c_3 \alpha_3^2 \mu^2 k_c c_{w0}^2 c^2}{\lambda^2 n},$$

where  $c_3$  is a constant,  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^h$ ,  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . The proof is completed.  $\square$

### C.3.3 Proof of Lemma 21

*Proof.* Lemma 19 proves that when  $m = \mathcal{O}\left(\frac{\rho k_c^2 c_{w0}^2 \|\mathbf{y} - \mathbf{u}(0)\|_2^2}{\lambda^2 n} (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^{2h}\right)$ , then with probability at least  $1 - \delta/2$  we have

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq \left(1 - \frac{\eta\lambda}{2}\right) \|\mathbf{y} - \mathbf{u}(k-1)\|_2^2 \leq \left(1 - \frac{\eta\lambda}{2}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|_2^2,$$

where  $\lambda$  is smallest eigenvalue of the Gram matrix  $\mathbf{G}(t)$  ( $t = 1, \dots, k-1$ ). Lemma 20 shows that if  $m$  satisfies  $m \geq \frac{c_3 \alpha_3^2 \mu^2 k_c c_{w0}^2 c^2}{\lambda^2 n}$ , where  $c_3$  is a constant,  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^h$ ,  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ , then we have

$$\|\mathbf{G}(k) - \mathbf{G}(0)\|_2 \leq \frac{\lambda_{\min}(\mathbf{G}(0))}{2},$$

where  $\lambda_{\min}(\mathbf{G}(0))$  is the smallest eigenvalue of  $\mathbf{G}(0)$ . So we have

$$\lambda_{\min}(\mathbf{G}(t)) \geq \frac{\lambda_{\min}(\mathbf{G}(0))}{2}.$$

So combining these results, we have

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq \left(1 - \frac{\eta\lambda_{\min}(\mathbf{G}(0))}{4}\right) \|\mathbf{y} - \mathbf{u}(k-1)\|_2^2 \leq \left(1 - \frac{\eta\lambda_{\min}(\mathbf{G}(0))}{4}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|_2^2,$$

when  $m$  satisfies  $m \geq \frac{c'_m \rho c^2 k_c^2 c_{w0}^2 \mu^2}{\lambda^2 n}$  and  $\eta \leq \frac{c'_\eta \lambda}{\sqrt{m} \mu^4 h^3 k_c^2 c^4}$ , where  $c'_m, c'_\eta$  are constants,  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c c_{w0}})^h$ ,  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . The proof is completed.  $\square$

## C.4 Step 2 Lower Bound of Eigenvalue of Gram Matrix

Here we define some necessary notations for this subsection first. By Gaussian distribution  $\mathcal{P}$  over a  $q$ -dimensional subspace  $\mathcal{W}$ , it means that for a basis  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_q\}$  of  $\mathcal{W}$  and  $(v_1, v_2, \dots, v_q) \sim \mathcal{N}(0, \mathbf{I})$  such that  $\sum_{i=1}^q v_i \mathbf{e}_i \sim \mathcal{P}$ . Then we equip one Gaussian distribution  $\mathcal{P}^{(i)}$  with each linear subspace  $\mathcal{W}$ . Based on these, we define a transform  $\mathcal{W}$  as

$$\mathcal{W}_{tq}^{(ls)}(\mathbf{K}) = \begin{cases} \mathbb{E}_{\mathbf{W}_t^{(l)} \sim \mathcal{P}}[\mathbf{W}_t^{(l)} \mathbf{K} (\mathbf{W}_t^{(l)})^\top], & \text{if } l = s \text{ and } t = q \\ \mathbb{E}_{\mathbf{W}_t^{(l)} \sim \mathcal{P}, \mathbf{W}_q^{(s)} \sim \mathcal{P}}[\mathbf{W}_t^{(l)} \mathbf{K} (\mathbf{W}_q^{(s)})^\top], & \text{otherwise} \end{cases},$$

where  $\mathbf{K} \in \mathbb{R}^{p \times p}$  and  $\mathbf{W}_t^{(l)}$  denotes the parameters in convolution.

Then we define the population Gram matrix as follows. For brevity, let  $\bar{\mathbf{X}} = \Phi(\mathbf{X}) \in \mathbb{R}^{k_c m \times p}$ . We first define the case where  $l = 0$ :

$$\begin{aligned} \mathbf{b}_i^{(-1)} &= \mathbf{0} \in \mathbb{R}^p, & \mathbf{K}_{ij}^{(-1)} &= \mathbf{X}_i^\top \mathbf{X}_j, & \mathbf{Q}_{ij}^{(-1)} &= \bar{\mathbf{X}}_i^\top \bar{\mathbf{X}}_j \in \mathbb{R}^{p \times p}, \\ \mathbf{A}^{(00)} &= \begin{bmatrix} \mathcal{W}^{(0)}(\mathbf{Q}_{ij}^{(-1)}), \mathcal{W}^{(0)}(\mathbf{Q}_{ij}^{(-1)}) \\ \mathcal{W}^{(0)}(\mathbf{Q}_{ji}^{(-1)}), \mathcal{W}^{(0)}(\mathbf{Q}_{jj}^{(-1)}) \end{bmatrix}, & (\mathbf{M}^{(00)}, \mathbf{N}^{(00)}) &\sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{(00)}) \\ \mathbf{b}_i^{(0)} &= \tau \mathbb{E}_{\mathbf{M}^{(00)}} \sigma(\mathbf{M}^{(00)}), & \mathbf{K}_{ij}^{(00)} &= \mathbb{E}_{(\mathbf{M}^{(00)}, \mathbf{N}^{(00)})} \left( \sigma(\mathbf{M}^{(00)}) \sigma(\mathbf{N}^{(00)})^\top \right), \\ \mathbf{Q}_{ij,ab}^{(00)} &= \text{Tr} \left( \mathbf{K}_{ij, S_a^{(l)}, S_b^{(s)}}^{(00)} \right), \end{aligned}$$

where  $\mathcal{W}^{(0)}(\mathbf{K}) = \mathbb{E}_{\mathbf{W}^{(0)} \sim \mathcal{P}}[\mathbf{W}^{(0)} \mathbf{K} (\mathbf{W}^{(0)})^\top]$ ,  $\mathbf{Q}_{ij}^{(00)} \in \mathbb{R}^{p \times p}$ ,  $\mathbf{K}_{ij,ab}^{(00)}$  denotes the  $(a, b)$ -th entry in  $\mathbf{K}_{ij}^{(00)}$ , and  $S_a^{(0)} = \{j \mid \mathbf{X}_{:,j} \in \text{the } a\text{-th patch for convolution}\}$ .

Then for  $1 \leq l \leq h, 1 \leq s \leq l$ , we can recurrently define

$$\begin{aligned} \mathbf{A}_{tq}^{(ls)} &= \begin{bmatrix} \mathcal{W}_{tq}^{(ls)}(\mathbf{Q}_{ii}^{(tq)}), \mathcal{W}_{tq}^{(ls)}(\mathbf{Q}_{ij}^{(tq)}) \\ \mathcal{W}_{tq}^{(ls)}(\mathbf{Q}_{ji}^{(tq)}), \mathcal{W}_{tq}^{(ls)}(\mathbf{Q}_{jj}^{(tq)}) \end{bmatrix}, \quad (\mathbf{M}_{tq}^{(ls)}, \mathbf{N}_{tq}^{(ls)}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_{tq}^{(ls)}), \quad (0 \leq t, q \leq l-1), \\ \mathbf{b}_i^{(l)} &= \sum_{t=1}^{l-1} \left( \alpha_{t,2}^{(l)} \mathbf{b}_i^{(t)} + \tau \alpha_{t,3}^{(l)} \mathbb{E}_{\mathbf{M}_{tt}^{(ll)}} \sigma(\mathbf{M}_{tt}^{(ll)}) \right); \\ \mathbf{K}_{ij}^{(ls)} &= \sum_{t=1}^{l-1} \sum_{q=1}^{s-1} \left[ \alpha_{t,2}^{(l)} \alpha_{q,2}^{(s)} \mathbf{K}_{ij}^{(tq)} + \tau \mathbb{E}_{(\mathbf{M}_{tq}^{(ls)}, \mathbf{N}_{tq}^{(ls)})} \left( \alpha_{t,3}^{(l)} \alpha_{q,2}^{(s)} \sigma(\mathbf{M}_{tq}^{(ls)}) (\mathbf{b}_j^{(q)})^\top + \alpha_{t,2}^{(l)} \alpha_{q,3}^{(s)} \mathbf{b}_i^{(t)} \sigma(\mathbf{N}_{tq}^{(ls)})^\top \right. \right. \\ &\quad \left. \left. + \tau \alpha_{t,3}^{(l)} \alpha_{q,3}^{(s)} \sigma(\mathbf{M}_{tq}^{(ls)}) \sigma(\mathbf{N}_{tq}^{(ls)})^\top \right) \right], \\ \mathbf{Q}_{ij,ab}^{(ls)} &= \text{Tr} \left( \mathbf{K}_{ij, S_a^{(l)}, S_b^{(s)}}^{(ls)} \right), \end{aligned}$$

where  $\mathbf{K}_{ij}^{(ls)} \in \mathbb{R}^{p \times p}$ ,  $\mathbf{Q}_{ij,ab}^{(ls)}$  denotes the  $(a, b)$ -th entry in  $\mathbf{Q}_{ij}^{(ls)}$ , and  $S_a^{(s)} = \{j \mid \mathbf{X}_{:,j}^{(s-1)} \in \text{the } a\text{-th patch for convolution}\}$ . Finally, we define

$$\begin{aligned} \mathbf{A}^{(s)} &= \begin{bmatrix} \mathcal{W}_{ss}^{(hs)}(\mathbf{Q}_{ii}^{(ss)}), \mathcal{W}_{ss}^{(hs)}(\mathbf{Q}_{ij}^{(ss)}) \\ \mathcal{W}_{ss}^{(hs)}(\mathbf{Q}_{ji}^{(ss)}), \mathcal{W}_{ss}^{(hs)}(\mathbf{Q}_{jj}^{(ss)}) \end{bmatrix}, \\ \mathbf{Q}_{ij,ab}^{(s)} &= \mathbf{Q}_{ij,ab}^{(ss)} \mathbb{E}_{(\mathbf{M}, \mathbf{N}) \sim \bar{\mathbf{A}}^{(s)}} \sigma'(\mathbf{M}) \sigma'(\mathbf{N})^\top, \quad \mathbf{K}_{ij,ab}^{(s)} = \text{Tr} \left( \mathbf{Q}_{ij}^{(s)} \right), \quad (s = 0, h-1). \end{aligned}$$

For brevity, we first define

$$\widehat{\mathbf{K}}_{ij}^{(ls)} = \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} (\mathbf{X}_{j,t}^{(s)})^\top, \quad \widehat{\mathbf{b}}_i^{(l)} = \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)}.$$

Then we prove that  $\mathbf{K}^{(s)}$  is very close to the randomly generated gram matrix  $\widehat{\mathbf{K}}_{ij}^{(ls)}$ .

**Lemma 22.** *With probability at least  $1 - \delta$  over the convolution parameters  $\mathbf{W}$  in each layer, then for  $0 \leq t \leq h, 0 \leq s \leq h$ , it holds*

$$\left\| \frac{1}{m} \sum_{s=1}^m (\mathbf{X}_{i,s}^{(t)})^\top \mathbf{X}_{j,s}^{(q)} - \mathbf{K}_{ij}^{(tq)} \right\|_\infty \leq C \sqrt{\frac{\log(n^2 p^2 h^2 / \delta)}{m}},$$

and

$$\left\| \frac{1}{m} \sum_{s=1}^m \mathbf{X}_{i,s}^{(t)} - \mathbf{b}_i^{(t)} \right\|_\infty \leq C \sqrt{\frac{\log(n^2 p^2 h^2 / \delta)}{m}},$$

where  $C$  is a constant which depends on the activation function  $\sigma(\cdot)$ , namely  $C \sim \sigma(0) + \sup_x \sigma'(x)$ .

See its proof in Appendix C.4.1.

**Lemma 23.** *Suppose Assumptions 1, 2 and 3 hold. Then if  $m \geq \frac{c_4 \mu^2 p^2 n^2 \log(n/\delta)}{\lambda^2}$ , we have*

$$\left\| \mathbf{G}^{hs}(0) - (\alpha_{s,3}^{(h)})^2 \mathbf{K}^{(s)} \right\|_{op} \leq \frac{\lambda}{4} \quad (s = 0, \dots, h),$$

where  $c_4$  and  $\lambda$  are constants.

See its proof in Appendix C.4.2.

**Lemma 24.** *Suppose Assumptions 1, 2 and 3 hold. Suppose  $\sigma$  is analytic and not a polynomial function. Consider data  $\{\mathbf{X}_{i=1}^n\}_{i=1}^n$  are not parallel, namely  $\text{vec}(\mathbf{X}_i) \notin \text{span}(\text{vec}(\mathbf{X}_j))$  for all  $i \neq j$ . Then if  $m \geq \frac{c_4 \mu^2 p^2 n^2 \log(n/\delta)}{\lambda^2}$ , it holds that with probability at least  $1 - \delta/2$ , the smallest eigenvalue the matrix  $\mathbf{G}$  satisfies*

$$\lambda_{\min}(\mathbf{G}(0)) \geq \frac{3c\sigma}{4} \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) \lambda_{\min}(\mathbf{K}).$$

where  $\lambda = 3c_\sigma \sum_{s=0}^{h-1} (\alpha_{s,3})^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2})^2 \right) \lambda_{\min}(\mathbf{K})$ ,  $c_\sigma$  is a constant that only depends on  $\sigma$  and the input data,  $\lambda_{\min}(\mathbf{K}) = \min_{i,j} \lambda_{\min}(\mathbf{K}_{ij})$  is larger than zero in which  $\lambda_{\min}(\mathbf{K}_{ij})$  is the smallest eigenvalue of  $\mathbf{K}_{ij} = \begin{bmatrix} \mathbf{X}_i^\top \mathbf{X}_j & \mathbf{X}_i^\top \mathbf{X}_j \\ \mathbf{X}_j^\top \mathbf{X}_i & \mathbf{X}_j^\top \mathbf{X}_j \end{bmatrix}$ .

See its proof in C.4.3.

#### C.4.1 Proof of Lemma 22

*Proof.* We use mathematical induction to prove these results. For brevity, let  $\bar{\mathbf{X}} = \Phi(\mathbf{X}) \in \mathbb{R}^{k_c m \times p}$  and  $\mathbf{X}_{i,s} = \mathbf{X}_{i,s}^\top \in \mathbb{R}^p$ . For the first layer ( $l = 0$ ), we have

$$\mathbf{X}_{i,s}^{(0)} = \tau\sigma \left( \sum_{t=1}^m \mathbf{W}_{ts}^{(0)} \bar{\mathbf{X}}_{i,t} \right) \quad (18)$$

Then let

$$\mathbf{A}_{i,s}^{(0)} = \sum_{t=1}^m \mathbf{W}_{ts}^{(0)} \bar{\mathbf{X}}_{i,t}. \quad (19)$$

Since the convolution parameter  $\mathbf{W}$  satisfies Gaussian distribution,  $\mathbf{A}_{i,s}^{(0)}$  is a mean-zero Gaussian variable with covariance matrix as follows

$$\mathbb{E} \left[ (\mathbf{A}_{i,s}^{(0)})^\top \mathbf{A}_{j,q}^{(0)} \right] = \mathbb{E} \sum_{t,t'} \mathbf{W}_{ts}^{(0)} \bar{\mathbf{X}}_{i,t}^{(0)} (\bar{\mathbf{X}}_{j,t'})^\top (\mathbf{W}_{t'q}^{(0)})^\top = \delta_{st} \mathcal{W}^{(0)} \left( \sum_t \bar{\mathbf{X}}_{i,t} \bar{\mathbf{X}}_{j,t}^\top \right) = \delta_{st} \mathcal{W}^{(0)} (\mathbf{Q}_{ij}^{(-1)}),$$

where  $\delta_{st}$  is a random variable with  $\delta_{st} = \pm 1$  with both probability 0.5. Therefore, we have

$$\mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \mathbf{X}_{i,t}^{(0)} (\mathbf{X}_{j,t}^{(0)})^\top \right] = \mathbf{K}_{ij}^{(00)}, \quad \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \mathbf{X}_{i,t}^{(0)} \right] = \mathbf{b}_i^{(0)}.$$

In this way, following [4] we can apply Hoeffding and Bernstein bounds and obtain the following results:

$$\mathbb{P} \left( \max_{ij} \left\| \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(0)} (\mathbf{X}_{j,t}^{(0)})^\top - \mathbf{K}_{ij}^{(00)} \right\|_\infty \leq \sqrt{\frac{16(1 + 2C_1^2/\sqrt{\pi})M^2 \log(4n^2 p^2 h^2 / \delta)}{m}} \right) \geq 1 - \frac{\delta}{h^2},$$

where we use  $\|\mathbf{X}_{i,t}^{(0)} (\mathbf{X}_{j,t}^{(0)})^\top\|_2 \leq \|\mathbf{X}_{i,t}^{(0)} (\mathbf{X}_{j,t}^{(0)})^\top\|_F \leq 0.5(\|\mathbf{X}_{i,t}^{(0)}\|_F^2 + \|\mathbf{X}_{j,t}^{(0)}\|_F^2) \stackrel{\textcircled{1}}{\leq} c_{x0}^2$ ,  $M_1 = 1 + 100 \max_{i,j,s,t,l} |\mathcal{W}^0(\mathbf{Q}_{ij}^{(-1)})_{st}|$ . Here  $\textcircled{1}$  holds by using Lemma 10. Similarly, we can prove

$$\mathbb{P} \left( \left\| \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(1)} - \mathbf{b}_i^{(1)} \right\|_\infty \leq \sqrt{\frac{2C_1 M \log(2nph/\delta)}{m}} \right) \geq 1 - \delta/h^2.$$

Then we prove the results still hold when  $l \geq 1, l \geq s \geq 0$ . For brevity, we first define

$$\widehat{\mathbf{K}}_{ij}^{(ls)} = \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} (\mathbf{X}_{j,t}^{(s)})^\top, \quad \widehat{\mathbf{b}}_i^{(l)} = \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)}.$$

Suppose the results in our lemma holds for  $0 \leq l \leq k, 0 \leq q \leq l$  with probability at least  $1 - \frac{k^2}{h^2} \delta$ . For  $l = k + 1$ , we need to prove the results still hold with probability at least  $1 - \frac{2l-1}{h^2} \delta$ . Toward this goal, we have

$$\mathbf{X}_{i,s}^{(l)} = \sum_{0 \leq q \leq l-1} \left[ \mathbf{X}_{i,s}^{(q)} + \tau\sigma \left( \sum_{t=1}^m \mathbf{W}_{q,ts}^{(l)} \bar{\mathbf{X}}_{i,t}^{(q)} \right) \right],$$

where  $\tau = \frac{1}{\sqrt{m}}$ . Then let

$$\mathbf{A}_{i,s}^{(lq)} = \sum_{t=1}^m \mathbf{W}_{q,ts}^{(l)} \bar{\mathbf{X}}_{i,t}^{(q)}.$$

Similarly, we can obtain  $\mathbf{A}_{i,s}^{(lq)}$  is a mean-zero Gaussian variable with covariance matrix

$$\mathbb{E} \left[ \mathbf{A}_{i,s}^{(lq)} (\mathbf{A}_{i,s}^{(lr)})^\top \right] = \delta_{st} \mathcal{W}_{qr}^{(l)} \left( \sum_t \bar{\mathbf{X}}_{i,t}^{(q)} (\bar{\mathbf{X}}_{j,t}^{(r)})^\top \right) = \delta_{st} \mathcal{W}_{qr}^{(l)} (\widehat{\mathbf{Q}}_{ij}^{qr}).$$

Note that since for convolution networks, each element in the output involves several elements in the input (implemented by the operation  $\Phi(\cdot)$ ), we need to consider this by combining the involved elements. Therefore, we can conclude

$$\widehat{\mathbf{Q}}_{ij,ab}^{(ls)} = \text{Tr} \left( \widehat{\mathbf{K}}_{ij, S_a^{(l)}, S_b^{(s)}}^{(ls)} \right) \quad (1 \leq s \leq l)$$

where  $\widehat{\mathbf{K}}_{ij,ab}^{(ls)}$  denotes the  $(a, b)$ -th entry in  $\widehat{\mathbf{K}}_{ij}^{(ls)}$ , and  $S_a^{(s)} = \{j \mid \mathbf{X}_{:,j}^{(s-1)} \in \text{the } a\text{-th patch}\}$ . Moreover, we can easily obtain

$$\mathbb{E} \left[ \widehat{\mathbf{b}}_i^{(l)} \right] = \sum_{t=1}^{l-1} \left( \alpha_{t,2}^{(l)} \widehat{\mathbf{b}}_i^{(t)} + \tau \alpha_{t,3}^{(l)} \mathbb{E}_{\widehat{\mathbf{M}}_{tt}^{(l)}} \sigma(\widehat{\mathbf{M}}_{tt}^{(l)}) \right).$$

In this way, we can further obtain

$$\widehat{\mathbf{A}}_{tq}^{(l)} = \begin{bmatrix} \mathcal{W}_{tq}^{(l)}(\widehat{\mathbf{Q}}_{ii}^{(tq)}), \mathcal{W}_{tq}^{(l)}(\widehat{\mathbf{Q}}_{ij}^{(tq)}) \\ \mathcal{W}_{tq}^{(l)}(\widehat{\mathbf{Q}}_{ji}^{(tq)}), \mathcal{W}_{tq}^{(l)}(\widehat{\mathbf{Q}}_{jj}^{(tq)}) \end{bmatrix}, \quad (\widehat{\mathbf{M}}_{tq}^{(l)}, \widehat{\mathbf{N}}_{tq}^{(l)}) \sim \mathcal{N} \left( \mathbf{0}, \widehat{\mathbf{A}}_{tq}^{(l)} \right), \quad (0 \leq t, q \leq l-1),$$

$$\mathbb{E} \left[ \widehat{\mathbf{K}}_{ij}^{(ls)} \right] = \sum_{t=1}^{l-1} \sum_{q=1}^{s-1} \left[ \alpha_{t,2}^{(l)} \alpha_{q,2}^{(s)} \widehat{\mathbf{K}}_{ij}^{(tq)} + \tau \mathbb{E}_{(\widehat{\mathbf{M}}_{tq}^{(l)}, \widehat{\mathbf{N}}_{tq}^{(l)})} \left( \alpha_{t,3}^{(l)} \alpha_{q,2}^{(s)} \sigma(\widehat{\mathbf{M}}_{tq}^{(l)}) (\widehat{\mathbf{b}}_j^{(q)})^\top + \alpha_{t,2}^{(l)} \alpha_{q,3}^{(s)} \widehat{\mathbf{b}}_i^{(t)} \sigma(\widehat{\mathbf{N}}_{tq}^{(l)})^\top \right. \right. \\ \left. \left. + \tau \alpha_{t,3}^{(l)} \alpha_{q,3}^{(s)} \sigma(\widehat{\mathbf{M}}_{tq}^{(l)}) \sigma(\widehat{\mathbf{N}}_{tq}^{(l)})^\top \right) \right] \in \mathbb{R}^{p \times p}.$$

Then we also apply the concentration inequality and obtain that for  $1 \leq s \leq l$

$$\mathbb{P} \left( \max_{ij} \left\| \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} (\mathbf{X}_{j,t}^{(s)})^\top - \mathbb{E} \widehat{\mathbf{K}}_{ij}^{(ls)} \right\|_\infty \leq \sqrt{\frac{16(1+2C_1^2/\sqrt{\pi})M^2 \log(4n^2 p^2 h^2 / \delta)}{m}} \right) \geq 1 - \delta/h^2$$

where we use  $\|\mathbf{X}_{i,t}^{(0)} (\mathbf{X}_{j,t}^{(0)})^\top\|_2 \leq \|\mathbf{X}_{i,t}^{(0)} (\mathbf{X}_{j,t}^{(0)})^\top\|_F \leq 0.5(\|\mathbf{X}_{i,t}^{(0)}\|_F^2 + \|\mathbf{X}_{j,t}^{(0)}\|_F^2) \leq c_{x0}^2$ ,  $M_1 = 1 + 100 \max_{i,j,s,t,l} |\mathcal{W}^l(\widehat{\mathbf{K}}_{ij}^{(l-1)})_{st}|$ . Similarly, we can prove

$$\mathbb{P} \left( \left\| \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} - \mathbb{E} \widehat{\mathbf{b}}_i^{(l)} \right\|_\infty \leq \sqrt{\frac{2C_1 M \log(2nph/\delta)}{m}} \right) \geq 1 - \delta/h^2.$$

According to the definition

$$\widehat{\mathbf{K}}_{ij}^{(ls)} = \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} (\mathbf{X}_{j,t}^{(s)})^\top, \quad \widehat{\mathbf{b}}_i^{(l)} = \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)}.$$

we have

$$\left\| \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} (\mathbf{X}_{j,t}^{(s)})^\top - \mathbf{K}_{ij}^{(ls)} \right\|_\infty \leq \left\| \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} (\mathbf{X}_{j,t}^{(s)})^\top - \mathbb{E} \widehat{\mathbf{K}}_{ij}^{(ls)} \right\|_\infty + \left\| \mathbb{E} \widehat{\mathbf{K}}_{ij}^{(ls)} - \mathbf{K}_{ij}^{(ls)} \right\|_\infty, \\ \left\| \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} - \mathbf{b}_i^{(l)} \right\|_\infty \leq \left\| \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(l)} - \mathbb{E} \widehat{\mathbf{b}}_i^{(l)} \right\|_\infty + \left\| \mathbb{E} \widehat{\mathbf{b}}_i^{(l)} - \mathbf{b}_i^{(l)} \right\|_\infty.$$

Then we only need to bound

$$\left\| \mathbb{E} \widehat{\mathbf{K}}_{ij}^{(ls)} - \mathbf{K}_{ij}^{(ls)} \right\|_\infty \quad \text{and} \quad \left\| \mathbb{E} \widehat{\mathbf{b}}_i^{(l)} - \mathbf{b}_i^{(l)} \right\|_\infty.$$

In the following content, we bound these two terms in turn. To begin with, we have

$$\left\| \mathbb{E} \widehat{\mathbf{K}}_{ij}^{(ls)} - \mathbf{K}_{ij}^{(ls)} \right\|_\infty = \left\| \text{Tr} \left( \widehat{\mathbf{Q}}_{ij, S_a^{(l)}, S_b^{(s)}}^{(ls)} \right) - \text{Tr} \left( \mathbf{Q}_{ij, S_a^{(s)}, S_b^{(ls)}}^{(ls)} \right) \right\|_\infty \leq \left\| \widehat{\mathbf{Q}}_{ij}^{(l)} - \mathbf{Q}_{ij}^{(l)} \right\|_\infty \\ \leq \sum_{t=1}^{l-1} \sum_{q=1}^{s-1} \left[ \alpha_{t,2}^{(l)} \alpha_{q,2}^{(s)} \left\| \widehat{\mathbf{K}}_{ij}^{(tq)} - \mathbf{K}_{ij}^{(tq)} \right\|_\infty \right. \\ \left. + \tau \alpha_{t,3}^{(l)} \alpha_{q,2}^{(s)} \left\| \mathbb{E}_{(\widehat{\mathbf{M}}^{(tq)}, \widehat{\mathbf{N}}^{(tq)})} \sigma(\widehat{\mathbf{M}}^{(tq)}) (\widehat{\mathbf{b}}_j^{(q)})^\top - \mathbb{E}_{((\mathbf{M}^{(tq)}, \mathbf{N}^{(tq)})} \sigma(\mathbf{M}^{(tq)}) (\mathbf{b}_j^{(q)})^\top \right\|_\infty \right. \\ \left. + \tau \alpha_{t,2}^{(l)} \alpha_{q,3}^{(s)} \left\| \mathbb{E}_{(\widehat{\mathbf{M}}^{(tq)}, \widehat{\mathbf{N}}^{(tq)})} \widehat{\mathbf{b}}_i^{(t)} \sigma(\widehat{\mathbf{N}}^{(tq)})^\top - \mathbb{E}_{((\mathbf{M}^{(tq)}, \mathbf{N}^{(tq)})} \mathbf{b}_i^{(t)} \sigma(\mathbf{N}^{(tq)})^\top \right\|_\infty \right. \\ \left. + \tau \alpha_{t,3}^{(l)} \alpha_{q,3}^{(s)} \left\| \mathbb{E}_{(\widehat{\mathbf{M}}^{(tq)}, \widehat{\mathbf{N}}^{(tq)})} \sigma(\widehat{\mathbf{M}}^{(tq)}) \sigma(\widehat{\mathbf{N}}^{(tq)})^\top - \mathbb{E}_{((\mathbf{M}^{(tq)}, \mathbf{N}^{(tq)})} \sigma(\mathbf{M}^{(tq)}) \sigma(\mathbf{N}^{(tq)})^\top \right\|_\infty \right]$$

Then we bound

$$\begin{aligned}
& \left\| \mathbb{E}_{((\widehat{\mathbf{M}}^{(tq)}, \widehat{\mathbf{N}}^{(tq)})} \sigma(\widehat{\mathbf{M}}^{(tq)}) (\widehat{\mathbf{b}}_j^{(q)})^\top - \mathbb{E}_{((\mathbf{M}^{(tq)}, \mathbf{N}^{(tq)})} \sigma(\mathbf{M}^{(tq)}) (\mathbf{b}_j^{(q)})^\top \right\|_\infty \\
&= \left\| \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \widehat{\mathbf{A}}^{(tq)})} \sigma(\mathbf{M}) (\widehat{\mathbf{b}}_j^{(q)})^\top - \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \mathbf{A}^{(tq)})} \sigma(\mathbf{M}) (\mathbf{b}_j^{(q)})^\top \right\|_\infty \\
&\leq \left\| \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \widehat{\mathbf{A}}^{(tq)})} \sigma(\mathbf{M}) (\widehat{\mathbf{b}}_j^{(q)} - \mathbf{b}_j^{(q)})^\top \right\|_\infty + \left\| \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \widehat{\mathbf{A}}^{(tq)})} \sigma(\mathbf{M}) - \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \mathbf{A}^{(tq)})} \sigma(\mathbf{M}) \right\| (\mathbf{b}_j^{(q)})^\top \right\|_\infty
\end{aligned}$$

Next, we bound the above inequality by bound each term:

$$\begin{aligned}
& \left\| \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \widehat{\mathbf{A}}^{(tq)})} \sigma(\mathbf{M}) - \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \mathbf{A}^{(tq)})} \sigma(\mathbf{M}) \right\| (\mathbf{b}_j^{(q)})^\top \right\|_\infty \\
&\leq \max_i \|\mathbf{b}_j^{(q)}\|_\infty (\sigma(0) + \sup_x \sigma'(x)) \|\widehat{\mathbf{A}}^{(tq)} - \mathbf{A}^{(tq)}\|_\infty \\
&\leq c_1 c_2 c_3 \|\widehat{\mathbf{Q}}_{ij}^{(tq)} - \mathbf{Q}_{ij}^{(tq)}\|_\infty \\
&= c_1 c_2 c_3 \max_{a,b} \left\| \text{Tr} \left( \widehat{\mathbf{K}}_{ij, S_a^{(l)}, S_b^{(s)}}^{(ls)} \right) - \text{Tr} \left( \mathbf{K}_{ij, S_a^{(l)}, S_b^{(s)}}^{(ls)} \right) \right\|_\infty \\
&\leq c_1 c_2 c_3 q \left\| \widehat{\mathbf{K}}_{ij}^{(l)} - \mathbf{K}_{ij}^{(l)} \right\|_\infty,
\end{aligned}$$

where  $c_1 = \max_l 1 + \|\mathcal{W}_{tq}^{(l)}\|_{L^\infty \rightarrow L^\infty}$ ,  $c_2 = \sigma(0) + \sup_x \sigma'(x)$ ,  $c_3 = \max_{i,q} \|\mathbf{b}_i^{(q)}\|_\infty$ . Similarly, we can bound

$$\left\| \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \widehat{\mathbf{A}}^{(tq)})} \sigma(\mathbf{M}) (\widehat{\mathbf{b}}_j^{(q)} - \mathbf{b}_j^{(q)})^\top \right\|_\infty \leq c_2 \sqrt{c_1 c_4} \|\mathbf{b}_j^{(q)} - \widehat{\mathbf{b}}_j^{(q)}\|_\infty$$

where  $c_4 = \max_{ij} \|\widehat{\mathbf{Q}}_{ij}^{(tq)}\|_\infty \leq q \max_{ij} \|\widehat{\mathbf{K}}_{ij}^{(tq)}\|_\infty \leq q c_{x0}^2$  and  $1 \leq q \leq l-1$ . Therefore we have

$$\begin{aligned}
& \left\| \mathbb{E}_{((\widehat{\mathbf{M}}^{(tq)}, \widehat{\mathbf{N}}^{(tq)})} \sigma(\widehat{\mathbf{M}}^{(tq)}) (\widehat{\mathbf{b}}_j^{(q)})^\top - \mathbb{E}_{((\mathbf{M}^{(tq)}, \mathbf{N}^{(tq)})} \sigma(\mathbf{M}^{(tq)}) (\mathbf{b}_j^{(q)})^\top \right\|_\infty \\
&= (c_1 c_2 c_3 q + c_2 \sqrt{c_1 c_4}) \max \left( \|\widehat{\mathbf{K}}_{ij}^{(tq)} - \mathbf{K}_{ij}^{(tq)}\|_\infty, \|\mathbf{b}_j^{(q)} - \widehat{\mathbf{b}}_j^{(q)}\|_\infty \right).
\end{aligned}$$

By using the same method, we can upper bound

$$\begin{aligned}
& \left\| \mathbb{E}_{((\widehat{\mathbf{M}}^{(tq)}, \widehat{\mathbf{N}}^{(tq)})} \widehat{\mathbf{b}}_i^{(t)} \sigma(\widehat{\mathbf{N}}^{(tq)})^\top - \mathbb{E}_{((\mathbf{M}^{(tq)}, \mathbf{N}^{(tq)})} \mathbf{b}_i^{(t)} \sigma(\mathbf{N}^{(tq)})^\top \right\|_\infty \\
&= (c_1 c_2 c_3 q + c_2 \sqrt{c_1 c_4}) \max \left( \|\widehat{\mathbf{K}}_{ij}^{(tq)} - \mathbf{K}_{ij}^{(tq)}\|_\infty, \|\mathbf{b}_j^{(q)} - \widehat{\mathbf{b}}_j^{(q)}\|_\infty \right).
\end{aligned}$$

Next, we can upper bound

$$\begin{aligned}
& \left\| \mathbb{E}_{((\widehat{\mathbf{M}}^{(tq)}, \widehat{\mathbf{N}}^{(tq)})} \sigma(\widehat{\mathbf{M}}^{(tq)}) \sigma(\widehat{\mathbf{N}}^{(tq)})^\top - \mathbb{E}_{((\mathbf{M}^{(tq)}, \mathbf{N}^{(tq)})} \sigma(\mathbf{M}^{(tq)}) \sigma(\mathbf{N}^{(tq)})^\top \right\|_\infty \\
&= \left\| \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \widehat{\mathbf{A}}^{(tq)})} \sigma(\mathbf{M}^{(tq)}) \sigma(\mathbf{N}^{(tq)})^\top - \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \mathbf{A}^{(tq)})} \sigma(\mathbf{M}^{(tq)}) \sigma(\mathbf{N}^{(tq)})^\top \right\|_\infty \\
&\leq c_\sigma \|\widehat{\mathbf{A}}^{(tq)} - \mathbf{A}^{(tq)}\|_\infty \leq c_\sigma c_1 \|\widehat{\mathbf{Q}}_{ij}^{(tq)} - \mathbf{Q}_{ij}^{(tq)}\|_\infty \leq c_\sigma c_1 q \|\widehat{\mathbf{K}}_{ij}^{(tq)} - \mathbf{K}_{ij}^{(tq)}\|_\infty,
\end{aligned}$$

where  $c_\sigma$  is a constant that only depends on  $\sigma$ . Combing all results yields

$$\begin{aligned}
& \left\| \mathbb{E} \widehat{\mathbf{K}}_{ij}^{(ls)} - \mathbf{K}_{ij}^{(ls)} \right\|_\infty \\
&\leq \sum_{t=1}^{l-1} \sum_{q=1}^{s-1} \left[ (\boldsymbol{\alpha}_{t,2}^{(l)} \boldsymbol{\alpha}_{q,2}^{(s)} + \tau^2 \boldsymbol{\alpha}_{t,3}^{(l)} \boldsymbol{\alpha}_{q,3}^{(s)} c_\sigma c_1 q) \|\widehat{\mathbf{K}}_{ij}^{(tq)} - \mathbf{K}_{ij}^{(tq)}\|_\infty \right. \\
&\quad \left. + \tau (\boldsymbol{\alpha}_{t,2}^{(l)} \boldsymbol{\alpha}_{q,2}^{(s)} + \boldsymbol{\alpha}_{t,3}^{(l)} \boldsymbol{\alpha}_{q,2}^{(s)}) (c_1 c_2 c_3 q + c_2 \sqrt{c_1 c_4}) \max \left( \|\widehat{\mathbf{K}}_{ij}^{(tq)} - \mathbf{K}_{ij}^{(tq)}\|_\infty, \|\mathbf{b}_j^{(q)} - \widehat{\mathbf{b}}_j^{(q)}\|_\infty \right) \right] \\
&\leq c \max_{1 \leq t \leq l-1, 1 \leq q \leq l-1} \left( \|\widehat{\mathbf{K}}_{ij}^{(tq)} - \mathbf{K}_{ij}^{(tq)}\|_\infty, \|\mathbf{b}_j^{(q)} - \widehat{\mathbf{b}}_j^{(q)}\|_\infty \right)
\end{aligned}$$

where  $c_l = \sum_{t=1}^{l-1} \sum_{q=1}^{s-1} \left[ \boldsymbol{\alpha}_{t,2}^{(l)} \boldsymbol{\alpha}_{q,2}^{(s)} + \tau^2 \boldsymbol{\alpha}_{t,3}^{(l)} \boldsymbol{\alpha}_{q,3}^{(s)} c_\sigma c_1 q + \tau (\boldsymbol{\alpha}_{t,2}^{(l)} \boldsymbol{\alpha}_{q,2}^{(s)} + \boldsymbol{\alpha}_{t,3}^{(l)} \boldsymbol{\alpha}_{q,2}^{(s)}) (c_1 c_2 c_3 q + c_2 \sqrt{c_1 c_4}) \right]$ . Since we have assumed that with probability  $1 - (l-1)^2 \delta / h^2$  for  $0 \leq t \leq l-1, 0 \leq s \leq l-1$ , it holds

$$\max \left( \left\| \frac{1}{m} \sum_{s=1}^m (\mathbf{X}_{i,s}^{(t)})^\top \mathbf{X}_{j,s}^{(q)} - \mathbf{K}_{ij}^{(tq)} \right\|_\infty, \left\| \frac{1}{m} \sum_{s=1}^m \mathbf{X}_{i,s}^{(t)} - \mathbf{b}_i^{(t)} \right\|_\infty \right) \leq C_{l-1} \sqrt{\frac{\log(n^2 p^2 h^2 / \delta)}{m}},$$

where  $C$  is a constant. Then with probability  $1 - (l-1)^2\delta/h^2$ , we have for all  $0 \leq s \leq l$

$$\left\| \mathbb{E} \widehat{\mathbf{K}}_{ij}^{(ls)} - \mathbf{K}_{ij}^{(ls)} \right\|_{\infty} \leq c_l C_{l-1} \sqrt{\frac{\log(n^2 p^2 h^2 / \delta)}{m}}.$$

Thus, with probability  $(1 - (l-1)^2\delta/h^2)(1 - \delta/h^2) \geq 1 - l^2\delta/h^2 \geq 1 - \delta$ , we have for all for  $0 \leq t \leq h, 0 \leq s \leq h$

$$\left\| \frac{1}{m} \sum_{s=1}^m (\mathbf{X}_{i,s}^{(t)})^{\top} \mathbf{X}_{j,s}^{(q)} - \mathbf{K}_{ij}^{(tq)} \right\|_{\infty} \leq C \sqrt{\frac{\log(n^2 p^2 h^2 / \delta)}{m}},$$

where  $C = C_0 \prod_{l=1}^h c_l$  is a constant.

Now we consider to bound

$$\begin{aligned} & \left\| \mathbb{E} \widehat{\mathbf{b}}_i^{(l)} - \mathbf{b}_i^{(l)} \right\|_{\infty} \\ &= \left\| \sum_{t=1}^{l-1} \left( \alpha_{t,2}^{(l)} (\widehat{\mathbf{b}}_i^{(t)} - \mathbf{b}_i^{(t)}) + \tau \alpha_{t,3}^{(l)} \left( \mathbb{E}_{M \sim \widehat{\mathbf{A}}^{lt}} \sigma(M) - \mathbb{E}_{M \sim \mathbf{A}^{lt}} \sigma(M) \right) \right) \right\|_{\infty} \\ &\leq \sum_{t=1}^{l-1} \left( \alpha_{t,2}^{(l)} \left\| \widehat{\mathbf{b}}_i^{(t)} - \mathbf{b}_i^{(t)} \right\|_{\infty} + \tau \alpha_{t,3}^{(l)} \left\| \mathbb{E}_{M \sim \widehat{\mathbf{A}}^{(l-1)t}} \sigma(M) - \mathbb{E}_{M \sim \mathbf{A}^{(l-1)t}} \sigma(M) \right\|_{\infty} \right) \\ &\leq \sum_{t=1}^{l-1} \left( \alpha_{t,2}^{(l)} \left\| \widehat{\mathbf{b}}_i^{(t)} - \mathbf{b}_i^{(t)} \right\|_{\infty} + \tau \alpha_{t,3}^{(l)} c_{\sigma} \left\| \widehat{\mathbf{A}}^{(l-1)t} - \mathbf{A}^{(l-1)t} \right\|_{\infty} \right) \\ &\leq \sum_{t=1}^{l-1} \left( \alpha_{t,2}^{(l)} \left\| \widehat{\mathbf{b}}_i^{(t)} - \mathbf{b}_i^{(t)} \right\|_{\infty} + \tau \alpha_{t,3}^{(l)} c_{\sigma} \left\| \widehat{\mathbf{Q}}^{(l-1)t} - \mathbf{Q}^{(l-1)t} \right\|_{\infty} \right) \\ &\leq \sum_{t=1}^{l-1} \left( \alpha_{t,2}^{(l)} + \tau \alpha_{t,3}^{(l)} c_{\sigma} c_1 q \right) \max \left( \left\| \widehat{\mathbf{b}}_i^{(t)} - \mathbf{b}_i^{(t)} \right\|_{\infty}, \left\| \widehat{\mathbf{K}}^{(l-1)t} - \mathbf{K}^{(l-1)t} \right\|_{\infty} \right) \end{aligned}$$

where  $c'_l = \sum_{t=1}^{l-1} \left( \alpha_{t,2}^{(l)} + \tau \alpha_{t,3}^{(l)} c_{\sigma} c_1 q \right)$ . Then with probability  $(1 - (l-1)^2\delta/h)(1 - \delta/h) \geq 1 - \delta$ , we have for all for  $0 \leq t \leq h$

$$\left\| \frac{1}{m} \sum_{s=1}^m \mathbf{X}_{i,s}^{(t)} - \mathbf{b}_i^{(t)} \right\|_{\infty} \leq C \sqrt{\frac{\log(n^2 p^2 h^2 / \delta)}{m}},$$

where  $C = C_0 \prod_{l=1}^h \max(c_l, c'_l)$  is a constant. The proof is completed.  $\square$

#### C.4.2 Proof of Lemma 23

*Proof.* For brevity, here we just use  $\mathbf{X}_i^{(s)}$ ,  $\mathbf{W}_s^{(h)}$ ,  $\mathbf{W}_h$ ,  $\bar{\mathbf{X}}_i^{(s)}$  to respectively denote  $X_{iii}(s)i(0)$ ,  $\mathbf{W}_s^{(h)}(0)$ ,  $\mathbf{W}_h(0)$ ,  $\Phi(\mathbf{X}_i^{(s)})$ , since here we only involve the initialization and does not update the variables. Let  $\bar{\mathbf{X}}_{i,t}^{(s)} = (\bar{\mathbf{X}}_{i,t}^{(s)})^{\top}$  and  $\mathbf{Z}_{i,tr} = (\mathbf{W}_{s,r}^{(h)})^{\top} \bar{\mathbf{X}}_{i,t}^{(s)}$ . Firstly according to the definition, we have

$$\begin{aligned} \mathbf{G}_{ij}^{hs}(0) &= \left\langle \frac{\partial \ell_i}{\partial \mathbf{W}_s^{(h)}(0)}, \frac{\partial \ell_j}{\partial \mathbf{W}_s^{(h)}(0)} \right\rangle \\ &= (\alpha_{s,3}^{(h)} \tau)^2 \left\langle \Phi(\mathbf{X}_i^{(s)}) \left( \sigma' \left( \mathbf{W}_s^{(l)} \Phi(\mathbf{X}_i^{(s)}) \right) \odot \mathbf{W}_h \right)^{\top}, \Phi(\mathbf{X}_j^{(s)}) \left( \sigma' \left( \mathbf{W}_s^{(l)} \Phi(\mathbf{X}_j^{(s)}) \right) \odot \mathbf{W}_h \right)^{\top} \right\rangle \\ &= (\alpha_{s,3}^{(h)} \tau)^2 \sum_{t=1}^p \sum_{q=1}^p \bar{\mathbf{X}}_{i,t}^{(s)} (\bar{\mathbf{X}}_{j,q}^{(s)})^{\top} \sum_{r=1}^m \mathbf{W}_{h,tr} \mathbf{W}_{h,qr} \sigma'(\mathbf{Z}_{i,tr}) \sigma'(\mathbf{Z}_{j,qr}). \end{aligned}$$

Then by taking expectation on  $\mathbf{W} \sim \mathcal{N}(0, \mathbf{I})$  and  $U \sim \mathcal{N}(0, \mathbf{I})$ , we have

$$\begin{aligned} \mathbf{G}_{ij}^{hs}(0) &= (\alpha_{s,3}^{(h)} \tau)^2 \sum_{t=1}^p \sum_{q=1}^p \bar{\mathbf{X}}_{i,t}^{(s)} (\bar{\mathbf{X}}_{j,q}^{(s)})^{\top} \sum_{r=1}^m \mathbb{E}_{\mathbf{W}_h} [\mathbf{W}_{h,tr} \mathbf{W}_{h,qr}] \mathbb{E}_{\mathbf{W}_s^{(h)}} [\sigma'(\mathbf{Z}_{i,tr}) \sigma'(\mathbf{Z}_{j,qr})] \\ &= (\alpha_{s,3}^{(h)} \tau)^2 \sum_{t=1}^p \bar{\mathbf{X}}_{i,t}^{(s)} (\bar{\mathbf{X}}_{j,t}^{(s)})^{\top} \sum_{r=1}^m \mathbb{E}_{\mathbf{W}_s^{(h)}} [\sigma'(\mathbf{Z}_{i,tr}) \sigma'(\mathbf{Z}_{j,qr})] \end{aligned} \tag{20}$$



where ① holds since  $\mathbb{E}_{\mathbf{W}_h} [\mathbf{W}_{h,tr} \mathbf{W}_{h,qr}] = 1$  if  $t = q$  and  $\mathbb{E}_{\mathbf{W}_h} [\mathbf{W}_{h,tr} \mathbf{W}_{h,qr}] = 0$  if  $t \neq q$ .

$$\mathbf{Z}_{i,r} = \sum_{t=1}^m (\mathbf{W}_{s,tr}^{(h)})^\top \bar{\mathbf{X}}_{i,t}^{(s)}.$$

Since the convolution parameter  $\mathbf{W}_s^{(h)}$  satisfies Gaussian distribution,  $\mathbf{Z}_{i,r}$  is a mean-zero Gaussian variable with covariance matrix as follows

$$\begin{aligned} \mathbb{E} \left[ (\mathbf{Z}_{i,r})^\top \mathbf{Z}_{j,q} \right] &= \mathbb{E} \sum_{t,t'} (\mathbf{W}_{s,t}^{(h)})^\top \bar{\mathbf{X}}_{i,t}^{(s)} (\bar{\mathbf{X}}_{j,t'}^{(s)})^\top (\mathbf{W}_{s,t'q}^{(h)})^\top = \delta_{st} \mathcal{W}^{(hs)} \left( \sum_t \bar{\mathbf{X}}_{i,t}^{(s)} (\bar{\mathbf{X}}_{j,t}^{(s)})^\top \right) \\ &= \delta_{st} \mathcal{W}^{(hs)} \left( \hat{\mathbf{Q}}_{ij}^{(s)} \right), \end{aligned} \quad (21)$$

where  $\delta_{st}$  is a random variable with  $\delta_{st} = \pm 1$  with both probability 0.5, and

$$\widehat{\mathbf{K}}_{ij}^{(ss)} = \frac{1}{m} \sum_{t=1}^m \mathbf{X}_{i,t}^{(s)} (\mathbf{X}_{j,t}^{(s)})^\top, \quad \hat{\mathbf{Q}}_{ij}^{(ss)} = \frac{1}{m} \sum_{t=1}^m \bar{\mathbf{X}}_{i,t}^{(s)} (\bar{\mathbf{X}}_{j,t}^{(s)})^\top.$$

According to this definition, we actually have

$$\hat{\mathbf{Q}}_{ij,ab}^{(ss)} = \text{Tr} \left( \widehat{\mathbf{K}}_{ij, S_a^{(s)}, S_b^{(s)}}^{(ss)} \right),$$

where  $\widehat{\mathbf{K}}_{ij}^{(ss)} \in \mathbb{R}^{p \times p}$ ,  $\hat{\mathbf{Q}}_{ij,ab}^{(ss)}$  denotes the  $(a, b)$ -th entry in  $\hat{\mathbf{Q}}_{ij}^{(ss)}$ , and  $S_a^{(s)} = \{j \mid \mathbf{X}_{:,j}^{(s-1)} \in \text{the } a\text{-th patch for convolution}\}$ . Then according to the following definitions

$$\begin{aligned} \hat{\mathbf{A}}^{(s)} &= \begin{bmatrix} \mathcal{W}_{ss}^{(h)}(\hat{\mathbf{Q}}_{ii}^{(ss)}), \mathcal{W}_{ss}^{(h)}(\hat{\mathbf{Q}}_{ij}^{(ss)}) \\ \mathcal{W}_{ss}^{(h)}(\hat{\mathbf{Q}}_{ji}^{(ss)}), \mathcal{W}_{ss}^{(h)}(\hat{\mathbf{Q}}_{jj}^{(ss)}) \end{bmatrix}, \\ \hat{\mathbf{Q}}_{ij,ab}^{(s)} &= \hat{\mathbf{Q}}_{ij,ab}^{(ss)} \mathbb{E}_{((M,N) \sim \hat{\mathbf{A}}^{(s)})} \sigma'(M) \sigma'(N)^\top, \quad \widehat{\mathbf{K}}_{ij,ab}^{(s)} = \text{Tr} \left( \hat{\mathbf{Q}}_{ij}^{(s)} \right), \quad (s = 0, h-1). \end{aligned}$$

and Eqns. (20) and (21), we have

$$\mathbb{E} \left[ \mathbf{G}_{ij}^{hs}(0) \right] = (\alpha_{s,3}^{(h)})^2 \widehat{\mathbf{K}}_{ij}^{(s)}, \quad \mathbb{E} \left[ \mathbf{G}^{hs}(0) \right] = (\alpha_{s,3}^{(h)})^2 \widehat{\mathbf{K}}^{(s)}.$$

In this way, we can apply the Hoeffding inequality and obtain that if  $m \geq \mathcal{O} \left( \frac{n^2 \log(n/\delta)}{\lambda^2} \right)$

$$\left\| \mathbf{G}^{hs}(0) - (\alpha_{s,3}^{(h)})^2 \widehat{\mathbf{K}}^{(s)} \right\|_{\text{op}} \leq \frac{\lambda}{8}.$$

On the other hand, Lemma 22 shows that with probability at least  $1 - \delta$

$$\left\| \widehat{\mathbf{K}}_{ij}^{(ss)} - \mathbf{K}_{ij}^{(ss)} \right\|_{\infty} \leq C \sqrt{\frac{\log(n^2 p^2 h^2 / \delta)}{m}} \stackrel{\text{①}}{\leq} \frac{C_3 \lambda}{n},$$

where ① holds by setting  $m \geq \mathcal{O} \left( \frac{C_3^2 n^2 \log(n^2 p^2 h^2 / \delta)}{\lambda^2} \right)$ . Moreover, Lemma 10 shows

$$\frac{1}{c_{x0}} \leq \|\mathbf{X}^{(l)}(0)\|_F \leq c_{x0}.$$

where  $c_{x0} \geq 1$  is a constant. So  $\|\widehat{\mathbf{K}}_{ij}^{(ss)}\|_{\infty}$  is upper bounded by  $c_{x0}^2$ .

Next, Lemma 7 shows if each diagonal entry in  $\mathbf{A}$  and  $\mathbf{B}$  is upper bounded by  $c$  and lower upper bounded by  $1/c$ , then

$$|g(\mathbf{A}) - g(\mathbf{B})| \leq c \|\mathbf{A} - \mathbf{B}\|_F \leq 2C_1 \|\mathbf{A} - \mathbf{B}\|_{\infty},$$

where  $g(\mathbf{A}) = \mathbb{E}_{(u,v) \sim \mathcal{N}(0,\mathbf{A})} \sigma(u)\sigma(v)$ ,  $C_1$  is a constant that only depends on  $c$  and the Lipschitz and smooth parameter of  $\sigma(\cdot)$ . By applying this lemma, we can obtain

$$\begin{aligned}
& \left| \widehat{\mathbf{Q}}_{ij,rq}^{(ss)} \mathbb{E}_{(M,N) \sim \widehat{\mathbf{A}}^{(s)}} [\sigma'(\mathbf{M}_r)\sigma'(\mathbf{N}_q)] - \mathbf{Q}_{ij,rq}^{(ss)} \mathbb{E}_{(M,N) \sim \bar{\mathbf{A}}^{(s)}} [\sigma'(\mathbf{M}_r)\sigma'(\mathbf{N}_q)] \right| \\
& \leq \left| \widehat{\mathbf{Q}}_{ij,rq}^{(ss)} \left( \mathbb{E}_{(M,N) \sim \widehat{\mathbf{A}}^{(s)}} [\sigma'(\mathbf{M}_r)\sigma'(\mathbf{N}_q)] - \mathbb{E}_{(M,N) \sim \bar{\mathbf{A}}^{(s)}} [\sigma'(\mathbf{M}_r)\sigma'(\mathbf{N}_q)] \right) \right| \\
& \quad + \left| (\widehat{\mathbf{Q}}_{ij,rq}^{(ss)} - \mathbf{Q}_{ij,rq}^{(ss)}) \mathbb{E}_{(M,N) \sim \bar{\mathbf{A}}^{(s)}} [\sigma'(\mathbf{M}_r)\sigma'(\mathbf{N}_q)] \right| \\
& \leq C_1 c_{x0}^2 |\widehat{\mathbf{A}}^{(s)} - \bar{\mathbf{A}}^{(s)}| + \mu^2 |\widehat{\mathbf{Q}}_{ij,rq}^{(ss)} - \mathbf{Q}_{ij,rq}^{(ss)}| \\
& \leq C_1 C_2 c_{x0}^2 \max_{i,j} |\widehat{\mathbf{Q}}_{ij,rq}^{(ss)} - \bar{\mathbf{Q}}_{ij,rq}^{(ss)}| + \mu^2 |\widehat{\mathbf{Q}}_{ij,rq}^{(ss)} - \mathbf{Q}_{ij,rq}^{(ss)}| \\
& \leq (C_1 C_2 c_{x0}^2 + \mu^2) \|\widehat{\mathbf{Q}}_{ij}^{(ss)} - \mathbf{Q}_{ij}^{(ss)}\|_\infty \\
& \leq (C_1 C_2 c_{x0}^2 + \mu^2) \max_{a,b} \left\| \text{Tr} \left( \widehat{\mathbf{K}}_{ij, S_a^{(s)}, S_b^{(s)}}^{(ss)} \right) - \text{Tr} \left( \mathbf{K}_{ij, S_a^{(s)}, S_b^{(s)}}^{(ss)} \right) \right\|_\infty \\
& \leq (C_1 C_2 c_{x0}^2 + \mu^2) p \left\| \widehat{\mathbf{K}}_{ij}^{(ss)} - \mathbf{K}_{ij}^{(ss)} \right\|_\infty,
\end{aligned}$$

where  $C_2 = 1 + \|\mathcal{W}_{ss}^{(h)}\|_{L^\infty \rightarrow L^\infty}$ .

Then we can bound

$$\begin{aligned}
\|\widehat{\mathbf{K}}^{(s)} - \bar{\mathbf{K}}^{(s)}\|_{\text{op}} & \leq \|\widehat{\mathbf{K}}^{(s)} - \bar{\mathbf{K}}^{(s)}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n \left[ \text{Tr} \left( \widehat{\mathbf{Q}}_{ij}^{(s)} \right) - \text{Tr} \left( \mathbf{Q}_{ij}^{(s)} \right) \right]^2} \\
& \leq \sqrt{\sum_{i=1}^n \sum_{j=1}^n p \sum_{r=1}^p \left[ \widehat{\mathbf{Q}}_{ij,rr}^{(s)} - \mathbf{Q}_{ij,rr}^{(s)} \right]^2} \\
& \leq \sqrt{\sum_{i=1}^n \sum_{j=1}^n p \sum_{r=1}^p \left[ \widehat{\mathbf{Q}}_{ij,rr}^{(ss)} \mathbb{E}_{(M,N) \sim \widehat{\mathbf{A}}^{(s)}} \sigma'(\mathbf{M}_r) \sigma'(\mathbf{N}_r)^\top - \mathbf{Q}_{ij,rr}^{(ss)} \mathbb{E}_{(M,N) \sim \bar{\mathbf{A}}^{(s)}} \sigma'(\mathbf{M}_r) \sigma'(\mathbf{N}_r)^\top \right]^2} \\
& \leq \sqrt{\sum_{i=1}^n \sum_{j=1}^n p^2 \sum_{r=1}^p (C_1 C_2 c_{x0}^2 + \mu^2)^2 \|\widehat{\mathbf{K}}_{ij}^{(ss)} - \bar{\mathbf{K}}_{ij}^{(ss)}\|_\infty^2} \\
& \leq (C_1 C_2 c_{x0}^2 + \mu^2) C_3 p^2 \lambda \\
& \stackrel{\textcircled{1}}{\leq} \frac{\lambda}{8},
\end{aligned}$$

where  $\textcircled{1}$  holds by setting  $C_3 \leq \frac{1}{(C_1 C_2 c_{x0}^2 + \mu^2) p^2}$ . In this way, we have

$$\left\| \mathbf{G}^{hs}(0) - (\alpha_{s,3}^{(h)})^2 \bar{\mathbf{K}}^{(s)} \right\|_{\text{op}} \leq \left\| \mathbf{G}^{hs}(0) - (\alpha_{s,3}^{(h)})^2 \widehat{\mathbf{K}}^{(s)} \right\|_{\text{op}} + (\alpha_{s,3}^{(h)})^2 \|\widehat{\mathbf{K}}^{(s)} - \bar{\mathbf{K}}^{(s)}\|_{\text{op}} \leq \frac{\lambda}{4}.$$

The proof is completed.  $\square$

### C.4.3 Proof of Lemma 24

*Proof.* To begin with, according to the definition, we have

$$\begin{aligned}
\mathbf{K}_{ij}^{(ls)} - \mathbf{b}_i^{(l)} (\mathbf{b}_i^{(s)})^\top & = \sum_{t=1}^{l-1} \sum_{q=1}^{s-1} \left[ \alpha_{t,2}^{(l)} \alpha_{q,2}^{(s)} \left( \mathbf{K}_{ij}^{(tq)} - \mathbf{b}_i^{(t)} (\mathbf{b}_i^{(q)})^\top \right) \right] \\
& \quad + \tau^2 \alpha_{t,3}^{(l)} \alpha_{q,3}^{(s)} \left[ \mathbb{E}_{(M_{tq}^{(ls)}, N_{tq}^{(ls)})} \sigma(\mathbf{M}_{tq}^{(ls)}) \sigma(\mathbf{N}_{tq}^{(ls)})^\top - \mathbb{E}_{M_{tq}^{(ls)}} \sigma(\mathbf{M}_{tq}^{(ls)}) \mathbb{E}_{N_{tq}^{(ls)}} \sigma(\mathbf{N}_{tq}^{(ls)})^\top \right].
\end{aligned}$$

By defining

$$\begin{aligned}
\mathbf{R}_{tq}^{(ls)} & := \mathbb{E}_{(M_{tq}^{(ls)}, N_{tq}^{(ls)})} \begin{bmatrix} \sigma(\mathbf{M}_{tq}^{(ls)}) \sigma(\mathbf{M}_{tq}^{(ls)})^\top, \sigma(\mathbf{M}_{tq}^{(ls)}) \sigma(\mathbf{N}_{tq}^{(ls)})^\top \\ \sigma(\mathbf{N}_{tq}^{(ls)}) \sigma(\mathbf{M}_{tq}^{(ls)})^\top, \sigma(\mathbf{N}_{tq}^{(ls)}) \sigma(\mathbf{N}_{tq}^{(ls)})^\top \end{bmatrix} \\
& \quad - \mathbb{E}_{(M_{tq}^{(ls)}, N_{tq}^{(ls)})} \begin{bmatrix} \sigma(\mathbf{M}_{tq}^{(ls)}) \\ \sigma(\mathbf{N}_{tq}^{(ls)}) \end{bmatrix} \mathbb{E}_{(M_{tq}^{(ls)}, N_{tq}^{(ls)})} \left[ (\sigma(\mathbf{M}_{tq}^{(ls)})^\top, \sigma(\mathbf{N}_{tq}^{(ls)})^\top) \right],
\end{aligned}$$

we can further obtain

$$\begin{aligned} & \begin{bmatrix} \mathbf{K}_{ii}^{(ls)}, \mathbf{K}_{ij}^{(ls)} \\ \mathbf{K}_{ji}^{(ls)}, \mathbf{K}_{jj}^{(ls)} \end{bmatrix} - \begin{bmatrix} \mathbf{b}_i^{(l)} \\ \mathbf{b}_j^{(l)} \end{bmatrix} \left[ (\mathbf{b}_i^{(s)})^\top, (\mathbf{b}_j^{(s)})^\top \right] \\ &= \sum_{t=1}^{l-1} \sum_{q=1}^{s-1} \left[ \alpha_{t,2}^{(l)} \alpha_{q,2}^{(s)} \begin{bmatrix} \mathbf{K}_{ii}^{(tq)}, \mathbf{K}_{ij}^{(tq)} \\ \mathbf{K}_{ji}^{(tq)}, \mathbf{K}_{jj}^{(tq)} \end{bmatrix} - \begin{bmatrix} \mathbf{b}_i^{(t)} \\ \mathbf{b}_j^{(t)} \end{bmatrix} \left[ (\mathbf{b}_i^{(q)})^\top, (\mathbf{b}_j^{(q)})^\top \right] \right] + \tau^2 \alpha_{t,3}^{(l)} \alpha_{q,3}^{(l)} \mathbf{R}_{tq}^{(ls)}. \end{aligned}$$

Let

$$\bar{\mathbf{R}}_{tq}^{(ls)} = \begin{bmatrix} \sigma(\mathbf{M}_{tq}^{(ls)}) \\ \sigma(\mathbf{N}_{tq}^{(ls)}) \end{bmatrix} - \mathbb{E}_{(\mathbf{M}_{tq}^{(ls)}, \mathbf{N}_{tq}^{(ls)})} \begin{bmatrix} \sigma(\mathbf{M}_{tq}^{(ls)}) \\ \sigma(\mathbf{N}_{tq}^{(ls)}) \end{bmatrix}.$$

Then we have

$$\mathbf{R}_{tq}^{(ls)} = \mathbb{E}_{(\mathbf{M}_{tq}^{(ls)}, \mathbf{N}_{tq}^{(ls)})} \left[ \bar{\mathbf{R}}_{tq}^{(ls)} (\bar{\mathbf{R}}_{tq}^{(ls)})^\top \right] \succeq \mathbf{0}.$$

Therefore, by induction, we can conclude

$$\begin{aligned} & \begin{bmatrix} \mathbf{K}_{ii}^{(ls)}, \mathbf{K}_{ij}^{(ls)} \\ \mathbf{K}_{ji}^{(ls)}, \mathbf{K}_{jj}^{(ls)} \end{bmatrix} - \begin{bmatrix} \mathbf{b}_i^{(l)} \\ \mathbf{b}_j^{(l)} \end{bmatrix} \left[ (\mathbf{b}_i^{(s)})^\top, (\mathbf{b}_j^{(s)})^\top \right] \succeq a \begin{bmatrix} \mathbf{K}_{ii}^{(-1)}, \mathbf{K}_{ij}^{(-1)} \\ \mathbf{K}_{ji}^{(-1)}, \mathbf{K}_{jj}^{(-1)} \end{bmatrix} - \begin{bmatrix} \mathbf{b}_i^{(-1)} \\ \mathbf{b}_j^{(-1)} \end{bmatrix} \left[ (\mathbf{b}_i^{(-1)})^\top, (\mathbf{b}_j^{(-1)})^\top \right] \\ & \succeq a \begin{bmatrix} \mathbf{K}_{ii}^{(-1)}, \mathbf{K}_{ij}^{(-1)} \\ \mathbf{K}_{ji}^{(-1)}, \mathbf{K}_{jj}^{(-1)} \end{bmatrix} \stackrel{\textcircled{1}}{\succeq} \mathbf{0}, \end{aligned}$$

where  $a$  is a constant that depends on  $\alpha_{t,2}^{(l)}$  ( $\forall l, t$ ),  $\textcircled{1}$  holds by using Lemma 5 which shows that  $\mathbf{K}_{ii}^{(00)} \succ \mathbf{0}$ . Based on this result, we can estimate

$$\begin{aligned} & \begin{bmatrix} \mathbf{K}_{ii}^{(ll)}, \mathbf{K}_{ij}^{(ll)} \\ \mathbf{K}_{ji}^{(ll)}, \mathbf{K}_{jj}^{(ll)} \end{bmatrix} - \begin{bmatrix} \mathbf{b}_i^{(l)} \\ \mathbf{b}_j^{(l)} \end{bmatrix} \left[ (\mathbf{b}_i^{(l)})^\top, (\mathbf{b}_j^{(l)})^\top \right] \\ &= \sum_{t=1}^{l-1} \sum_{q=1}^{l-1} \left[ \alpha_{t,2}^{(l)} \alpha_{q,2}^{(s)} \begin{bmatrix} \mathbf{K}_{ii}^{(tq)}, \mathbf{K}_{ij}^{(tq)} \\ \mathbf{K}_{ji}^{(tq)}, \mathbf{K}_{jj}^{(tq)} \end{bmatrix} - \begin{bmatrix} \mathbf{b}_i^{(t)} \\ \mathbf{b}_j^{(t)} \end{bmatrix} \left[ (\mathbf{b}_i^{(q)})^\top, (\mathbf{b}_j^{(q)})^\top \right] \right] + \tau^2 \alpha_{t,3}^{(l)} \alpha_{q,3}^{(l)} \mathbf{R}_{tq}^{(ls)} \\ &\succeq \sum_{t=1}^{l-1} \left[ (\alpha_{t,2}^{(l)})^2 \begin{bmatrix} \mathbf{K}_{ii}^{(tt)}, \mathbf{K}_{ij}^{(tt)} \\ \mathbf{K}_{ji}^{(tt)}, \mathbf{K}_{jj}^{(tt)} \end{bmatrix} - \begin{bmatrix} \mathbf{b}_i^{(t)} \\ \mathbf{b}_j^{(t)} \end{bmatrix} \left[ (\mathbf{b}_i^{(t)})^\top, (\mathbf{b}_j^{(t)})^\top \right] \right] + \tau^2 (\alpha_{t,3}^{(l)})^2 \mathbf{R}_{tt}^{(ll)} \\ &\succeq \left( \prod_{t=1}^{l-1} (\alpha_{t,2}^{(l)})^2 \right) \begin{bmatrix} \mathbf{K}_{ii}^{(-1)}, \mathbf{K}_{ij}^{(-1)} \\ \mathbf{K}_{ji}^{(-1)}, \mathbf{K}_{jj}^{(-1)} \end{bmatrix} - \begin{bmatrix} \mathbf{b}_i^{(-1)} \\ \mathbf{b}_j^{(-1)} \end{bmatrix} \left[ (\mathbf{b}_i^{(-1)})^\top, (\mathbf{b}_j^{(-1)})^\top \right] \\ &\succeq \left( \prod_{t=1}^{l-1} (\alpha_{t,2}^{(l)})^2 \right) \begin{bmatrix} \mathbf{K}_{ii}^{(-1)}, \mathbf{K}_{ij}^{(-1)} \\ \mathbf{K}_{ji}^{(-1)}, \mathbf{K}_{jj}^{(-1)} \end{bmatrix}. \end{aligned}$$

Then there must exist a constant  $c$  such that

$$\lambda_{\min}(\mathbf{K}^{(ll)}) \geq \left( \prod_{t=0}^{l-1} (\alpha_{t,2}^{(l)})^2 \right) \lambda_{\min}(\widehat{\mathbf{K}}).$$

where  $\widehat{\mathbf{K}} = \begin{bmatrix} \mathbf{K}_{ii}^{(-1)}, \mathbf{K}_{ij}^{(-1)} \\ \mathbf{K}_{ji}^{(-1)}, \mathbf{K}_{jj}^{(-1)} \end{bmatrix}$ . On the other hand, we have

$$\mathbf{Q}_{ij,ab}^{(ll)} = \text{Tr} \left( \mathbf{K}_{ij, S_a^{(l)}, S_b^{(l)}}^{(ll)} \right),$$

where  $S_a^{(s)} = \{j \mid \mathbf{X}_{:,j}^{(s-1)} \in \text{the } a\text{-th patch for convolution}\}$ . This actually means that we can obtain  $\mathbf{Q}_{ij}^{(ll)}$  by using (adding) linear transformation on  $\mathbf{K}_{ij}^{(ll)}$ . Since for all  $\mathbf{Q}_{ij}^{(ll)}$  we use the same linear transformation which means that  $\mathbf{Q}^{(ll)}$  by using (adding) linear transformation on  $\mathbf{K}^{(ll)}$ . Since linear transformation does not change the eigenvalue property of a matrix, we can further obtain

$$\lambda_{\min}(\mathbf{Q}^{(ll)}) \geq \left( \prod_{t=0}^{l-1} (\alpha_{t,2}^{(l)})^2 \right) \lambda_{\min}(\widehat{\mathbf{K}}).$$

Finally, let  $\mathbf{Q} = \mathbf{B}\mathbf{S}\mathbf{B}^\top$  be the SVD of  $\mathbf{Q}$  and  $\mathbf{Z} = \mathbf{S}^{1/2}\mathbf{B}^\top$  denotes  $n$  samples (each column denotes one). Since  $\mathbf{Q}$  is full rank, the samples in  $\mathbf{Z}$  are not parallel. In this way, we can apply Lemma 5 and obtain that  $\mathbf{Q}^{(s)}$  which is defined below, is full rank

$$\mathbf{A}^{(l)} = \begin{bmatrix} \mathcal{W}_l^{(h)}(\mathbf{Q}_{ii}^{(ll)}), \mathcal{W}_l^{(h)}(\mathbf{Q}_{ij}^{(ll)}) \\ \mathcal{W}_l^{(h)}(\mathbf{Q}_{ji}^{(ll)}), \mathcal{W}_l^{(h)}(\mathbf{Q}_{jj}^{(ll)}) \end{bmatrix},$$

$$\mathbf{Q}_{ij,ab}^{(l)} = \mathbf{Q}_{ij,ab}^{(ll)} \mathbb{E}_{((\mathbf{M}, \mathbf{N}) \sim \mathbf{A}^{(l)})} \sigma'(\mathbf{M}) \sigma'(\mathbf{N})^\top, \quad \mathbf{K}_{ij,ab}^{(l)} = \text{Tr}(\mathbf{Q}_{ij}^{(s)}), \quad (s = l, \dots, h-1).$$

Recall that Lemma 10 shows

$$\frac{1}{c_{x0}} \leq \|\mathbf{X}^{(l)}(0)\|_F \leq c_{x0}.$$

where  $c_{x0} \geq 1$  is a constant. Therefore, we have  $\mathbf{K}_{ii}^{ll} = \langle \mathbf{X}^{(l)}(0), \mathbf{X}^{(l)}(0) \rangle \in [1/c_{x0}^2, c_{x0}^2]$  and thus  $\mathbf{Q}_{ii}^{ll} = \langle \Phi(\mathbf{X}^{(l)}(0)), \Phi(\mathbf{X}^{(l)}(0)) \rangle \geq \langle \mathbf{X}^{(l)}(0), \mathbf{X}^{(l)}(0) \rangle \geq 1/c_{x0}^2$  and  $\mathbf{Q}_{ii}^{ll} = \langle \Phi(\mathbf{X}^{(l)}(0)), \Phi(\mathbf{X}^{(l)}(0)) \rangle \leq k_c \langle \mathbf{X}^{(l)}(0), \mathbf{X}^{(l)}(0) \rangle \geq k_c/c_{x0}^2$ . Then we have

$$\mathbf{Q}_{ij}^{(l)} = \mathbf{Q}_{ij}^{ll} \mathbb{E}_{(\mathbf{M} \sim \mathcal{N}(0, \mathbf{I}))} \sigma'(\mathbf{M}\mathbf{Z}_i) \sigma'(\mathbf{M}\mathbf{Z}_j)^\top$$

where  $\mathbf{Z} = \mathbf{S}^{1/2}\mathbf{B}^\top$  and  $\mathbf{Z}_i = \mathbf{Z}_{:i}$  in which  $\mathbf{Q}^{ll} = \mathbf{B}\mathbf{S}\mathbf{B}^\top$  is the SVD of  $\mathbf{Q}^{ll}$ . Since  $\mathbf{Q}^{ll}$  is full rank, the samples in  $\mathbf{Z}$  are not parallel. Then we can apply Lemma 6 and obtain

$$\lambda_{\min}(\mathbf{Q}^{(l)}) \geq c_\sigma \left( \prod_{t=0}^{l-1} (\alpha_{t,2}^{(l)})^2 \right) \lambda_{\min}(\widehat{K}),$$

where  $c_\sigma$  is a constant that only depends on  $\sigma$  and input data. Since

$$\mathbf{K}_{ij,ab}^{(s)} = \text{Tr}(\mathbf{Q}_{ij}^{(s)}), \quad (s = 0, h-1)$$

which means that  $\mathbf{K}^{(s)}$  can be obtained by using adding linear transformation on  $\mathbf{Q}^{(s)}$ . So the eigenvalue of  $\mathbf{K}^{(s)}$  also satisfies

$$\lambda_{\min}(\mathbf{K}^{(l)}) \geq c_\sigma \left( \prod_{t=0}^{l-1} (\alpha_{t,2}^{(l)})^2 \right) \lambda_{\min}(\widehat{K}),$$

In this way, we can further establish

$$\begin{aligned} \lambda_{\min}(\mathbf{G}(0)) &\geq \sum_{s=0}^{h-1} \lambda_{\min}(\mathbf{G}^{hs}(0)) \stackrel{\textcircled{1}}{\geq} \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \lambda_{\min}(\mathbf{K}^{(s)}(0)) - \frac{\lambda}{4} \\ &\geq \frac{3c_\sigma}{4} \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) \lambda_{\min}(\widehat{K}), \end{aligned}$$

where  $\textcircled{1}$  holds since we set  $\lambda = c_\sigma \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) \lambda_{\min}(\widehat{K})$  and Lemma 23 shows

$$\left\| \mathbf{G}^{hs}(0) - (\alpha_{s,3}^{(h)})^2 \mathbf{K}^{(s)} \right\|_{\text{op}} \leq \frac{\lambda}{4} \quad (s = 0, \dots, h).$$

where  $\lambda$  is a constant. The proof is completed.  $\square$

## D Proofs of Results in Sec. 4

### D.1 Proof of Theorem 2

*Proof.* We first prove the first result. Suppose except one gate  $\mathbf{g}_{s,t}^{(l)}$ , all remaining stochastic gates  $\mathbf{g}_{s',t}^{(l')}$  are fixed. Then we discuss the type of the gate  $\mathbf{g}_{s,t}^{(l)}$ . Note  $\mathbf{g}_{s,t}^{(l)}$  denotes one operation in the operation set  $\mathcal{O} = \{\mathcal{O}_t\}_{t=1}^s$ , including zero operation, skip connection, pooling, and convolution with any kernel size, between nodes  $\mathbf{X}^{(s)}$  and  $\mathbf{X}^{(l)}$ . Now we discuss different kinds of operations.

If the gate  $\mathbf{g}_{s,t}^{(l)}$  is for zero operation, it is easily to check that the loss  $F_{\text{val}}(\mathbf{W}^*(\beta), \beta)$  in (2) will not change, since zero operation does not delivery any information to subsequent node  $\mathbf{X}^{(l)}$ .

If the gate  $g_{s,t}^{(l)}$  is for skip connection, there are two cases. Firstly, increasing the weight  $g_{s,t}^{(l)}$  gives smaller loss. For this case, it directly obtain our result. Secondly, increasing the weight  $g_{s,t}^{(l)}$  gives larger loss. For this case, suppose we increase  $g_{s,t}^{(l)}$  to  $g_{s,t}^{(l)} + \epsilon$ . Then node  $\mathbf{X}^{(l)}$  will become  $\mathbf{X}^{(l)} + \epsilon\mathbf{X}^{(s)} = \mathbf{X}_{\text{conv}}^{(l)} + \mathbf{X}_{\text{nonconv}}^{(l)} + \epsilon\mathbf{X}^{(s)}$  if we fix the remaining operations, where  $\mathbf{X}_{\text{conv}}^{(l)}$  denotes the output of convolution and  $\mathbf{X}_{\text{nonconv}}^{(l)}$  denotes the sum of all remaining operations. Now suppose the convolution operation between node  $\mathbf{X}^{(l)}$  and  $\mathbf{X}^{(s)}$  is  $g_{s,t}^{(l)}\text{conv}(\mathbf{W}_s^{(l)}; \mathbf{X}^{(s)}) = g_{s,t}^{(l)}\sigma(\mathbf{W}_s^{(l)}\Phi(\mathbf{X}^{(s)}))$  where  $t$  denotes the index of convolution in the operation set. Then we consider a function

$$g_{s,t}^{(l)}\sigma(\bar{\mathbf{W}}_s^{(l)}\Phi(\mathbf{X}^{(s)})) = -\epsilon\mathbf{X}^{(s)}. \quad (22)$$

Since for the almost activation functions are monotone increasing, this means that  $\sigma(\cdot)$  does not change the rank of  $\bar{\mathbf{W}}_s^{(l)}\Phi(\mathbf{X}^{(s)})$ . At the same time, the linear transformation  $\Phi(\mathbf{X}^{(s)})$  has the same rank as  $\mathbf{X}^{(s)}$ . Then when  $g_{s,t}^{(l)} \neq 0$  there exist a  $\bar{\mathbf{W}}_s^{(l)}$  such that Eqn. (22) holds. On the other hand, we already have

$$g_{s,t}^{(l)}\sigma(\mathbf{W}_s^{(l)}\Phi(\mathbf{X}^{(s)})) = \mathbf{X}_{\text{conv}}^{(l)}.$$

Since we assume the function  $\sigma(\cdot)$  is Lipschitz and smooth and the constant  $\epsilon$  is sufficient small, then by using mean value theorem, there must exist  $g_{s,t}^{(l)}\sigma(\bar{\mathbf{W}}_s^{(l)}\Phi(\mathbf{X}^{(s)})) = \mathbf{X}_{\text{conv}}^{(l)} - \epsilon\mathbf{X}^{(s)}$ . So the convolution can counteract the increment  $\epsilon\mathbf{X}^{(s)}$  brought by increasing the weight of skip connection. In this way, the whole network remains the same, leading the same loss. When the weight of convolution satisfies  $g_{s,t}^{(l)} = 0$ , we only need to increase  $g_{s,t}^{(l)}$  to a positive constant, then we use the same method and can prove the same result. In this case, we actually increase the weights of skip connection and convolution at the same time, which also accords with our results in the Proposition 2.

If the gate  $g_{s,t}^{(l)}$  is for pooling connection, we can use the same method for skip connection to prove our result, since pooling operation is also a linear transformation.

If the gate  $g_{s,t}^{(l)}$  is for convolution, then we increase it to  $g_{s,t}^{(l)} + \epsilon g_{s,t}^{(l)}$  and obtain the new output  $(1 + \epsilon)\mathbf{X}_{\text{conv}}^{(l)}$  because of  $g_{s,t}^{(l)}\sigma(\mathbf{W}_s^{(l)}\Phi(\mathbf{X}^{(s)})) = \mathbf{X}_{\text{conv}}^{(l)}$ . If the new feature map can lead to smaller loss, then we directly obtain our results. If the new feature map can lead to larger loss we only need to find a new parameter  $\bar{\mathbf{W}}_s^{(l)}$  such that  $g_{s,t}^{(l)}\sigma(\bar{\mathbf{W}}_s^{(l)}\Phi(\mathbf{X}^{(s)})) = \frac{1}{1+\epsilon}\mathbf{X}_{\text{conv}}^{(l)}$ . Since for most activation  $\sigma(0) = 0$ , we have  $g_{s,t}^{(l)}\sigma(\bar{\mathbf{W}}_s^{(l)}\Phi(\mathbf{X}^{(s)})) = 0$  when  $\bar{\mathbf{W}}_s^{(l)} = 0$ . On the other hand, we have  $g_{s,t}^{(l)}\sigma(\mathbf{W}_s^{(l)}\Phi(\mathbf{X}^{(s)})) = \mathbf{X}_{\text{conv}}^{(l)}$ . Moreover since we assume the function  $\sigma(\cdot)$  is Lipschitz and smooth and the constant  $\epsilon$  is sufficient small, then by using mean value theorem, there must exist  $\bar{\mathbf{W}}_s^{(l)}$  such that  $g_{s,t}^{(l)}\sigma(\bar{\mathbf{W}}_s^{(l)}\Phi(\mathbf{X}^{(s)})) = \frac{1}{1+\epsilon}\mathbf{X}_{\text{conv}}^{(l)}$ .

Then we prove the results in the second part. From Theorem 1, we know that for the  $k$ -th iteration in the search phase, increasing the weights  $g_{s,t_1}^{(l)}$  ( $l \neq h$ ) of skip connects and the weights  $g_{s,t_2}^{(h)}$  of convolutions can reduce the loss  $F_{\text{train}}(\mathbf{W}^*(\beta), \beta)$  in (2), where  $t_1$  and  $t_2$  respectively denote the indexes of skip connection and convolution in the operation set  $\mathcal{O} = \{O_t\}_{t=1}^s$ . Specifically, Theorem 1 proves for the training loss

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq \left(1 - \frac{\eta\lambda}{4}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|_2^2,$$

where  $\lambda = \frac{3c\sigma}{4} \lambda_{\min}(\widehat{\mathbf{K}}) \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2$ . Moreover, since  $F(\Omega) = \frac{1}{2n} \sum_{i=1}^n (u_i - y_i)^2 = \frac{1}{2n} \|\mathbf{u} - \mathbf{y}\|_2^2$ , increasing the weights  $g_{s,t_1}^{(l)}$  ( $l \neq h$ ) of skip connects and the weights  $g_{s,t_2}^{(h)}$  of convolutions can reduce the loss  $F_{\text{train}}(\mathbf{W}^*(\beta), \beta)$ . Since the samples for training and validation are drawn from the same distribution which means that  $\mathbb{E}[F_{\text{train}}(\Omega)] = \mathbb{E}[F_{\text{val}}(\Omega)]$ , increasing weights of skip connections and convolution can reduce  $F_{\text{val}}(\Omega)$  in expectation. Then by using first-order extension, we can obtain

$$\mathbb{E} \left[ F_{\text{val}}(g_{s,t_1}^{(l)} + \epsilon) - F_{\text{val}}(g_{s,t_1}^{(l)}) \right] = \epsilon \mathbb{E} \left[ \nabla_{g_{s,t_1}^{(l)}} F_{\text{val}}(g_{s,t_1}^{(l)}) \right].$$

where  $g_{s,t_1}^{(l)} \in \bar{g}_{s,t_1}^{(l)} \leq g_{s,t_1}^{(l)} + \epsilon$ . Since as above analysis, increasing the weights  $g_{s,t_1}^{(l)}$  ( $l \neq h$ ) of skip connects will reduce the current loss  $F_{\text{val}}(g_{s,t_1}^{(l)})$  in expectation, which means that  $\mathbb{E} \left[ \nabla_{g_{s,t_1}^{(l)}} F_{\text{val}}(g_{s,t_1}^{(l)}) \right]$  is positive. Since when the algorithm does not converge, we have  $0 < C \leq \mathbb{E} \left[ \nabla_{g_{s,t_1}^{(l)}} F_{\text{val}}(g_{s,t_1}^{(l)}) \right]$ . In

this way, we have

$$\mathbb{E} \left[ F_{\text{val}}(\mathbf{g}_{s,t_1}^{(l)} + \epsilon) - F_{\text{val}}(\mathbf{g}_{s,t_1}^{(l)}) \right] \geq C\epsilon.$$

Similarly, for convolution we can obtain

$$\mathbb{E} \left[ F_{\text{val}}(\mathbf{g}_{s,t_2}^{(l)} + \epsilon) - F_{\text{val}}(\mathbf{g}_{s,t_2}^{(l)}) \right] \geq C\epsilon.$$

The proof is completed.  $\square$

## D.2 Proof of Theorem 3

*Proof.* For the results in the first part, it is easily to check according to the definitions. Now we focus on proving the results in the second part. When  $\tilde{\mathbf{g}}_{s,t}^{(l)} \leq -\frac{a}{b-a}$ , then  $\mathbf{g}_{s,t}^{(l)} = 0$ . Meanwhile, the cumulative distribution of  $\tilde{\mathbf{g}}_{s,t}^{(l)}$  is  $\Theta(\tau(\ln \delta - \ln(1 - \delta)) - \beta_{s,t}^{(l)})$  [7]. In this way, we can easily compute

$$\begin{aligned} \mathbb{P}(\mathbf{g}_{s,t}^{(l)} \neq 0) &= 1 - \mathbb{P}\left(\tilde{\mathbf{g}}_{s,t}^{(l)} \leq -\frac{a}{b-a}\right) \\ &= 1 - \Theta\left(\tau\left(\ln\left(-\frac{a}{b-a}\right) - \ln\left(1 + \frac{a}{b-a}\right)\right) - \beta_{s,t}^{(l)}\right) \\ &= \Theta\left(\beta_{s,t}^{(l)} - \tau \ln \frac{-a}{b}\right). \end{aligned}$$

The proof is completed.  $\square$

## D.3 Proof of Theorem 4

*Proof.* Here we first prove the convergence rate of the shallow network with two branches. The proof is very similar to Theorem C.1. By using the totally same method, we can follow Lemma 21 to prove

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq \left(1 - \frac{\eta \lambda_{\min}(\mathbf{G}(0))}{4}\right) \|\mathbf{y} - \mathbf{u}(k-1)\|_2^2.$$

Here  $\mathbf{G}(0)$  denotes the Gram matrix of the shallow network and have the same definition as the Gram matrix of deep network with one branch. Please refer to the definition of Gram matrix in Appendix B.

The second step is to prove the smallest least eigenvalue of  $\mathbf{G}(0)$  is lower bounded. For this step, the analysis method is also the same as the method to lower bounding smallest least eigenvalue of  $\mathbf{G}(0)$  in DARTS. Specifically, by following Lemma 24, we can obtain

$$\lambda_{\min}(\mathbf{G}(0)) \geq \frac{3c_\sigma}{4} \left[ \sum_{s=1}^{\frac{h}{2}-1} (\alpha_{s,3}^{(h/2)})^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) + \sum_{s=\frac{h}{2}}^{h-1} (\alpha_{s,3}^h)^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) \right] \lambda_{\min}(\mathbf{K}).$$

where  $c_\sigma$  is a constant that only depends on  $\sigma$  and the input data,  $\lambda_{\min}(\mathbf{K}) > 0$  is given in Theorem 1.

From Theorem 1, we know that for deep cell with one branch, the loss satisfies

$$\|\mathbf{y} - \mathbf{u}(k)\|_2^2 \leq \left(1 - \frac{\eta\lambda}{4}\right)^k \|\mathbf{y} - \mathbf{u}(0)\|_2^2,$$

where  $\lambda = \frac{3c_\sigma}{4} \lambda_{\min}(\mathbf{K}) \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2$ .

Since all weights  $\alpha_{s,t}^{(l)}$  belong to the range  $[0, 1]$ , by comparison, the convergence rate  $\lambda'$  of shallow cell with two branch is large than the convergence rate  $\lambda$  of shallow cell with two branch:

$$\begin{aligned} \lambda' &= \frac{3c_\sigma}{4} \left[ \sum_{s=1}^{\frac{h}{2}-1} (\alpha_{s,3}^{(h/2)})^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) + \sum_{s=\frac{h}{2}}^{h-1} (\alpha_{s,3}^h)^2 \left( \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2 \right) \right] \lambda_{\min}(\mathbf{K}) \\ &> \lambda = \frac{3c_\sigma}{4} \lambda_{\min}(\mathbf{K}) \sum_{s=0}^{h-1} (\alpha_{s,3}^{(h)})^2 \prod_{t=0}^{s-1} (\alpha_{t,2}^{(s)})^2. \end{aligned}$$

This completes the proof.  $\square$

## E Proofs of Auxiliary Lemmas

### E.1 Proof of Lemma 8

*Proof.* We use chain rule to obtain the following gradients:

$$\begin{aligned}
\frac{\partial \ell}{\partial \mathbf{X}^{(h-1)}} &= (u - y) \mathbf{W}_h \in \mathbb{R}^{m \times p}; \\
\frac{\partial \ell}{\partial \mathbf{X}^{(l)}} &= (u - y) \mathbf{W}_l + \sum_{s=l+1}^h \frac{\partial \ell}{\partial \mathbf{X}^{(s)}} \frac{\partial \mathbf{X}^{(s)}}{\partial \mathbf{X}^{(l)}} \quad (l = 0, \dots, h-2) \\
&= (u - y) \mathbf{W}_l + \sum_{s=l+1}^h \left( \alpha_{l,2}^{(s)} \frac{\partial \ell}{\partial \mathbf{X}^{(s)}} + \alpha_{l,3}^{(s)} \tau \Psi \left( (\mathbf{W}_l^{(s)})^\top \left( \sigma' \left( \mathbf{W}_l^{(s)} \Phi(\mathbf{X}^{(l)}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(s)}} \right) \right) \right) \in \mathbb{R}^{m \times p}; \\
\frac{\partial \ell}{\partial \mathbf{X}} &= \frac{\partial \ell}{\partial \mathbf{X}^{(1)}} \frac{\partial \mathbf{X}^{(1)}}{\partial \mathbf{X}^{(0)}} = \tau \Psi \left( (\mathbf{W}^{(0)})^\top \left( \sigma' \left( \mathbf{W}^{(0)} \Phi(\mathbf{X}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(0)}} \right) \right) \in \mathbb{R}^{m \times p}, \\
\frac{\partial \ell}{\partial \mathbf{W}_s^{(l)}} &= \frac{\partial \ell}{\partial \mathbf{X}^{(l)}} \frac{\partial \mathbf{X}^{(l)}}{\partial \mathbf{W}_s^{(l)}} = \alpha_{s,3}^{(l)} \tau \Phi(\mathbf{X}^{(s)}) \left( \sigma' \left( \mathbf{W}_s^{(l)} \Phi(\mathbf{X}^{(s)}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(l)}} \right)^\top \in \mathbb{R}^{m \times p} \\
&\quad (1 \leq l \leq h, 1 \leq s \leq l-1); \\
\frac{\partial \ell}{\partial \mathbf{W}^{(0)}} &= \frac{\partial \ell}{\partial \mathbf{X}^{(0)}} \frac{\partial \mathbf{X}^{(0)}}{\partial \mathbf{W}^{(0)}} = \tau \Phi(\mathbf{X}) \left( \sigma' \left( \mathbf{W}^{(0)} \Phi(\mathbf{X}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(0)}} \right)^\top \in \mathbb{R}^{m \times p}, \\
\frac{\partial \ell}{\partial \mathbf{W}_s} &= (u - y) \mathbf{X}^{(l)} \in \mathbb{R}^{m \times p},
\end{aligned}$$

where  $\odot$  denotes the dot product. □

### E.2 Proof of Lemma 9

*Proof.* We use chain rule to obtain the following gradients:

$$\begin{aligned}
\frac{\partial u}{\partial \mathbf{X}^{(h-1)}} &= \mathbf{W}_{h-1} \in \mathbb{R}^{m \times p}; \\
\frac{\partial u}{\partial \mathbf{X}^{(l)}} &= \mathbf{W}_l + \sum_{s=l+1}^h \frac{\partial u}{\partial \mathbf{X}^{(s)}} \frac{\partial \mathbf{X}^{(s)}}{\partial \mathbf{X}^{(l)}} \quad (l = 0, \dots, h-2) \\
&= \mathbf{W}_l + \sum_{s=l+1}^h \left( \alpha_{l,2}^{(s)} \frac{\partial u}{\partial \mathbf{X}^{(s)}} + \alpha_{l,3}^{(s)} \tau \Psi \left( (\mathbf{W}_l^{(s)})^\top \left( \sigma' \left( \mathbf{W}_l^{(s)} \Phi(\mathbf{X}^{(l)}) \right) \odot \frac{\partial u}{\partial \mathbf{X}^{(s)}} \right) \right) \right) \in \mathbb{R}^{m \times p}; \\
&\quad (0 \leq l \leq h-1, 0 \leq s \leq l-1), \\
\frac{\partial u}{\partial \mathbf{X}} &= \frac{\partial u}{\partial \mathbf{X}^{(1)}} \frac{\partial \mathbf{X}^{(1)}}{\partial \mathbf{X}^{(0)}} = \tau \Psi \left( (\mathbf{W}^{(0)})^\top \left( \sigma' \left( \mathbf{W}^{(0)} \Phi(\mathbf{X}) \right) \odot \frac{\partial u}{\partial \mathbf{X}^{(0)}} \right) \right) \in \mathbb{R}^{m \times p}, \\
\frac{\partial u}{\partial \mathbf{W}_s^{(l)}} &= \frac{\partial u}{\partial \mathbf{X}^{(l)}} \frac{\partial \mathbf{X}^{(l)}}{\partial \mathbf{W}_s^{(l)}} = \alpha_{s,3}^{(l)} \tau \Phi(\mathbf{X}^{(s)}) \left( \sigma' \left( \mathbf{W}_s^{(l)} \Phi(\mathbf{X}^{(s)}) \right) \odot \frac{\partial u}{\partial \mathbf{X}^{(l)}} \right)^\top \in \mathbb{R}^{m \times p} \\
&\quad (0 \leq l \leq h-1, 1 \leq s \leq l-1); \\
\frac{\partial u}{\partial \mathbf{W}^{(0)}} &= \frac{\partial u}{\partial \mathbf{X}^{(0)}} \frac{\partial \mathbf{X}^{(0)}}{\partial \mathbf{W}^{(0)}} = \tau \Phi(\mathbf{X}) \left( \sigma' \left( \mathbf{W}^{(0)} \Phi(\mathbf{X}) \right) \odot \frac{\partial u}{\partial \mathbf{X}^{(0)}} \right)^\top \in \mathbb{R}^{m \times p},
\end{aligned}$$

where  $\odot$  denotes the dot product. □

### E.3 Proof of Lemma 10

*Proof.* We each layer in turn. Our proof follows the proof framework in [4]. Note for notation simplicity, we have assumed that the input  $\mathbf{X}$  is of size  $m \times p$  in Sec. B. To begin with, we look at the

first layer. For brevity, let  $\mathbf{H} = \Phi(\mathbf{X})$ . According to the definition, we have

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{X}^{(0)}(0)\|_F^2 \right] &= \tau^2 \mathbb{E} \left[ \|\sigma(\mathbf{W}^{(0)}(0)\Phi(\mathbf{X}))\|_F^2 \right] = \tau^2 \sum_{i=1}^m \sum_{j=1}^p \mathbb{E} \left[ \sigma^2(\mathbf{W}_{i:}^{(0)}(0)\mathbf{H}_{:j}) \right] \\ &\stackrel{\textcircled{1}}{=} \sum_{j=1}^p \mathbb{E}_{\omega \sim \mathcal{N}(0,1)} \left[ \sigma^2(\|\mathbf{H}_{:j}\|_F \omega) \right] \stackrel{\textcircled{2}}{\geq} \mathbb{E}_{\omega \sim \mathcal{N}(0,1)} \left[ \sigma^2(\|\mathbf{H}_{:j'}\|_F \omega) \right] \\ &\geq \mathbb{E}_{\omega \sim \mathcal{N}(0, \frac{1}{\sqrt{p}})} \left[ \sigma^2(\omega) \right] := c > 0, \end{aligned}$$

where  $\textcircled{1}$  holds since  $\tau = 1/\sqrt{m}$  and the entries in  $\mathbf{W}^{(0)}(0)$  obeys i.i.d. Gaussian distribution which gives  $\sum_{i=1}^n a_i \omega_i \sim \mathcal{N}(0, \sum_{i=1}^n a_i^2)$  with  $\omega_i \sim \mathcal{N}(0, 1)$ ;  $\textcircled{2}$  holds since  $\|\mathbf{X}\| = 1$  which means there must exist one  $j'$  such that  $\|\mathbf{H}_{:j'}\|_F \geq \frac{1}{\sqrt{p}}$ .

Next, we can bound the variance

$$\begin{aligned} &\text{Var} \left[ \|\mathbf{X}^{(0)}(0)\|_F^2 \right] \\ &= \tau^4 \text{Var} \left[ \|\sigma(\mathbf{W}^{(0)}(0)\Phi(\mathbf{X}))\|_F^2 \right] = \tau^4 \text{Var} \left[ \sum_{i=1}^m \sum_{j=1}^p \mathbb{E} \left[ \sigma^2(\mathbf{W}_{i:}^{(0)}(0)\mathbf{H}_{:j}) \right] \right] \\ &\stackrel{\textcircled{1}}{=} \tau^2 \text{Var} \left[ \sum_{j=1}^p \mathbb{E} \left[ \sigma^2(\mathbf{W}_{i:}^{(0)}(0)\mathbf{H}_{:j}) \right] \right] \stackrel{\textcircled{2}}{\leq} \tau^2 \mathbb{E}_{\omega \sim \mathcal{N}(0,1)} \left[ \left( \sum_{j=1}^p (\sigma(0) + \|\mathbf{H}_{:j}\|_F |\omega|)^2 \right)^2 \right] \\ &\leq \frac{p^2}{m} c_1, \end{aligned}$$

where  $\textcircled{1}$  holds since  $\tau = 1/\sqrt{m}$  and the entries in  $\mathbf{W}^{(0)}(0)$  obeys i.i.d. Gaussian distribution,  $\textcircled{2}$  holds since  $\text{Var}(x) \leq \mathbb{E}[x^2] - [\mathbb{E}(x)]^2$ ,  $\textcircled{3}$  holds since  $\|\mathbf{H}_{:j}\| \leq 1$  and  $c_1 = \sigma^4(0) + 4|\sigma^3(0)|\mu\sqrt{2/\pi} + 8|\sigma(0)|\mu^3\sqrt{2/\pi} + 32\mu^4$ . Then by using Chebyshev's inequality in Lemma 1, we have

$$\mathbb{P} \left( \left| \|\mathbf{X}^{(0)}(0)\|_F^2 - \mathbb{E}[\|\mathbf{X}^{(0)}(0)\|_F^2] \right| \geq \frac{c}{2} \right) \leq \frac{4\text{Var}(\|\mathbf{X}^{(0)}(0)\|_F^2)}{c^2} \leq \frac{4p^2}{mc^2} c_1.$$

By setting  $m \geq \frac{4c_1 p^2}{c^2 \delta}$ , we have with probability at least  $1 - \frac{\delta}{n}$ ,

$$\|\mathbf{X}^{(0)}(0)\|_F^2 \geq \frac{c}{2}.$$

Meanwhile, we can upper bound  $\|\mathbf{X}^{(0)}(0)\|_F^2$  as follows:

$$\|\mathbf{X}^{(0)}(0)\|_F^2 \leq \tau^2 \|\sigma(\mathbf{W}^{(0)}(0)\Phi(\mathbf{X}))\|_F^2 \leq \tau^2 \mu^2 \|\mathbf{W}^{(0)}(0)\Phi(\mathbf{X})\|_F^2 \stackrel{\textcircled{1}}{\leq} \mu^2 c_{w0}^2 \|\Phi(\mathbf{X})\|_F^2 \stackrel{\textcircled{2}}{\leq} k_c \mu^2 c_{w0}^2,$$

where  $\textcircled{1}$  holds since  $\|\mathbf{W}_s^{(l)}(0)\|_2 \leq \sqrt{m} c_{w0}$ , and  $\textcircled{2}$  uses  $\|\Phi(\mathbf{X})\|_F^2 \leq k_c \|\mathbf{X}\|_F^2$ .

Next we consider the cases where  $l \geq 1$ . According to the definition, we can obtain

$$\begin{aligned} \|\mathbf{X}^{(l)}(0)\|_F &= \left\| \sum_{s=0}^{l-1} \left( \alpha_{s,2}^{(l)} \mathbf{X}^{(s)}(0) + \alpha_{s,3}^{(l)} \tau \sigma(\mathbf{W}_s^{(l)}(0)\Phi(\mathbf{X}^{(s)}(0))) \right) \right\|_F \\ &\leq \sum_{s=0}^{l-1} \left( \alpha_{s,2}^{(l)} \|\mathbf{X}^{(s)}(0)\|_F + \alpha_{s,3}^{(l)} \tau \|\sigma(\mathbf{W}_s^{(l)}(0)\Phi(\mathbf{X}^{(s)}(0)))\|_F \right) \\ &\stackrel{\textcircled{1}}{\leq} \left( \alpha_{s,2}^{(l)} + \alpha_{s,3}^{(l)} \sqrt{k_c} \mu c_{w0} \right) \sum_{s=0}^{l-1} \|\mathbf{X}^{(s)}(0)\|_F \\ &\stackrel{\textcircled{2}}{\leq} \frac{c_2^{l+1} - 1}{c_2 - 1} c_2 \sqrt{k_c} \mu c_{w0}, \end{aligned}$$

where  $\textcircled{1}$  uses the fact that  $\|\sigma(\mathbf{W}_s^{(l)}(0)\Phi(\mathbf{X}^{(s)}(0)))\|_F \leq \mu \|\mathbf{W}_s^{(l)}(0)\Phi(\mathbf{X}^{(s)}(0))\|_F \leq \sqrt{m} \mu c_{w0} \|\Phi(\mathbf{X}^{(s)}(0))\|_F \leq \sqrt{m} \mu \sqrt{k_c} c_{w0} \|\mathbf{X}^{(s)}(0)\|_F$ ,  $\textcircled{2}$  holds by setting  $c_2 = \alpha_{s,2}^{(l)} + \alpha_{s,3}^{(l)} \sqrt{k_c} \mu c_{w0}$ .



Similarly, we can obtain

$$\begin{aligned}
\|\mathbf{X}^{(l)}(0)\|_F &= \left\| \sum_{s=0}^{l-1} \left( \alpha_{s,2}^{(l)} \mathbf{X}^{(s)}(0) + \alpha_{s,3}^{(l)} \tau \sigma(\mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0))) \right) \right\|_F \\
&\geq \min_{0 \leq s \leq l-1} \left| \alpha_{s,2}^{(l)} \|\mathbf{X}^{(s)}(0)\|_F - \alpha_{s,3}^{(l)} \tau \|\sigma(\mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)))\|_F \right| \\
&\geq \min_{0 \leq s \leq l-1} \left| \alpha_{s,2}^{(l)} - \alpha_{s,3}^{(l)} \sqrt{k_c} \mu c_{w0} \right| \|\mathbf{X}^{(s)}(0)\|_F \\
&\geq \left| \alpha_{s,2}^{(l)} - \alpha_{s,3}^{(l)} \sqrt{k_c} \mu c_{w0} \right|^{l-1} \sqrt{k_c} \mu c_{w0} > 0.
\end{aligned}$$

Therefore, we can obtain that there exists a constant  $c_{x0}$  such that for all  $l \in [0, 1, \dots, h-1]$ ,

$$\frac{1}{c_{x0}} \leq \|\mathbf{X}^{(l)}(0)\|_F \leq c_{x0}.$$

The proof is completed.  $\square$

#### E.4 Proof of Lemma 11

*Proof.* For this proof, we will respectively bound each layer. We first consider the first layer, namely  $l = 1$ .

**Step 1. Case where  $l = 0$ : upper bound of  $\|\mathbf{X}^{(0)}(k) - \mathbf{X}^{(0)}(0)\|_F$ .** According to the definition, we have  $\mathbf{X}^{(0)}(k) = \tau \sigma(\mathbf{W}^{(0)}(k) \Phi(\mathbf{X}))$  which yields

$$\begin{aligned}
\|\mathbf{X}^{(0)}(k) - \mathbf{X}^{(0)}(0)\|_F &= \tau \|\sigma(\mathbf{W}^{(0)}(k) \Phi(\mathbf{X})) - \sigma(\mathbf{W}^{(0)}(0) \Phi(\mathbf{X}))\|_F \\
&\stackrel{\textcircled{1}}{\leq} \tau \mu \|\mathbf{W}^{(0)}(k) \Phi(\mathbf{X}) - \mathbf{W}^{(0)}(0) \Phi(\mathbf{X})\|_F \\
&\stackrel{\textcircled{2}}{\leq} \tau \mu \sqrt{k_c} \|\mathbf{W}^{(0)}(k) - \mathbf{W}^{(0)}(0)\|_F \\
&\stackrel{\textcircled{3}}{\leq} \mu \sqrt{k_c} r,
\end{aligned}$$

where  $\textcircled{1}$  uses the  $\mu$ -Lipschitz of  $\sigma(\cdot)$ ,  $\textcircled{2}$  uses  $\|\Phi(\mathbf{X})\| \leq \sqrt{k_c} \|\mathbf{X}\| \leq \sqrt{k_c}$ ,  $\textcircled{3}$  uses the assumption  $\|\mathbf{W}^{(0)}(k) - \mathbf{W}^{(0)}(0)\|_2 \leq \sqrt{mr}$ .

**Step 2. Case where  $l \geq 1$ : upper bound of  $\|\mathbf{X}^{(l)}(k) - \mathbf{X}^{(l)}(0)\|_F$ .** According to the definition, we have

$$\begin{aligned}
&\|\mathbf{X}^{(l)}(k) - \mathbf{X}^{(l)}(0)\|_F \\
&= \left\| \sum_{s=0}^{l-1} \left[ \alpha_{s,2}^{(l)} \left( \mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0) \right) + \alpha_{s,3}^{(l)} \tau \left( \sigma(\mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k))) - \sigma(\mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0))) \right) \right] \right\|_F \\
&= \sum_{s=0}^{l-1} \left[ \alpha_{s,2}^{(l)} \|\mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0)\|_F + \alpha_{s,3}^{(l)} \tau \left\| \sigma(\mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k))) - \sigma(\mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0))) \right\|_F \right] \\
&\leq \sum_{s=0}^{l-1} \left[ \alpha_{s,2}^{(l)} \|\mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0)\|_F + \alpha_{s,3}^{(l)} \tau \mu \left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k)) - \mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F \right]
\end{aligned}$$

Then we first bound the second term as follows:

$$\begin{aligned}
&\left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k)) - \mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F \\
&\leq \left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k)) - \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F + \left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(0)) - \mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F \\
&\leq \|\mathbf{W}_s^{(l)}(k)\| \left\| \Phi(\mathbf{X}^{(s)}(k)) - \Phi(\mathbf{X}^{(s)}(0)) \right\|_F + \left\| \mathbf{W}_s^{(l)}(k) - \mathbf{W}_s^{(l)}(0) \right\|_F \|\Phi(\mathbf{X}^{(s)}(0))\|_F \\
&\leq \sqrt{k_c} \|\mathbf{W}_s^{(l)}(k)\| \left\| \mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0) \right\|_F + \sqrt{k_c} \left\| \mathbf{W}_s^{(l)}(k) - \mathbf{W}_s^{(l)}(0) \right\|_F \|\mathbf{X}^{(s)}(0)\|_F \\
&\stackrel{\textcircled{1}}{\leq} \sqrt{k_c} \sqrt{m} (r + c_{w0}) \left\| \mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0) \right\|_F + \sqrt{k_c} m c_{x0} \tilde{r},
\end{aligned}$$

where in  $\textcircled{1}$  we use  $\|\mathbf{W}_s^{(l)}(k)\|_F \leq \|\mathbf{W}_s^{(l)}(k) - \mathbf{W}_s^{(l)}(0)\|_F + \|\mathbf{W}_s^{(l)}(0)\|_F \leq \sqrt{m}(r + c_{w0})$ ,  $\left\| \mathbf{W}_s^{(l)}(k) - \mathbf{W}_s^{(l)}(0) \right\|_F \leq \sqrt{m} \tilde{r}$ , and the results in Lemma 10 that  $\frac{1}{c_{x0}} \leq \|\mathbf{X}^{(l)}(0)\|_F \leq c_{x0}$ . Plugging

this result into the above inequality gives

$$\begin{aligned}
& \|\mathbf{X}^{(l)}(k) - \mathbf{X}^{(l)}(0)\|_F \\
& \leq \sum_{s=0}^{l-1} \left[ \alpha_{s,2}^{(l)} \|\mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0)\|_F + \alpha_{s,3}^{(l)} \tau \mu \left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k)) - \mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F \right] \\
& \leq \sum_{s=0}^{l-1} \left[ \left( \alpha_{s,2}^{(l)} + \alpha_{s,3}^{(l)} \mu \sqrt{k_c} (r + c_{w0}) \right) \|\mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0)\|_F + \alpha_{s,3}^{(l)} \mu \sqrt{k_c} c_{x0} \tilde{r} \right] \\
& \leq \sum_{s=0}^{l-1} \left[ \left( \alpha_{s,2}^{(l)} + \alpha_{s,3}^{(l)} \mu \sqrt{k_c} (r + c_{w0}) \right) \|\mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0)\|_F + \alpha_{s,3}^{(l)} \mu \sqrt{k_c} c_{x0} \tilde{r} \right] \tag{23} \\
& \leq \sum_{s=0}^{l-1} \left[ \left( \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right) \|\mathbf{X}^{(s)}(k) - \mathbf{X}^{(s)}(0)\|_F + \alpha_{s,3}^{(l)} \mu \sqrt{k_c} c_{x0} \tilde{r} \right] \\
& \leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right) \|\mathbf{X}^{(l-1)}(k) - \mathbf{X}^{(l-1)}(0)\|_F \\
& \leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right)^l \|\mathbf{X}^{(0)}(k) - \mathbf{X}^{(0)}(0)\|_F \\
& \leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right)^l \mu \sqrt{k_c} r,
\end{aligned}$$

where  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ .

By using Eqn. (23), we have

$$\left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k)) - \mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F \leq \frac{1}{\alpha_3} \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0}) \right)^l \sqrt{k_c} m r,$$

The proof is completed.  $\square$

## E.5 Proof of Lemma 12

*Proof.* According to definition, we have

$$\frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(h)}(t)} \right\|_F = \frac{1}{n} \sum_{i=1}^n \|(u_i(t) - y_i) \mathbf{W}_h(t)\|_F \stackrel{\textcircled{1}}{\leq} \frac{1}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_F \|\mathbf{W}_l(t)\|_F \stackrel{\textcircled{2}}{\leq} c_y c_u, \tag{24}$$

where  $\textcircled{1}$  holds since  $\sum_{i=1}^n |u_i - y_i| \leq \sqrt{n} \|\mathbf{u} - \mathbf{y}\|_2 = \sqrt{n} \sqrt{\sum_i (u_i - y_i)^2}$ ,  $\textcircled{2}$  holds by assuming  $\frac{1}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_F = c_y$  and  $\|\mathbf{W}_h(t)\|_F \leq c_u$ .

Then for  $0 \leq l < h$ , we have

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)} \right\|_F = \frac{1}{n} \sum_{i=1}^n \|(u_i(t) - y_i) \mathbf{W}_l(t) \\
& + \sum_{s=l+1}^{h-1} \left( \alpha_{l,2}^{(s)} \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} + \alpha_{l,3}^{(s)} \tau \Psi \left( (\mathbf{W}_l^{(s)}(t))^\top \left( \sigma' \left( \mathbf{W}_l^{(s)}(t) \Phi(\mathbf{X}_i^{(l)}(t)) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right) \right) \right) \Big\|_F \\
& \leq \frac{1}{n} \sum_{i=1}^n \|(u_i(t) - y_i) \mathbf{W}_l(t)\|_F \\
& + \sum_{s=l+1}^{h-1} \frac{1}{n} \sum_{i=1}^n \left\| \alpha_{l,2}^{(s)} \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} + \alpha_{l,3}^{(s)} \tau \Psi \left( (\mathbf{W}_l^{(s)}(t))^\top \left( \sigma' \left( \mathbf{W}_l^{(s)}(t) \Phi(\mathbf{X}_i^{(l)}(t)) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right) \right) \right\|_F
\end{aligned}$$

The main task is to bound

$$\begin{aligned}
& \left\| \alpha_{i,2}^{(s)} \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} + \alpha_{i,3}^{(s)} \tau \Psi \left( (\mathbf{W}_l^{(s)}(t))^\top \left( \sigma' \left( \mathbf{W}_l^{(s)}(t) \Phi(\mathbf{X}_i^{(l)}(t)) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right) \right) \right\|_F \\
& \leq \alpha_{i,2}^{(s)} \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right\|_F + \alpha_{i,3}^{(s)} \tau \left\| \Psi \left( (\mathbf{W}_l^{(s)}(t))^\top \left( \sigma' \left( \mathbf{W}_l^{(s)}(t) \Phi(\mathbf{X}_i^{(l)}(t)) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right) \right) \right\|_F \\
& \stackrel{\textcircled{1}}{\leq} \alpha_{i,2}^{(s)} \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right\|_F + \alpha_{i,3}^{(s)} \tau \mu \sqrt{k_c} \|\mathbf{W}_l^{(s)}(t)\|_F \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right\|_F \\
& \stackrel{\textcircled{1}}{\leq} \left( \alpha_{i,2}^{(s)} + \alpha_{i,3}^{(s)} \mu \sqrt{k_c} (c_{w0} + r) \right) \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right\|_F,
\end{aligned}$$

where  $\textcircled{1}$  holds since  $\|\Psi(\mathbf{X})\|_F \leq \sqrt{k_c} \|\mathbf{X}\|_F$  and the activation function  $\sigma(\cdot)$  is  $\mu$ -Lipschitz,  $\textcircled{2}$  holds since  $\|\mathbf{W}_l^{(s)}(t)\|_F \leq \|\mathbf{W}_l^{(s)}(t) - \mathbf{W}_l^{(s)}(0)\|_F + \|\mathbf{W}_l^{(s)}(0)\|_F \leq \sqrt{m} (c_{w0} + r)$ . Similar to (24), we can prove

$$\frac{1}{n} \sum_{i=1}^n \|(u_i(t) - y_i) \mathbf{W}_l(t)\|_F \leq \frac{1}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_F \|\mathbf{W}_l(t)\|_F \leq c_y c_u,$$

Combining the above results yields

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)} \right\|_F & \leq c_y c_u + \sum_{s=l+1}^{h-1} \left( \alpha_{i,2}^{(s)} + \alpha_{i,3}^{(s)} \mu \sqrt{k_c} (c_{w0} + r) \right) \frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right\|_F \\
& \stackrel{\textcircled{1}}{\leq} c_y c_u + \sum_{s=l+1}^{h-1} \left( \alpha_2 + \alpha_3 \mu \sqrt{k_c} (c_{w0} + r) \right) \frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(t)} \right\|_F \\
& \leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (c_{w0} + r) \right) \frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l-1)}(t)} \right\|_F \\
& \leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (c_{w0} + r) \right)^l \frac{1}{n} \sum_{i=1}^n \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(0)}(t)} \right\|_F \\
& \leq \left( 1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (c_{w0} + r) \right)^l c_y c_u,
\end{aligned}$$

where  $\textcircled{1}$  uses  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . The proof is completed.  $\square$

## E.6 Proof of Lemma 13

*Proof.* Here we use mathematical induction to prove these results in turn. We first consider  $t = 0$ . The following results hold:

$$\|\mathbf{W}_s^{(l)}(t) - \mathbf{W}_s^{(l)}(0)\|_F \leq \sqrt{m} \tilde{r}, \quad \|\mathbf{W}_s(t) - \mathbf{W}_s(0)\|_F \leq \sqrt{m} \tilde{r}. \quad (25)$$

Now we assume (25) holds for  $t = 1, \dots, k$ . We only need to prove it hold for  $t + 1$ . According to the definitions, we can establish

$$\begin{aligned}
\|\mathbf{W}_s^{(l)}(t+1) - \mathbf{W}_s^{(l)}(t)\|_F & = \eta \alpha_{s,3}^{(l)} \tau \left\| \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{X}_i^{(s)}(t)) \left( \sigma' \left( \mathbf{W}_s^{(l)}(t) \Phi(\mathbf{X}_i^{(s)}(t)) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)} \right) \right\|_F^\top \\
& \leq \eta \alpha_{s,3}^{(l)} \tau \frac{1}{n} \sum_{i=1}^n \left\| \Phi(\mathbf{X}_i^{(s)}(t)) \left( \sigma' \left( \mathbf{W}_s^{(l)}(t) \Phi(\mathbf{X}_i^{(s)}(t)) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)} \right) \right\|_F^\top \\
& \stackrel{\textcircled{1}}{\leq} \eta \alpha_{s,3}^{(l)} \tau \sqrt{k_c} \frac{1}{n} \sum_{i=1}^n \|\mathbf{X}_i^{(s)}(t)\| \left\| \sigma' \left( \mathbf{W}_s^{(l)}(t) \Phi(\mathbf{X}_i^{(s)}(t)) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)} \right\|_F \\
& \stackrel{\textcircled{2}}{\leq} 2\eta \alpha_{s,3}^{(l)} \tau \sqrt{k_c} c_{x0} \frac{1}{n} \sum_{i=1}^n \left\| \sigma' \left( \mathbf{W}_s^{(l)}(t) \Phi(\mathbf{X}_i^{(s)}(t)) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)} \right\|_F
\end{aligned}$$

where ① holds since  $\|\Phi(\mathbf{X}^{(s)})\|_F \leq \sqrt{k_c}\|\mathbf{X}^{(s)}\|_F$ ; ② holds since in Lemma 11 and Lemma 10, we have

$$\begin{aligned}\|\mathbf{X}^{(l)}(t)\| &\leq \|\mathbf{X}^{(l)}(t) - \mathbf{X}^{(l)}(0)\|_F + \|\mathbf{X}^{(l)}(0)\|_F \\ &\leq c_{x0} + \left(1 + \alpha_2 + \alpha_3\mu\sqrt{k_c}(r + c_{w0})\right)^l \mu\sqrt{k_c}r \\ &\stackrel{\textcircled{1}}{\leq} 2c_{x0},\end{aligned}\tag{26}$$

where  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ , and  $c_{x0} \geq 1$  is given in Lemma 10. The inequality holds by setting  $r$  small enough, namely  $r \leq \min\left(\frac{c_{x0}}{(1+\alpha_2+2\alpha_3\mu\sqrt{k_c}c_{w0})^l \mu\sqrt{k_c}}, c_{w0}\right)$ . This condition will be satisfied by setting enough large  $m$  and will be discussed later.

Since the activation function  $\sigma(\cdot)$  is  $\mu$ -Lipschitz, we have

$$\left\|\sigma'(\mathbf{W}_s^{(l)}(t)\Phi(\mathbf{X}^{(s)}(t))) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(l)}(t)}\right\|_F \leq \mu \left\|\frac{\partial \ell}{\partial \mathbf{X}^{(l)}(t)}\right\|_F.$$

So the remaining task is to upper bound  $\left\|\frac{\partial \ell}{\partial \mathbf{X}^{(l)}(t)}\right\|_F$ . Towards this goal, we have  $\frac{1}{\sqrt{n}}\|\mathbf{u}(t) - \mathbf{y}\|_F \leq c_y = \frac{1}{\sqrt{n}}(1 - \frac{\eta\lambda}{2})^{t/2}\|\mathbf{y} - \mathbf{u}(0)\|_2$ ,  $\|\mathbf{W}_h(t)\|_F \leq \|\mathbf{W}_h(t) - \mathbf{W}_h(0)\|_F + \|\mathbf{W}_h(0)\|_F \leq c_u = \sqrt{m}(\tilde{r} + c_{w0})$ ,  $\|\mathbf{W}_l^{(s)}(t) - \mathbf{W}_l^{(s)}(0)\|_F \leq \sqrt{m}r$ , and  $\|\mathbf{W}_l^{(s)}(0)\|_F \leq c_{w0}$ . In this way, we can use Lemma Lemma 12 and obtain

$$\frac{1}{n} \sum_{i=1}^n \left\|\frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(t)}\right\|_F \leq c_1 c_y c_u = \frac{c_1(\tilde{r} + c_{w0})}{\sqrt{n}} \left(1 - \frac{\eta\lambda}{2}\right)^{t/2} \|\mathbf{y} - \mathbf{u}(0)\|_2,$$

where  $c_1 = (1 + \alpha_2 + \alpha_3\tau\mu\sqrt{k_c}(\tilde{r} + c_{w0}))^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ .

By combining the above results, we can directly obtain

$$\begin{aligned}\|\mathbf{W}_s^{(l)}(t+1) - \mathbf{W}_s^{(l)}(t)\|_F &\leq \frac{2c_1\eta\alpha_{s,3}^{(l)}\mu\sqrt{k_c}c_{x0}(\tilde{r} + c_{w0})}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_F \\ &\leq \frac{2c_1\eta\alpha_{s,3}^{(l)}\mu\sqrt{k_c}c_{x0}(\tilde{r} + c_{w0})}{\sqrt{n}} \left(1 - \frac{\eta\lambda}{2}\right)^{t/2} \|\mathbf{y} - \mathbf{u}(0)\|_2.\end{aligned}$$

Therefore, we have

$$\begin{aligned}\|\mathbf{W}_s^{(l)}(t+1) - \mathbf{W}_s^{(l)}(0)\|_F &\leq \|\mathbf{W}_s^{(l)}(t+1) - \mathbf{W}_s^{(l)}(t)\|_F + \|\mathbf{W}_s^{(l)}(t) - \mathbf{W}_s^{(l)}(0)\|_F \\ &\leq \frac{8c_1\alpha_{s,3}^{(l)}\mu\sqrt{k_c}c_{x0}(\tilde{r} + c_{w0})}{\lambda\sqrt{n}} \|\mathbf{y} - \mathbf{u}(0)\|_2 \stackrel{\textcircled{1}}{\leq} \sqrt{m}\tilde{r},\end{aligned}$$

where ① holds by setting  $\tilde{r} = \frac{16(1+\alpha_2+2\alpha_3\mu\sqrt{k_c}c_{w0})^l \alpha_{s,3}^{(l)}\mu\sqrt{k_c}c_{x0}c_{w0}}{\lambda\sqrt{m}n} \|\mathbf{y} - \mathbf{u}(0)\|_2 \leq c_{w0}$ . By using the same way, we can prove

$$\begin{aligned}\|\mathbf{W}^{(0)}(t+1) - \mathbf{W}^{(0)}(t)\|_F &\leq \frac{2c_1\eta\mu\sqrt{k_c}c_{x0}(\tilde{r} + c_{w0})}{\sqrt{n}} \left(1 - \frac{\eta\lambda}{2}\right)^{t/2} \|\mathbf{y} - \mathbf{u}(0)\|_2, \\ \|\mathbf{W}_s^{(l)}(t+1) - \mathbf{W}_s^{(l)}(0)\|_F &\leq \sqrt{m}\tilde{r}.\end{aligned}$$

Then similarly, we can obtain

$$\begin{aligned}\|\mathbf{W}_s(t+1) - \mathbf{W}_s(t)\|_F &= \eta \left\|\frac{1}{n} \sum_{i=1}^n (u_i - y_i) \mathbf{X}_i^{(s)}(t)\right\|_F \leq \eta \frac{1}{n} \sum_{i=1}^n |u_i(t) - y_i| \|\mathbf{X}_i^{(s)}(t)\|_F \\ &\stackrel{\textcircled{1}}{\leq} \frac{2\eta c_{x0}}{\sqrt{n}} \|\mathbf{u}(t) - \mathbf{y}\|_2 \leq \frac{2\eta c_{x0}}{\sqrt{n}} \left(1 - \frac{\eta\lambda}{2}\right)^{t/2} \|\mathbf{y} - \mathbf{u}(0)\|_2,\end{aligned}$$

where ① holds since  $\sum_{i=1}^n |u_i - y_i| \leq \sqrt{n}\|\mathbf{u} - \mathbf{y}\|_2$ , and  $\|\mathbf{X}_i^{(s)}(t)\|_F \leq 2c_{x0}$  in (E.7). Then we establish

$$\begin{aligned}\|\mathbf{W}_s(t+1) - \mathbf{W}_s(0)\|_F &\leq \|\mathbf{W}_s(t+1) - \mathbf{W}_s(t)\|_F + \|\mathbf{W}_s(t) - \mathbf{W}_s(0)\|_F \\ &\leq \frac{8c_{x0}\|\mathbf{y} - \mathbf{u}(0)\|_2}{\lambda\sqrt{n}} \stackrel{\textcircled{1}}{\leq} \sqrt{m}\tilde{r},\end{aligned}$$

where ① holds by setting  $\tilde{r} = \frac{8c_{x0}\|\mathbf{y}-\mathbf{u}(0)\|_2}{\lambda\sqrt{mn}}$ . Finally, combining the value of  $\tilde{r}$ , we have  $\tilde{r} = \max\left(\frac{8c_{x0}\|\mathbf{y}-\mathbf{u}(0)\|_2}{\lambda\sqrt{mn}}, \frac{16(1+\alpha_2+2\alpha_3\mu\sqrt{k_c}c_{w0})^l\alpha_{s,3}^{(l)}\mu\sqrt{k_c}c_{x0}c_{w0}}{\lambda\sqrt{mn}}\|\mathbf{y}-\mathbf{u}(0)\|_2\right) \leq c_{w0}$ . Under this setting, we have

$$\begin{aligned}\|\mathbf{W}_s^{(l)}(t+1) - \mathbf{W}_s^{(l)}(t)\|_F &\leq \frac{4c\eta\alpha_{s,3}^{(l)}\mu c_{x0}c_{w0}\sqrt{k_c}}{\sqrt{n}}\|\mathbf{u}(t) - \mathbf{y}\|_F \\ &\leq \frac{4c\eta\alpha_{s,3}^{(l)}\mu c_{x0}c_{w0}\sqrt{k_c}}{\sqrt{n}}\left(1 - \frac{\eta\lambda}{2}\right)^{t/2}\|\mathbf{y} - \mathbf{u}(0)\|_2, \\ \|\mathbf{W}^{(0)}(t+1) - \mathbf{W}^{(0)}(t)\|_F &\leq \frac{4c\eta\mu c_{x0}c_{w0}\sqrt{k_c}}{\sqrt{n}}\|\mathbf{u}(t) - \mathbf{y}\|_F \\ &\leq \frac{4c\eta\mu c_{x0}c_{w0}\sqrt{k_c}}{\sqrt{n}}\left(1 - \frac{\eta\lambda}{2}\right)^{t/2}\|\mathbf{y} - \mathbf{u}(0)\|_2,\end{aligned}$$

where  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l$  with  $\alpha_2 = \max_{s,l}\alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l}\alpha_{s,3}^{(l)}$ . The proof is completed.  $\square$

## E.7 Proof of Lemma 14

*Proof.* We use mathematical induction to prove the results. We first consider  $h = 0$ . According to the definition, we have

$$\begin{aligned}\|\mathbf{X}^{(0)}(k+1) - \mathbf{X}^{(0)}(k)\|_F &= \tau\left\|\sigma(\mathbf{W}^{(0)}(k+1)\Phi(\mathbf{X})) - \sigma(\mathbf{W}^{(0)}(k)\Phi(\mathbf{X}))\right\|_F \\ &\leq \tau\mu\left\|\mathbf{W}^{(0)}(k+1) - \mathbf{W}^{(0)}(k)\right\|_F\|\Phi(\mathbf{X})\|_F \\ &\stackrel{\textcircled{1}}{\leq} \tau\mu\sqrt{k_c}\left\|\mathbf{W}^{(0)}(k+1) - \mathbf{W}^{(0)}(k)\right\|_F \\ &\stackrel{\textcircled{2}}{\leq} \frac{4c\tau\eta\mu^2c_{x0}c_{w0}k_c}{\sqrt{n}}\|\mathbf{u}(k) - \mathbf{y}\|_F,\end{aligned}$$

where ① uses  $\|\Phi(\mathbf{X})\|_F \leq \sqrt{k_c}\|\mathbf{X}\|_F \leq \sqrt{k_c}$  where the sample  $\mathbf{X}$  obeys  $\|\mathbf{X}\|_F = 1$ ; ② uses the result in Lemma 13 that  $\|\mathbf{W}^{(0)}(t+1) - \mathbf{W}^{(0)}(t)\|_F \leq \frac{4c\eta\mu c_{x0}c_{w0}\sqrt{k_c}}{\sqrt{n}}\|\mathbf{u}(t) - \mathbf{y}\|_F$ .

Then we first consider  $h \geq 1$ .

$$\begin{aligned}&\left\|\mathbf{X}^{(l)}(k+1) - \mathbf{X}^{(l)}(k)\right\|_F \\ &= \left\|\sum_{s=0}^{l-1}\left(\alpha_{s,2}^{(l)}(\mathbf{X}^{(s)}(k+1) - \mathbf{X}^{(s)}(k)) + \alpha_{s,3}^{(l)}\tau\left(\sigma(\mathbf{W}_s^{(l)}(k+1)\Phi(\mathbf{X}^{(s)}(k+1))) - \sigma(\mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}^{(s)}(k)))\right)\right)\right\|_F \\ &\leq \sum_{s=0}^{l-1}\left[\alpha_{s,2}^{(l)}\left\|\mathbf{X}^{(s)}(k+1) - \mathbf{X}^{(s)}(k)\right\|_F + \alpha_{s,3}^{(l)}\tau\left\|\sigma(\mathbf{W}_s^{(l)}(k+1)\Phi(\mathbf{X}^{(s)}(k+1))) - \sigma(\mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}^{(s)}(k)))\right\|_F\right] \\ &\leq \sum_{s=0}^{l-1}\left[\alpha_{s,2}^{(l)}\left\|\mathbf{X}^{(s)}(k+1) - \mathbf{X}^{(s)}(k)\right\|_F + \alpha_{s,3}^{(l)}\tau\mu\left\|\mathbf{W}_s^{(l)}(k+1)\Phi(\mathbf{X}^{(s)}(k+1)) - \mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}^{(s)}(k))\right\|_F\right]\end{aligned}$$

Then we bound the second term carefully:

$$\begin{aligned}&\left\|\mathbf{W}_s^{(l)}(k+1)\Phi(\mathbf{X}^{(s)}(k+1)) - \mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}^{(s)}(k))\right\|_F \\ &= \left\|\mathbf{W}_s^{(l)}(k+1)(\Phi(\mathbf{X}^{(s)}(k+1)) - \Phi(\mathbf{X}^{(s)}(k)))\right\|_F + \left\|(\mathbf{W}_s^{(l)}(k+1) - \mathbf{W}_s^{(l)}(k))\Phi(\mathbf{X}^{(s)}(k))\right\|_F \\ &\leq \sqrt{k_c}\left\|\mathbf{W}_s^{(l)}(k+1)\right\|_F\left\|\mathbf{X}^{(s)}(k+1) - \mathbf{X}^{(s)}(k)\right\|_F + \sqrt{k_c}\left\|\mathbf{W}_s^{(l)}(k+1) - \mathbf{W}_s^{(l)}(k)\right\|_F\left\|\mathbf{X}^{(s)}(k)\right\|_F\end{aligned}$$

By using Lemma 11 and Lemma 10, we have

$$\begin{aligned}\|\mathbf{X}^{(s)}(k)\| &\leq \|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0)\|_F + \|\mathbf{X}_i^{(l)}(0)\|_F \\ &\leq c_{x0} + \left(1 + \alpha_2 + \alpha_3\mu\sqrt{k_c}(\tilde{r} + c_{w0})\right)^l\mu\sqrt{k_c}\tilde{r} \stackrel{\textcircled{1}}{\leq} 2c_{x0},\end{aligned}$$

where  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ , and  $c_{x0} \geq 1$  is given in Lemma 10. ① holds since in Lemma 13, we set  $m$  large enough such that  $\tilde{r}$  is enough small.

Besides, Lemma E.7 shows that

$$\|\mathbf{W}_s^{(l)}(k+1) - \mathbf{W}_s^{(l)}(k)\|_F \leq \frac{4c\eta\alpha_{s,3}^{(l)}\mu c_{x0}c_{w0}\sqrt{k_c}}{\sqrt{n}} \|\mathbf{u}(k) - \mathbf{y}\|_F,$$

where  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . Combing all results yields

$$\begin{aligned} & \left\| \mathbf{W}_s^{(l)}(k+1)\Phi(\mathbf{X}^{(s)}(k+1)) - \mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}^{(s)}(k)) \right\|_F \\ & \leq 2\sqrt{k_c}m c_{w0} \left\| \mathbf{X}^{(s)}(k+1) - \mathbf{X}^{(s)}(k) \right\|_F + \frac{8c\eta\alpha_{s,3}^{(l)}\mu^2 c_{x0}^2 c_{w0} k_c}{\sqrt{n}} \|\mathbf{u}(k) - \mathbf{y}\|_F. \end{aligned}$$

Thus, we can further obtain

$$\begin{aligned} & \left\| \mathbf{X}^{(l)}(k+1) - \mathbf{X}^{(l)}(k) \right\|_F \\ & \leq \sum_{s=0}^{l-1} \left[ (\alpha_{s,2}^{(l)} + 2\sqrt{k_c}c_{w0}\alpha_{s,3}^{(l)}) \left\| \mathbf{X}^{(s)}(k+1) - \mathbf{X}^{(s)}(k) \right\|_F + \frac{8\tau c\eta(\alpha_{s,3}^{(l)})^2 \mu^2 c_{x0}^2 c_{w0} k_c}{\sqrt{n}} \|\mathbf{u}(k) - \mathbf{y}\|_F \right] \\ & \stackrel{\text{①}}{\leq} \sum_{s=0}^{l-1} \left[ (\alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu) \left\| \mathbf{X}^{(s)}(k+1) - \mathbf{X}^{(s)}(k) \right\|_F + \frac{8\tau c\eta(\alpha_3)^2 \mu^2 c_{x0}^2 c_{w0} k_c}{\sqrt{n}} \|\mathbf{u}(k) - \mathbf{y}\|_F \right] \\ & \leq (1 + \alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu)^l \left( \left\| \mathbf{X}^{(0)}(k+1) - \mathbf{X}^{(0)}(k) \right\|_F + \frac{8\tau c\eta(\alpha_3)^2 \mu^2 c_{x0}^2 c_{w0} k_c}{(\alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu)\sqrt{n}} \|\mathbf{u}(k) - \mathbf{y}\|_F \right) \\ & \leq (1 + \alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu)^l \left( \frac{4c\tau\eta\mu^2 c_{x0}c_{w0}k_c}{\sqrt{n}} + \frac{8\tau c\eta(\alpha_3)^2 \mu^2 c_{x0}^2 c_{w0} k_c}{(\alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu)\sqrt{n}} \right) \|\mathbf{u}(k) - \mathbf{y}\|_F \\ & \leq (1 + \alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu)^l \left( 1 + \frac{2(\alpha_3)^2 c_{x0}}{(\alpha_2 + 2\sqrt{k_c}c_{w0}\alpha_3\mu)\sqrt{n}} \right) \frac{4c\tau\eta\mu^2 c_{x0}c_{w0}k_c}{\sqrt{n}} \|\mathbf{u}(k) - \mathbf{y}\|_F. \end{aligned}$$

The proof is completed.  $\square$

## E.8 Proof of Lemma 15

*Proof.* In Lemma 13, we have show

$$\max \left( \|\mathbf{W}^{(0)}(t) - \mathbf{W}^{(0)}(0)\|_F, \|\mathbf{W}_s^{(l)}(t) - \mathbf{W}_s^{(l)}(0)\|_F, \|\mathbf{W}_s(t) - \mathbf{W}_s(0)\|_F \right) \leq \sqrt{m}\tilde{r} \leq \sqrt{m}c_{w0}. \quad (27)$$

Note =  $\frac{1}{\sqrt{m}}$ . In this way, from Lemma 13, we have

$$\left\| \mathbf{W}^{(0)}(t) \right\|_F \leq \left\| \mathbf{W}^{(0)}(t) - \mathbf{W}^{(0)}(0) \right\|_F + \left\| \mathbf{W}^{(0)}(0) \right\|_F \leq 2\sqrt{m}c_{w0},$$

$$\left\| \mathbf{W}_s^{(l)}(t) \right\|_F \leq \left\| \mathbf{W}_s^{(l)}(t) - \mathbf{W}_s^{(l)}(0) \right\|_F + \left\| \mathbf{W}_s^{(l)}(0) \right\|_F \leq 2\sqrt{m}c_{w0},$$

$$\left\| \mathbf{W}_h(t) \right\|_F \leq \left\| \mathbf{W}_h(t) - \mathbf{W}_h(0) \right\|_F + \left\| \mathbf{W}_h(0) \right\|_F \leq 2\sqrt{m}c_{w0}$$

In Lemma 10, we show that when Eqn. (27) holds which is proven in Lemma 13, then  $\|\mathbf{X}_i^{(l)}(0)\|_F \leq c_{x0}$ . Under Eqn. (10), Lemma 11 shows

$$\left\| \mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0) \right\|_F \leq \left( 1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0} \right)^l \mu\sqrt{k_c}\tilde{r} \stackrel{\text{①}}{\leq} c_{x0},$$

where ① holds since in Lemma 13, we set  $m = \mathcal{O} \left( \frac{k_c^2 c_{w0}^2 \|\mathbf{y} - \mathbf{u}(0)\|_2^2}{\lambda^2 n} (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^{4h} \right)$  such that

$$\begin{aligned} \tilde{r} &= \frac{8c_{x0}\|\mathbf{y} - \mathbf{u}(0)\|_2}{\lambda\sqrt{mn}} \max \left( 1, 2 \left( 1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0} \right)^l \alpha_{s,3}^{(l)} \mu\sqrt{k_c}c_{w0} \right) \\ &\leq \frac{c_{x0}}{\left( 1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0} \right)^l \mu\sqrt{k_c}}. \end{aligned}$$

Therefore, we have

$$\left\| \mathbf{X}_i^{(l)}(k) \right\|_F \leq \left\| \mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0) \right\|_F + \left\| \mathbf{X}_i^{(l)}(0) \right\|_F \leq 2c_{x0}.$$

The proof is completed.  $\square$

## E.9 Proof of Lemma 16

*Proof.* We first consider  $l = 0$ . Specifically, we have

$$\begin{aligned} \|\mathbf{X}_i^{(0)}(k) - \mathbf{X}_i^{(0)}(0)\|_F &= \tau \left\| \sigma(\mathbf{W}^{(0)}(k)\Phi(\mathbf{X}_i)) - \sigma(\mathbf{W}^{(0)}(0)\Phi(\mathbf{X}_i)) \right\|_F \\ &\leq \tau \mu \left\| \mathbf{W}^{(0)}(k) - \mathbf{W}^{(0)}(0) \right\|_F \|\Phi(\mathbf{X}_i)\|_F \\ &\stackrel{\textcircled{1}}{\leq} \tau \mu \sqrt{k_c} \left\| \mathbf{W}^{(0)}(k) - \mathbf{W}^{(0)}(0) \right\|_F \\ &\stackrel{\textcircled{2}}{\leq} \mu \sqrt{k_c} \tilde{r}, \end{aligned}$$

where  $\textcircled{1}$  holds since  $\|\Phi(\mathbf{X}_i)\|_F \leq \sqrt{k_c} \|\mathbf{X}_i\|_F \leq \sqrt{k_c}$  and the results in Lemma 13 that  $\left\| \mathbf{W}^{(0)}(k) - \mathbf{W}^{(0)}(0) \right\|_F \leq \sqrt{m} \tilde{r}$ .

Then we consider  $l \geq 1$ . According to the definition, we have

$$\begin{aligned} &\|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0)\|_F \\ &= \left\| \sum_{s=0}^{l-1} \left( \alpha_{s,2}^{(l)} (\mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0)) + \alpha_{s,3}^{(l)} \tau \left( \sigma(\mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}_i^{(s)}(k))) - \sigma(\mathbf{W}_s^{(l)}(0)\Phi(\mathbf{X}_i^{(s)}(0))) \right) \right) \right\|_F \\ &\leq \sum_{s=0}^{l-1} \left[ \alpha_{s,2}^{(l)} \left\| \mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0) \right\|_F + \alpha_{s,3}^{(l)} \tau \left\| \sigma(\mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}_i^{(s)}(k))) - \sigma(\mathbf{W}_s^{(l)}(0)\Phi(\mathbf{X}_i^{(s)}(0))) \right\|_F \right] \\ &\leq \sum_{s=0}^{l-1} \left[ \alpha_{s,2}^{(l)} \left\| \mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0) \right\|_F + \alpha_{s,3}^{(l)} \tau \mu \left\| \mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}_i^{(s)}(k)) - \mathbf{W}_s^{(l)}(0)\Phi(\mathbf{X}_i^{(s)}(0)) \right\|_F \right]. \end{aligned}$$

Then we bound

$$\begin{aligned} &\left\| \mathbf{W}_s^{(l)}(k)\Phi(\mathbf{X}_i^{(s)}(k)) - \mathbf{W}_s^{(l)}(0)\Phi(\mathbf{X}_i^{(s)}(0)) \right\|_F \\ &\leq \left\| (\mathbf{W}_s^{(l)}(k) - \mathbf{W}_s^{(l)}(0))\Phi(\mathbf{X}_i^{(s)}(k)) \right\|_F + \left\| \mathbf{W}_s^{(l)}(0)(\Phi(\mathbf{X}_i^{(s)}(k)) - \Phi(\mathbf{X}_i^{(s)}(0))) \right\|_F \\ &\leq \left\| \mathbf{W}_s^{(l)}(k) - \mathbf{W}_s^{(l)}(0) \right\|_F \left\| \Phi(\mathbf{X}_i^{(s)}(k)) \right\|_F + \left\| \mathbf{W}_s^{(l)}(0) \right\|_F \left\| \Phi(\mathbf{X}_i^{(s)}(k)) - \Phi(\mathbf{X}_i^{(s)}(0)) \right\|_F \\ &\stackrel{\textcircled{1}}{\leq} 2\sqrt{k_c m} c_{x0} \tilde{r} + 2\sqrt{k_c m} c_{w0} \left\| \mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0) \right\|_F, \end{aligned}$$

where  $\textcircled{1}$  holds since Lemma 13 shows  $\left\| \mathbf{W}^{(0)}(k) - \mathbf{W}^{(0)}(0) \right\|_F \leq \sqrt{m} \tilde{r}$  and Lemma 15 shows  $\left\| \mathbf{X}_i^{(s)}(k) \right\|_F \leq 2c_{x0}$  and  $\left\| \mathbf{W}_s^{(l)}(0) \right\|_F \leq 2\sqrt{m} c_{w0}$ .

In this way, we have

$$\begin{aligned} &\|\mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0)\|_F \\ &\leq \sum_{s=0}^{l-1} \left[ \left( \alpha_{s,2}^{(l)} + 2\alpha_{s,3}^{(l)} \mu \sqrt{k_c} c_{w0} \right) \left\| \mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0) \right\|_F + 2\alpha_{s,3}^{(l)} \mu \sqrt{k_c} c_{x0} \tilde{r} \right] \\ &\stackrel{\textcircled{1}}{\leq} \sum_{s=0}^{l-1} \left[ \left( \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0} \right) \left\| \mathbf{X}_i^{(s)}(k) - \mathbf{X}_i^{(s)}(0) \right\|_F + 2\alpha_3 \mu \sqrt{k_c} c_{x0} \tilde{r} \right] \\ &\stackrel{\textcircled{2}}{\leq} c \left[ \left\| \mathbf{X}_i^{(0)}(k) - \mathbf{X}_i^{(0)}(0) \right\|_F + 2\alpha_3 \mu \sqrt{k_c} c_{x0} \tilde{r} \right] \\ &= c(1 + 2\alpha_3 c_{x0}) \mu \sqrt{k_c} \tilde{r} \end{aligned}$$

where  $\textcircled{1}$  and  $\textcircled{2}$  hold by using  $c = (1 + \alpha_2 + 2\alpha_3 \mu \sqrt{k_c} c_{w0})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . The proof is completed.  $\square$

## E.10 Proof of Lemma 17

*Proof.* For this proof, we need to use the results in other lemmas. Specifically, Lemma 13

$$\left\| \mathbf{W}^{(0)}(t) - \mathbf{W}^{(0)}(0) \right\|_F \leq \sqrt{m} \tilde{r}, \quad \left\| \mathbf{W}_s^{(l)}(t) - \mathbf{W}_s^{(l)}(0) \right\|_F \leq \sqrt{m} \tilde{r}, \quad \left\| \mathbf{W}_s(t) - \mathbf{W}_s(0) \right\|_F \leq \sqrt{m} \tilde{r}, \quad (28)$$

where  $c = (1 + \alpha_2 + 2\alpha_3\mu\sqrt{k_c}c_{w0})^l$  with  $\alpha_2 = \max_{s,l} \alpha_{s,2}^{(l)}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}^{(l)}$ . Based on this, Lemma 15 further shows

$$\left\| \mathbf{W}^{(0)}(k) \right\|_F \leq 2\sqrt{m}c_{w0}, \quad \left\| \mathbf{W}_s^{(l)}(k) \right\|_F \leq 2\sqrt{m}c_{w0}, \quad \left\| \mathbf{W}_s(k) \right\|_F \leq 2\sqrt{m}c_{w0}, \quad \left\| \mathbf{X}_i^{(l)}(k) \right\|_F \leq 2c_{x0}. \quad (29)$$

Next, Lemma 16 also proves

$$\left\| \mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0) \right\|_F \leq c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c}\tilde{r}.$$

Then we can easily obtain our result:

$$\begin{aligned} |u_i(k) - u_i(0)| &= \left| \sum_{s=1}^h \langle \mathbf{W}_s(k), \mathbf{X}_i^{(l)}(k) \rangle - \langle \mathbf{W}_s(0), \mathbf{X}_i^{(l)}(0) \rangle \right| \\ &\leq \sum_{s=1}^h \left| \langle \mathbf{W}_s(k) - \mathbf{W}_s(0), \mathbf{X}_i^{(l)}(k) \rangle + \langle \mathbf{W}_s(0), \mathbf{X}_i^{(l)}(k) - \mathbf{X}_i^{(l)}(0) \rangle \right| \\ &\leq \sum_{s=1}^h 2\sqrt{m}\tilde{r}c_{x0} + 2\sqrt{m}c_{w0}c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c}\tilde{r} \\ &= 2\sqrt{m}h \left( c_{x0} + c_{w0}c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c} \right) \tilde{r}. \end{aligned}$$

Then we look at the second part. We first look at  $l = h$ :

$$\begin{aligned} \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(0)} \right\|_F &= \left\| (u_i(k) - y_i)\mathbf{W}_l(k) - (u_i(0) - y_i)\mathbf{W}_l(0) \right\|_F \\ &= |u_i(k) - y_i| \left\| \mathbf{W}_l(k) \right\|_F + |u_i(0) - y_i| \left\| \mathbf{W}_l(0) \right\|_F \\ &\leq \left\| (u_i(k) - u_i(0))\mathbf{W}_l(k) \right\|_F + \left\| (u_i(0) - y_i)(\mathbf{W}_l(k) - \mathbf{W}_l(0)) \right\|_F \\ &\leq |u_i(k) - u_i(0)| \left\| \mathbf{W}_l(k) \right\|_F + |u_i(0) - y_i| \left\| (\mathbf{W}_l(k) - \mathbf{W}_l(0)) \right\|_F \\ &\leq 4\sqrt{m}\tilde{r} \left( c_{w0}\sqrt{m}h \left( c_{x0} + c_{w0}c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c} \right) + |u_i(0) - y_i| \right). \end{aligned} \quad (30)$$

Then we consider  $l < h$ . According to the definitions in Lemma 8, we have

$$\frac{\partial \ell}{\partial \mathbf{X}^{(l)}} = (u - \mathbf{y})\mathbf{W}_l + \sum_{s=l+1}^h \left( \alpha_{i,2}^{(s)} \frac{\partial \ell}{\partial \mathbf{X}^{(s)}} + \alpha_{i,3}^{(s)} \tau \Psi \left( (\mathbf{W}_l^{(s)})^\top \left( \sigma' \left( \mathbf{W}_l^{(s)} \Phi(\mathbf{X}^{(l)}) \right) \odot \frac{\partial \ell}{\partial \mathbf{X}^{(s)}} \right) \right) \right).$$

In this way, we can upper bound

$$\begin{aligned} &\left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(0)} \right\|_F \\ &= \left\| (u_i(k) - y_i)\mathbf{W}_l(k) - (u_i(0) - y_i)\mathbf{W}_l(0) \right\|_F + \sum_{s=l+1}^h \alpha_{i,2}^{(s)} \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(0)} \right\|_F + \sum_{s=l+1}^h \alpha_{i,3}^{(s)} \tau \sqrt{k_c} D, \end{aligned}$$

where  $D = \left\| \mathbf{A}_k^\top (\mathbf{B}_k \odot \mathbf{C}_k) - \mathbf{A}_0^\top (\mathbf{B}_0 \odot \mathbf{C}_0) \right\|_F$  in which  $\mathbf{A}_k = \mathbf{W}_l^{(s)}(k)$ ,  $\mathbf{B}_k = \sigma' \left( \mathbf{W}_l^{(s)}(k) \Phi(\mathbf{X}_i^{(l)}(k)) \right)$ ,  $\mathbf{C}_k = \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(k)}$ . Similar to Eqn. (30), we have

$$\begin{aligned} &\left\| (u_i(k) - y_i)\mathbf{W}_l(k) - (u_i(0) - y_i)\mathbf{W}_l(0) \right\|_F \\ &\leq 4\sqrt{m}\tilde{r} \left( c_{w0}\sqrt{m}h \left( c_{x0} + c_{w0}c(1 + 2\alpha_3c_{x0})\mu\sqrt{k_c} \right) + |u_i(0) - y_i| \right). \end{aligned}$$

Then, we can bound  $D$  as follows:

$$\begin{aligned} D &= \left\| (\mathbf{A}_k - \mathbf{A}_0)^\top (\mathbf{B}_0 \odot \mathbf{C}_0) \right\|_F + \left\| \mathbf{A}_k^\top (\mathbf{B}_k \odot \mathbf{C}_k - \mathbf{B}_0 \odot \mathbf{C}_0) \right\|_F \\ &\leq \left\| \mathbf{A}_k - \mathbf{A}_0 \right\|_F \left\| \mathbf{B}_0 \odot \mathbf{C}_0 \right\|_F + \left\| \mathbf{A}_k \right\|_F \left\| \mathbf{B}_k \odot \mathbf{C}_k - \mathbf{B}_0 \odot \mathbf{C}_0 \right\|_F \\ &\stackrel{\textcircled{1}}{\leq} \mu\sqrt{m}\tilde{r} \left\| \mathbf{C}_0 \right\|_2 + 2\sqrt{m}c_{w0} \left\| \mathbf{B}_k \odot \mathbf{C}_k - \mathbf{B}_0 \odot \mathbf{C}_0 \right\|_F \end{aligned}$$



where ① uses the results in Eqns. (29) and (28). The remaining work is to bound

$$\begin{aligned} \|\mathbf{B}_k \odot \mathbf{C}_k - \mathbf{B}_0 \odot \mathbf{C}_0\|_F &= \|\mathbf{B}_k \odot (\mathbf{C}_k - \mathbf{C}_0)\|_F + \|(\mathbf{B}_k - \mathbf{B}_0) \odot \mathbf{C}_0\|_F \\ &\leq \mu \|\mathbf{C}_k - \mathbf{C}_0\|_F + \rho \left\| \mathbf{W}_l^{(s)}(k) \Phi(\mathbf{X}_i^{(l)}(k)) - \mathbf{W}_l^{(s)}(0) \Phi(\mathbf{X}_i^{(l)}(0)) \right\|_F \|\mathbf{C}_0\|_\infty \end{aligned}$$

where ① uses the assumption that the activation function  $\sigma(\cdot)$  is  $\mu$ -Lipschitz and  $\rho$ -smooth. Note  $\|\mathbf{C}_0\|_\infty$  is a constant, since it is the gradient norm at the initialization which does not involve the algorithm updating. Recall Lemma 11 shows

$$\left\| \mathbf{W}_s^{(l)}(k) \Phi(\mathbf{X}^{(s)}(k)) - \mathbf{W}_s^{(l)}(0) \Phi(\mathbf{X}^{(s)}(0)) \right\|_F \leq \frac{1}{\alpha_3} \left(1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0})\right)^l \sqrt{k_c m \tilde{r}},$$

where  $\alpha_2 = \max_{s,l} \alpha_{s,2}$  and  $\alpha_3 = \max_{s,l} \alpha_{s,3}$ , and  $c_{x0} \geq 1$  is given in Lemma 10. Then we upper bound

$$\left\| \mathbf{W}_l^{(s)}(k) \Phi(\mathbf{X}_i^{(l)}(k)) - \mathbf{W}_l^{(s)}(0) \Phi(\mathbf{X}_i^{(l)}(0)) \right\|_F \leq \frac{1}{\alpha_3} \left(1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0})\right)^l \sqrt{k_c m \tilde{r}}.$$

Therefore, we have

$$D \leq \mu \sqrt{m \tilde{r}} \|\mathbf{C}_0\|_2 + 2\sqrt{m} c_{w0} \left( \mu \|\mathbf{C}_k - \mathbf{C}_0\|_F + \frac{\rho \|\mathbf{C}_0\|_\infty}{\alpha_3} \left(1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0})\right)^l \sqrt{k_c m \tilde{r}} \right)$$

By combining the above results, we have

$$\begin{aligned} &\left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(0)} \right\|_F \\ &\leq c_1 + \sum_{s=l+1}^h \left[ \left( \alpha_{l,2}^{(s)} + 2\alpha_{l,3}^{(s)} \sqrt{k_c} \mu c_{w0} \right) \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(0)} \right\|_F + c_2 \right] \\ &\leq c_1 + \sum_{s=l+1}^h \left[ \left( \alpha_2 + 2\alpha_3 \sqrt{k_c} \mu c_{w0} \right) \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(s)}(0)} \right\|_F + c_3 \right] \\ &\leq \left(1 + \alpha_2 + 2\alpha_3 \sqrt{k_c} \mu c_{w0}\right)^l \left[ \left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(h)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(h)}(0)} \right\|_F + c_3 \right] \end{aligned}$$

where  $c_1 = 4\sqrt{m \tilde{r}} (c_{w0} \sqrt{m} h (c_{x0} + c_{w0} c (1 + 2\alpha_3 c_{x0}) \mu \sqrt{k_c}) + |u_i(0) - y_i|)$ ,  $c_2 = \alpha_{l,3}^{(s)} \left( \mu \tilde{r} \|\mathbf{C}_0\|_2 + 2c_{w0} \frac{\rho \|\mathbf{C}_0\|_\infty}{\alpha_3} \left(1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0})\right)^l \sqrt{k_c m \tilde{r}} \right)$  and  $c_3 = \alpha_3 \left( \mu \tilde{r} \|\mathbf{C}_0\|_2 + 2c_{w0} \frac{\rho \|\mathbf{C}_0\|_\infty}{\alpha_3} \left(1 + \alpha_2 + \alpha_3 \mu \sqrt{k_c} (r + c_{w0})\right)^l \sqrt{k_c m \tilde{r}} \right)$ . Consider  $\|\mathbf{C}_0\|_2 = \mathcal{O}(\sqrt{m})$ , for brevity, we ignore constants and obtain

$$\left\| \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(k)} - \frac{\partial \ell}{\partial \mathbf{X}_i^{(l)}(0)} \right\|_F \leq c_1 c \alpha_3^2 c_{w0}^2 c_{x0} \rho k_c m \tilde{r},$$

where  $c = (1 + \alpha_2 + 2\alpha_3 \sqrt{k_c} \mu c_{w0})^l$  and  $c_1$  is a constant. The proof is completed.  $\square$

### E.11 Proof of Lemma 18

*Proof.* By Assumption 2, each entry for the initial parameter  $\mathbf{W}_s^{(l)}(0)$  obeys Gaussian distribution  $\mathcal{N}(0, 1)$ . Then  $\|\mathbf{W}_s^{(l)}(0)\|_F^2$  is chi-square variable with freedom degree  $k_c p m$ . In this way, by using Lemma 4, we have

$$\mathbb{P} \left( \|\mathbf{W}_s^{(l)}(0)\|_F^2 - k_c p m \geq 2\sqrt{k_c p m t} + 2t \right) \leq \exp(-t).$$

Therefore, with probability at least  $1 - \frac{\delta}{2h(h+3)}$ , we can obtain

$$\|\mathbf{W}_s^{(l)}(0)\|_F \leq \sqrt{k_c p m + 2\sqrt{k_c p m \log(2h(h+3)/\delta)} + 2\log(2h(h+3)/\delta)} \leq \sqrt{m} c_{w0},$$

where  $c_{w0} \sim \sqrt{k_c p}$  is a constant. Note here we focus on  $m$  more than  $p$  and  $k_c$ , since  $m$  is much larger than  $p$  and  $k_c$  which is introduced in subsequent analysis.

By using the same method, we can prove that with probability at least  $1 - \frac{\delta}{2h(h+3)}$ ,

$$\|\mathbf{W}^0(0)\|_F \leq \sqrt{m}c_{w0} \quad \text{and} \quad \|\mathbf{W}_s(0)\|_F \leq \sqrt{m}c_{w0}$$

In this way, with probability at least  $\left(1 - \frac{\delta}{2h(h+3)}\right)^{\frac{h(h+3)}{2}} \geq 1 - \frac{\delta}{2h(h+3)} \frac{h(h+3)}{2} = 1 - \delta/4$ , these results hold at the same time. The proof is completed.  $\square$

## References

- [1] I. Loshchilov and F. Hutter. SGDR: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [2] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *Int'l Conf. Learning Representations*, 2014.
- [3] J. Saw, M. Yang, and T. Mo. Chebyshev inequality with estimated mean and variance. *The American Statistician*, 38(2):130–132, 1984.
- [4] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai. Gradient descent finds global minima of deep neural networks. In *Proc. Int'l Conf. Machine Learning*, 2019.
- [5] S. Hwang. Cauchy's interlace theorem for eigenvalues of hermitian matrices. *The American Mathematical Monthly*, 111(2):157–159, 2004.
- [6] R. Alessandro. 36-755: Advanced statistics theory. *UC Berkeley Lecture*, [http://www.stat.cmu.edu/~arinaldo/Teaching/36755/F17/Scribed\\_Lectures/F17\\_0911.pdf](http://www.stat.cmu.edu/~arinaldo/Teaching/36755/F17/Scribed_Lectures/F17_0911.pdf), 2017.
- [7] C. Louizos, M. Welling, and D. Kingma. Learning sparse neural networks through  $\ell_0$  regularization. In *Int'l Conf. Learning Representations*, 2018.