

---

# How do students become teachers: A dynamical analysis for two-layer neural networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This paper investigates the fundamental regression task of learning  $k$  neurons  
2 (*a.k.a.* teachers) from Gaussian input, using two-layer ReLU neural networks  
3 with width  $m$  (*a.k.a.* students) and  $m, k = \mathcal{O}(1)$ , trained via gradient descent  
4 under proper initialization and a small step-size. Our analysis follows a three-  
5 phase structure: *alignment* after weak recovery, *tangential growth*, and *local*  
6 *convergence*, providing deeper insights into the learning dynamics of gradient  
7 descent (GD). We prove the global convergence at the rate of  $\mathcal{O}(T^{-3})$  for the zero  
8 loss of excess risk. Additionally, our results show that GD automatically groups  
9 and balances student neurons, revealing an implicit bias toward achieving the  
10 minimum “balanced”  $\ell_2$ -norm in the solution. Our work extends beyond previous  
11 studies in exact-parameterization setting ( $m = k = 1$ , [Yehudai and Ohad, 2020])  
12 and single-neuron setting ( $m \geq k = 1$ , [Xu and Du, 2023]). The key technical  
13 challenge lies in handling the interactions between multiple teachers and students  
14 during training, which we address by refining the alignment analysis in Phase 1 and  
15 introducing a new dynamic system analysis for tangential components in Phase 2.  
16 Our results pave the way for further research on optimizing neural network training  
17 dynamics and understanding implicit biases in more complex architectures.

## 18 1 Introduction

19 Learning a target function  $f^* : \mathbb{R}^d \rightarrow \mathbb{R}$  via neural networks through gradient descent or flow has  
20 received significant attention in machine learning theory. Research in this area primarily focuses on  
21 understanding the learnability and dynamics, aiming to theoretically explain the advantage of *feature*  
22 *learning* in neural networks. This problem is often studied under various assumptions about  $f^*$ . For  
23 instance,  $f^*$  is frequently (implicitly) assumed to be smooth in a kernel regime [Jacot et al., 2018,  
24 Allen-Zhu et al., 2019, Arora et al., 2019]. Additionally,  $f^*$  might possess further structures, such as  
25 being located on a low-dimensional subspace [Mousavi-Hosseini et al., 2023] or a manifold [Arora  
26 et al., 2022]. A typical example is assuming  $f^*$  is a sparse polynomial [Abbe et al., 2022]. In this  
27 setting, the separation between kernel methods and neural networks is well studied through metrics  
28 like the information exponent [Arous et al., 2021], leap complexity [Abbe et al., 2023], and generative  
29 exponent [Damian et al., 2024].

30 In contrast to smooth functions, another research direction focuses on non-smooth target functions,  
31 such as ReLU. This non-smoothness naturally highlights the difference between kernel methods and  
32 neural networks in terms of approximation ability [Bach, 2017]. As a result, researchers have turned  
33 their attention to studying the learning dynamics to gain a deeper understanding of convergence. For  
34 instance, they investigate learning with a single ReLU neuron [Wu et al., 2023, Xu and Du, 2023] or  
35 multiple ReLU neurons [Zhou et al., 2021, Akiyama and Suzuki, 2023].

36 We consider the problem of learning one-hidden-layer ReLU networks under the Gaussian measure.  
 37 The target function  $f^*$  is a sum of multiple ReLU neurons  $f^*(\mathbf{x}) := \sum_{l=1}^k \sigma(\langle \mathbf{v}_l, \mathbf{x} \rangle)$  with the  
 38 parameters  $\{\mathbf{v}_l\}_{l=1}^k$ , which can be learned from  $n$  i.i.d. samples  $\{(\mathbf{x}_i, f^*(\mathbf{x}_i))\}_{i=1}^n$  via a two-layer  
 39 neural network with  $m$  (student) neurons with random Gaussian initialization  $\{\mathbf{w}_i\}_{i=1}^m \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$   
 40 under the expected squared loss:

$$L(\mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} \left( \sum_{i=1}^m \sigma(\mathbf{w}_i^\top \mathbf{x}) - f^*(\mathbf{x}) \right)^2, \quad (1)$$

41 which aims to find a good approximation of  $f^*$  from the student network. To ensure learning  
 42 performance,  $m \geq k$  is needed.

43 This problem is identified as an additive model in statistics [Stone, 1985, Hastie and Tibshirani, 1987],  
 44 and it receives great attention in theoretical computer science [Chen et al., 2023] and machine learning  
 45 theory, especially on sample/time complexity as well as training dynamics [Boursier and Flammarion,  
 46 2024, Bietti et al., 2023]. However, understanding how gradient-based training algorithms recover  
 47 the teacher network and analyzing the entire training dynamics are still challenging. Therefore, most  
 48 current analyses are limited to non-gradient-based algorithms [Chen et al., 2023], or *local analysis*  
 49 for gradient-based algorithms, which assumes that the loss has already been minimized below a  
 50 very small threshold, or the angles between teacher neurons and their nearest student neurons are  
 51 already small (called *strong recovery*), e.g., [Zhou et al., 2021]. If we go beyond the local analysis,  
 52 previous result on GD training can only handle specific cases, such as [Yehudai and Ohad, 2020]  
 53 for  $m = k = 1$ , [Wu et al., 2018] for  $m = k = 2$ , and [Xu and Du, 2023, Chistikov et al., 2023] for  
 54  $m \geq k = 1$ . In fact, studying more general cases, such as  $m, k = \mathcal{O}(1)$ , remains unresolved, even in  
 55 local analysis. Accordingly, we aim to address the following question:

56 **How can gradient-based algorithms recover teacher neurons and learn useful features beyond**  
 57 **the local analysis?**

58 To better understand the learning dynamics in the above question, we follow the ‘‘align then fit’’  
 59 framework [Maennel et al., 2018, Boursier and Flammarion, 2024], which also helps to explain the  
 60 implicit bias of the learned solution. In this study, we run the gradient descent (GD) over Eq. (1).  
 61 Since analyzing the entire training dynamics is still challenging and is an open problem, so we assume  
 62 the *weak recovery*, where for each student neuron, exactly one teacher neuron exists that is not nearly  
 63 perpendicular to it. Note that the weak recovery condition is still much weaker than the condition  
 64 with local analysis and strong recovery that will be proved in our analysis. An informal version of  
 65 our theoretical results is given as below.

66 **Theorem 1** (Global Convergence after Weak Recovery: Informal). *Under proper assumptions*  
 67 *(e.g., teacher neurons are with same length  $\|\mathbf{v}\|$ , and orthogonal to each other), sufficiently small*  
 68 *initialization with  $\sigma = o(\text{poly}(d^{-\frac{1}{2}}))$ , and trained via gradient descent with sufficiently small step-size*  
 69  *$\eta = o(1)$  to minimize Eq. (1), after time  $T^* = \Omega(\frac{1}{\eta})$ , for any  $T \in \mathbb{N}$ , we have:*

$$L(\mathbf{W}(T^* + T)) \leq \mathcal{O}\left(\frac{\|\mathbf{v}\|^2}{\eta^3 T^3}\right), \quad \text{and} \quad \|\mathbf{w}_i(T^* + T)\| = \Theta\left(\frac{k \|\mathbf{v}\|}{m}\right) \quad \forall i \in [m], w.h.p.$$

70 Our result demonstrates that the Eq. (1) can be solved by GD in the polynomial time to find the  
 71 global minima and achieves the global convergence rate at  $\mathcal{O}(1/T^3)$ . We admit that the derived  
 72 sample/time complexity is not optimal, but to our knowledge, this is the first polynomial-time result  
 73 of GD training beyond the local analysis for Eq. (1) with  $m, k = \mathcal{O}(1)$ . Besides, our results also  
 74 indicate that the obtained solution will converge to a minimum ‘‘balanced’’  $\ell_2$  solution, where the  
 75 ‘‘balanced’’ is determined by student neurons and their respective nearest teacher neurons.

76 **Technical challenges.** We employ the similar proof framework of Xu and Du [2023] on  $m \geq k = 1$ .  
 77 The main challenge of this paper is how to address the coupling of different teacher neurons’ influences  
 78 on the student neurons, even though the teacher neurons are orthogonal to each other. For instance:

- 79 • In phase 1, single teacher neuron ( $k = 1$ ) [Xu and Du, 2023] allows for monotonic convergence
- 80 on the angular difference between the teacher and student neurons. However, this does not hold
- 81 for  $k > 1$ . In this case, we use approximations of sine and cosine values for decoupling when the

82 angle is very small or near perpendicular. Hence we can simplify the training dynamics and prove  
 83 that the sine of the minimum angle converges linearly to a tiny neighborhood.

- 84 • In phase 2, during the analysis of neuron growth, the tangential components of the student neurons  
 85 at each teacher neuron (and for more teacher neurons) are quite complex. Classical recursive  
 86 relationship in [Xu and Du, 2023] can not handle this. Instead, we develop a new technical tool by  
 87 building a dynamical system: we formulate the matrix iteration form, estimate the eigenvalues of  
 88 the transition matrix, and establish the upper and lower bounds of such a dynamical system.

## 89 2 Notations, problem setting, and assumptions

90 In this section, we give notations that are needed in this paper and then introduce our problem setting  
 91 as well as the required assumptions in our proof.

### 92 2.1 Notations

93 *Basic notations:* We use the shorthand  $[n] := \{1, 2, \dots, n\}$  for a positive integer  $n$ . We denote by  
 94  $a(n) \gtrsim b(n)$ : the inequality  $a(n) \geq cb(n)$  that hides a positive constant  $c$  that is independent of  $n$ .  
 95 Vectors (matrices) are denoted by boldface, lower-case (upper-case) letters. The used norm  $\|\cdot\|$  in  
 96 this paper is  $\ell_2$  norm if we do not specify. We follow the standard Bachmann–Landau notation in  
 97 complexity theory e.g.,  $\mathcal{O}$ ,  $o$ ,  $\Omega$ , and  $\Theta$  for order notation.

98 *Notations on angle:* The angle between any two non-zero vectors  $\mathbf{w}$  and  $\mathbf{v}$  is denoted as  $\angle(\mathbf{w}, \mathbf{v}) :=$   
 99  $\cos^{-1} \frac{\langle \mathbf{w}, \mathbf{v} \rangle}{\|\mathbf{w}\| \|\mathbf{v}\|}$ . Then we use the following notations for any  $i, j \in [m], l \in [k]$

- 100 •  $\theta_{il} \triangleq \angle(\mathbf{w}_i, \mathbf{v}_l)$ : the angle between a student neuron  $\mathbf{w}_i$  and a teacher neuron  $\mathbf{v}_l$ .
- 101 •  $\varphi_{ij} \triangleq \angle(\mathbf{w}_i, \mathbf{w}_j)$ : the angle between two neurons  $\mathbf{w}_i$  and  $\mathbf{w}_j$  for student model.
- 102 •  $\tau_i \triangleq \arg \min_j \angle(\mathbf{w}_i(0), \mathbf{v}_j(0))$ : the index of the teacher neuron with **the smallest angle** to  
 103 the  $\mathbf{w}_i$  at initialization, in which the smallest angle is denoted as  $\theta_{i^*} \triangleq \theta_{i\tau_i} = \angle(\mathbf{w}_i, \mathbf{v}_{\tau_i})$ .
- 104 •  $\mathbf{r}_j \triangleq \sum_{i:\tau_i=l} \mathbf{w}_i - \mathbf{v}_l$ : the difference of the teacher neuron  $\mathbf{v}_l$  and the sum of the student  
 105 neurons around  $\mathbf{v}_l$ .

106 For notational simplicity, by denoting  $\bar{\mathbf{a}} \triangleq \frac{\mathbf{a}}{\|\mathbf{a}\|}$ , we denote the tangential part  $h_{il} \triangleq \langle \mathbf{w}_i, \bar{\mathbf{v}}_l \rangle$  as the  
 107 projection of  $\mathbf{w}_i$  along with the direction of  $\mathbf{v}_l$ ; and a similar notation for  $h_{i^*} \triangleq \langle \mathbf{w}_i, \bar{\mathbf{v}}_{\tau_i} \rangle$ . Besides,  
 108 we denote  $\mathbf{w}_i(t)$  as the vector  $\mathbf{w}_i$  at time  $t$ , which also adapts to  $\theta_{ij}(t)$ , etc.

109 *Notations on loss:* The standard Gaussian distribution is  $\mathcal{N}(0, 1)$  with zero-mean and unit variance.  
 110 We denote  $\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0,1)}$  by  $\mathbb{E}_{\mathbf{x}}$  for simplicity. By defining the residuals of the neural network as:

$$R(\mathbf{x}) := \sum_{i=1}^m \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) - \sum_{i=1}^k \sigma(\langle \mathbf{v}_i, \mathbf{x} \rangle),$$

111 then the loss can be written as  $L(\mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} R(\mathbf{x})^2$ .

### 112 2.2 Closed form expressions of gradient of loss: $\nabla L$

113 To make our paper self-contained, we present the closed-form expressions for  $\nabla L$  when the input data  
 114 follows a Gaussian distribution, as given by Safran and Shamir [2018], see the details in Appendix B.

115 We denote  $\nabla_i \triangleq \frac{\partial L(\mathbf{W})}{\partial \mathbf{w}_i}$  as the gradient of loss to the  $\mathbf{w}_i$ , when  $\mathbf{w}_i \neq \mathbf{0}$ . Then for any  $i \in [m]$ , the  
 116 loss function is differentiable with gradient given by:

$$\nabla_i = \frac{1}{2} \sum_{j=1}^m \mathbf{w}_j - \frac{1}{2} \sum_{l=1}^k \mathbf{v}_l + \frac{1}{2\pi} \left[ \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|} \left( \sum_{j=1, j \neq i}^m \sin \varphi_{ij} \|\mathbf{w}_j\| - \sum_{l=1}^k \sin \theta_{il} \|\mathbf{v}_l\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij} \mathbf{w}_j + \sum_{l=1}^k \theta_{il} \mathbf{v}_l \right]. \quad (2)$$

117 We use random Gaussian initialization for neural network training, i.e.,  $\forall i \in [m], \mathbf{w}_i(0) \sim$   
 118  $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$  with the variance  $\sigma^2$ . Then we can prove that  $\|\mathbf{w}_i\|$  has bounded norm with high  
 119 probability if the dimension  $d$  is not small, see Lemma 1 in Appendix B.

## 120 2.3 Assumptions

121 We state the used assumptions in this paper.

122 **Assumption 1** (Weak recovery). *Regarding the angle  $\theta_{ij}(0)$  defined before for any  $i \in [m], j \in [k]$ ,  
123 at initialization, denote  $\theta_{i^*}(0)$  as the smallest angle between  $\mathbf{w}_i$  and its closet teacher neuron. The  
124 weak recovery assumes  $\theta_{i^*}(0) \ll \theta_{ij}(0)$  with  $j \in [k]$  and  $j \neq \tau_i$ . We mathematically formulate this  
125 as below.*

- 126 •  $\theta_{i^*}(0)$  is acute:  $0 < \frac{\pi}{2} - \theta_{i^*}(0) \triangleq \zeta_i = \Theta(1)$ , and  $\zeta_i \in (0, \frac{\pi}{2}]$ .
- 127 •  $\theta_{ij}(0)$  is close to orthogonal:  $|\frac{\pi}{2} - \theta_{ij}(0)| \leq \zeta = o(1)$  with  $j \in [k]$  and  $j \neq \tau_i$ .

128 **Remark:** The weak recovery assumption requires that a student neuron is not orthogonal to its closet  
129 student neuron but is nearly orthogonal to the remaining teacher neurons [Dandi et al., 2024]. If  
130 we focus on the single ReLU case like Xu and Du [2023], this assumption can be directly removed  
131 because there is only one teacher neuron. With only one teacher neuron, there are no competing  
132 neurons for alignment, and thus the angle between the student and teacher neuron is naturally the  
133 smallest.

134 **Assumption 2** (Orthogonal and same norm for teacher neurons). *The teacher neurons are given by  
135  $\{\mathbf{v}_i\}_{i=1}^k$ , and are assumed to be orthogonal to each other with the same norm, i.e.,  $\langle \mathbf{v}_i, \mathbf{v}_j \rangle = 0$  and  
136  $\|\mathbf{v}_i\| = \|\mathbf{v}_j\| = \|\mathbf{v}\|$ ,  $\forall i \neq j, i, j \in [m]$ . Clearly, we have  $k \leq d$  due to the orthogonality of  $k$   
137 teacher neurons.*

138 **Remark:** This assumption requires all teacher neurons pointing to different (orthogonal) directions,  
139 which is important for identifiability or recovery. It aligns with practical considerations by allowing  
140 diverse tasks such that the target feature directions do not significantly overlap. This assumption  
141 as well as its variant (e.g., separation among teacher neurons) has been widely used in previous  
142 theoretical results, e.g., [Zhou et al., 2021, Oko et al., 2024, Simsek et al., 2023]. We can relax this  
143 assumption where the teacher neurons are nearly orthogonal and have similar norms. However, such  
144 relaxation would require additional computations in our analysis. To avoid unnecessary complexity  
145 and focus on the core analysis, we concentrate on the basic assumptions.

146 **Assumption 3** (Balance condition at initialization). *At initialization, we record the number of student  
147 neurons  $\mathbf{w}_i$  with  $\tau_i = l$  as  $m_l$ . Then we assume  $\frac{m}{3k} \leq m_l \leq \frac{3m}{k}, \forall l \in [k]$ .*

148 **Remark:** This is a balance condition such that the number of merged student neurons among each  
149 teacher neuron is not extremely small or large. It is motivated by Boursier et al. [2022, Assumption  
150 3] and [Wojtowysch, 2020] that requires the student neurons to cover all directions of the teacher  
151 neurons. Our assumption requires student neurons coincide with teacher neurons in a *balanced* way.

## 152 3 Main Results

153 In this section, we will provide the main theoretical results. First, in Section 3.1, we provide the  
154 primary result on global convergence. Then, in the following subsections, we discuss the training  
155 dynamics of the three phases and provide proof sketches. In Section 3.2.1, we provide the main  
156 dynamics and final state results of the alignment process in the first phase. In Section 3.2.2, we  
157 provide the main dynamics and final state results of the tangential growth process in the second phase.  
158 In Section 3.2.3, we provide the results of the local convergence in the third phase and then achieve  
159 the final global convergence result.

### 160 3.1 Main Theorem

**Theorem 2.** *Assume  $d = \Omega(\log(m/\delta))$  with  $\delta \in (0, 1)$ , under Assumptions 1 2 and 3, let  
 $\sigma = o(\text{poly}(m^{-k^2}, d^{-\frac{1}{2}})) = o(\text{poly}(d^{-\frac{1}{2}}))$ , and trained via gradient descent with step-size,  $\eta =$   
 $o(\text{poly}(m^{-k^2})) = o(\text{poly}(1))$  to minimize Eq. (1), then there exists a  $T^* = \Omega(\frac{k \log k \log m}{m\eta}) = \Omega(\frac{1}{\eta})$   
such that with probability at least  $1 - \delta$  over the initialization, for any  $T \in \mathbb{N}$ , we have:*

$$L(\mathbf{W}(T^* + T)) \leq \mathcal{O}\left(\frac{k^{12} \|\mathbf{v}\|^2}{\eta^3 T^3}\right), \quad \text{and} \quad \frac{\|\mathbf{v}\|}{4m_{\tau_i}} \leq \|\mathbf{w}_i(T^* + T)\| \leq \frac{4\|\mathbf{v}\|}{m_{\tau_i}} \quad \forall i \in [m].$$

161 **Remark:** Theorem 2 provides a convergence rate of  $T^{-3}$ , which is consistent with previous results  
 162 [Xu and Du, 2023]. Moreover, it indicates that the more teacher neurons and the larger their norms,  
 163 the slower the convergence rate. This aligns with our intuition that when the initialization is very  
 164 small, a larger norm and more teacher neurons require student neurons to take more time to align  
 165 and converge to the teacher neurons. Unlike [Xu and Du, 2023], we present a stronger bound that  
 166 is independent of the number of student neurons and does not deteriorate as the number of student  
 167 neurons increases. Furthermore, our results indicate that the student neurons will implicitly converge  
 168 to a specific teacher neuron and maintain a balance among themselves.

### 169 3.2 Proof overview

170 In this section, we provide a sketch of the Theorem 2. The complete proof can be found in the  
 171 appendix. Our proof is primarily divided into three phases: alignment (Section 3.2.1), tangential  
 172 growth (Section 3.2.2), and global convergence (Section 3.2.3). Finally we can summarize the results  
 173 in these three phases for the main result.

#### 174 3.2.1 Phase 1 - Alignment

175 During this phase, each student neuron individually aligns with a specific teacher neuron. The  
 176 outcomes of this section are divided into two main parts: *i*) the upper and lower bounds on the  
 177 lengths of the student neurons, as well as the angle between student and teacher neurons during the  
 178 time Theorem 3. *ii*) the upper bound on the angle between student and teacher neurons at the end of  
 179 phase 1, as well as the balance of projection strength from different student neurons onto the teacher  
 180 neuron Corollary 1. The detailed derivation can be found in Appendix D.

181 **Theorem 3** (Phase 1: Alignment Process). *Assume  $d = \Omega(\log(m/\delta))$  with  $\delta \in (0, 1)$ , for any  $\epsilon_1 > 0$ ,  
 182 under Assumption 1 with  $\epsilon_1^2 = o(1)$  and Assumptions 2, 3 such that  $\sigma = o(\frac{\text{poly}(\epsilon_1)\|\mathbf{v}\|}{\sqrt{d}})$  in our random  
 183 Gaussian initialization, and the stepsize satisfies  $\eta = o(\frac{\sigma\sqrt{d}\epsilon_1^2}{\|\mathbf{v}\|})$ , then there exist a  $T_1 = \Theta(\frac{\epsilon_1^2}{\eta})$ , for  
 184  $0 \leq t \leq T_1$ , the following statements hold with probability at least  $1 - \delta$ :*

$$s_1 \leq \|\mathbf{w}_i(t)\| \leq s_2 + 2k\eta\|\mathbf{v}\|t, \quad \forall i \in [m], \quad \text{with } s_1 := \frac{1}{2}\sigma\sqrt{d}, \quad s_2 := 2\sigma\sqrt{d}, \quad (3)$$

185 and

$$\sin^2\left(\frac{\theta_{i^*}(t)}{2}\right) - \epsilon_1^2 \leq \left(1 + \frac{\eta k\|\mathbf{v}\|t}{s_2}\right)^{-\frac{1}{sk}} \left(\sin^2\left(\frac{\theta_{i^*}(0)}{2}\right) - \epsilon_1^2\right), \quad \forall i \in [m]. \quad (4)$$

186 **Remark:** Our theorem implies that, during phase 1 of the training, the norm of each student  
 187 neuron always has an immutable lower bound, while the upper bound increases linearly over time.  
 188 Additionally, for each student neuron, the angle with its nearest teacher neuron converges linearly  
 189 within an error range of  $\epsilon_1^2$ . Compared to the results of Xu and Du [2023], the upper bound of our  
 190 neuron norm increases  $k$  times faster because we have  $k$  different teacher neurons, which naturally  
 191 leads to this outcome. Taking  $k = 1$ , our convergence rate is faster by a constant factor compared to  
 192 the results of Xu and Du [2023], and our condition for  $\sigma$  is weaker. When the same  $\sigma$  is selected, the  
 193 total duration of phase 1, denoted as  $T_1$ , they are the same.

194 Then, we briefly introduce our proof technique, due to the presence of multiple teacher neurons, the  
 195 gradient expression in Eq. (15) contains  $2(k-1)$  cross terms including  $\theta_{il}$  with detailed interactions,  
 196 which do not exist in Xu and Du [2023]. To handle this challenge, we provide additional analysis on  
 197 alignment related to these cross terms in phase 1. Specifically, we prove these results by induction.

198 **Proof of Eq. (3):** This formula provides the upper and lower bounds of  $\|\mathbf{w}\|$  during the training. For  
 199 the lower bound, based on the gradient expression  $\nabla_i$  in Eq. (2), we prove  $\langle \mathbf{w}_i(t), \nabla_i(t) \rangle \leq 0$ , which  
 200 ensures that the norm  $\|\mathbf{w}\|$  increases monotonically such that  $\|\mathbf{w}_i(t)\| \geq s_1$ . For the upper bound,  
 201 we need to bound the norm of gradient using Eq. (2). Then, applying the triangle inequality, we can  
 202 obtain the desired result.

203 **Proof of Eq. (4):** This formula illustrates the angle dynamics (i.e., the alignment process) between  
 204 the student neuron and its closest teacher neuron during phase 1. For larger  $\theta_{i^*}$ , it is easy to prove  
 205 that the  $\theta_{i^*}$  decreases monotonically. However, when  $\theta_{i^*}$  is very small, we cannot guarantee the  
 206 monotonic decreasing property of  $\theta_{i^*}$ . To this end, we build the connection between  $\sin^2(\theta_{i^*}(t)/2)$

207 and the following angle difference

$$\cos(\theta_{i^*}(t+1)) - \cos(\theta_{i^*}(t)) := I_2 + I_3, \quad \text{the first-order term } I_2 \text{ and the second-order term } I_3.$$

208 By estimating  $I_2$  and  $I_3$ , we can track the dynamics of the angle difference and then prove that  $\sin \theta_{i^*}$   
209 converges linearly to a very small neighborhood (i.e.,  $\epsilon_1^2$ ).

210 At the end of phase 1, we conclude the following result:

211 **Corollary 1** (Final State of Phase 1). *Under the same conditions as Theorem 3, at time  $T_1$ , the*  
212 *following statements hold with probability at least  $1 - \delta$ :*

$$\theta_{i^*}(T_1) \leq 4\epsilon_1, \quad \text{and} \quad h_{i^*}(T_1) \leq 2h_{j^*}(T_1), \quad \forall i, j \in [m]. \quad (5)$$

213 **Remark:** By the end of phase 1, each student neuron will align with its nearest teacher neuron with  
214 the residual angle at the order of  $\mathcal{O}(\epsilon_1)$ . Additionally, the projection lengths of these student neurons  
215 in the direction of their corresponding teacher neurons are relatively balanced, with a rough upper  
216 bound of 2.

217 **Proof of Eq. (5):** For the first part, substituting the parameters from Theorem 3 into Eq. (4) will yield  
218 the result. For the second part, firstly, we derive the upper and lower bounds for  $h_{i^*}(t+1) - h_{i^*}(t)$   
219 and then accumulate these bounds. Next, we prove that before a certain time (e.g.,  $t := T_1/50$ ),  
220 the upper bound of  $h_{i^*}(t)$  is relatively small compared to this accumulated value. This allows us to  
221 establish the upper and lower bounds for all  $h_{i^*}(T_1)$  and thereby determine the maximum ratio of  
222  $h_{i^*}(T_1)$  among different student neurons.

### 223 3.2.2 Phase 2 - Tangential Growth

224 In this section, we present the results of the second phase, in which each student neuron grows along  
225 the tangential direction of the teacher neuron aligned in phase 1 as below. The detailed derivation can  
226 be found in Appendix E.

227 **Theorem 4** (Phase 2: Tangential Growth Process). *Assume  $d = \Omega(\log(m/\delta))$  with  $\delta \in (0, 1)$ ,*  
228 *for any  $\epsilon_1 > 0, \epsilon_2 > 0$ , under Assumption 1 with  $\epsilon_2 = o(1), \epsilon_1^2 = o(\text{poly}(\epsilon_2))$ , Assumptions 2, 3*  
229 *such that  $\sigma = o(\frac{\text{poly}(\epsilon_1)\|\mathbf{v}\|}{\sqrt{d}})$  in our random Gaussian initialization, and the stepsize satisfies  $\eta =$   
230  $o(\frac{\sigma\sqrt{d}\epsilon_1^2}{\|\mathbf{v}\|})$ , then by setting there exist a  $T_2 = T_1 + \Theta(\frac{1}{\eta} \ln(\frac{1}{\epsilon_2}))$ , then  $\forall T_1 \leq t \leq T_2$ , we define  
231  $H_l(t) := \|\mathbf{v}\| - \sum_{i=1}^m \mathbb{I}_{\tau_i=l} h_{i^*}(t)$  for  $l \in [k]$ , the following statements hold with probability at least  
232  $1 - \delta$ :*

$$h_{i^*}(t) \leq 2h_{j^*}(t), \quad \text{and} \quad \frac{2\|\mathbf{v}\|}{m_{\tau_i}} \geq h_{i^*}(t) \geq \frac{s_1}{2}, \quad \forall i, j \in [m] \text{ and } \tau_i = \tau_j. \quad (6)$$

$$\left(1 - \frac{\eta m}{9k}\right)^{t-T_1} \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\| \geq H_l(t) \geq \frac{2}{3} \|\mathbf{v}\| \left(1 - \frac{3\eta m}{2k}\right)^{t-T_1} - 8\pi\epsilon_2 \|\mathbf{v}\| \geq 24\pi\epsilon_2 \|\mathbf{v}\|, \quad \forall l \in [k], \quad (7)$$

234 and

$$\theta_{i^*}(t) \leq \epsilon_2, \quad \forall i \in [m]. \quad (8)$$

235 **Remark:** This theorem tells us that during phase 2:

236 1). The norm of student neurons close to the same teacher neuron remains relatively balanced, with  
237 each neuron having strict upper and lower bounds (Eq. (6)). It is worth noting that, unlike in phase  
238 one, see Eq. (5), this balance is not maintained for all neurons.

239 2). The projections of the student neurons near each teacher neuron will gradually increase, and the  
240 difference from  $\|\mathbf{v}\|$  will approach zero at a linear convergence rate (Eq. (7)). This result implies that  
241 as training progresses, the loss gradually decreases. We will further prove that by the end of phase 2,  
242 the loss has decreased to a sufficiently small value.

243 3). The angle between each student neuron and its nearest teacher neuron stays within a small range  
244 (Eq. (8)). However, the angle is slightly larger than that of Phase 1 because additional cost/movement  
245 is required to handle the convergence for tangential difference and the decrease of loss. For example,  
246 we have  $\|\nabla_i(t)\| \leq 2k\|\mathbf{v}\|$  in Phase 1 but it changes to  $\|\nabla_i(t)\| \leq 15k\|\mathbf{v}\|$  in Phase 2.

247 4). Taking  $k = 1$ , our condition for  $\epsilon_2$  is similar to that of Xu and Du [2023], but we have relaxed the  
 248 learning rate condition by a factor of  $m$ . And the total duration of phase 2, is reduced by a constant  
 249 factor of  $\frac{1}{2}$ .

250 Then, we briefly introduce our proof technique. Compared to one teacher setting [Xu and Du,  
 251 2023], the tangential analysis requires a new dynamical system analysis regarding the dynamics  
 252 of  $\{H_l(t)\}_{l=1}^k$  due to the coupling tangential components among student/teacher neurons. Besides,  
 253 the loss function becomes more complex Eq. (14) and we have to control the loss below a certain  
 254 threshold in the presence of these interactions, which requires additional quantities to estimate.  
 255 Specifically, we prove these results by induction.

256 **Proof of Eq. (6):** For the first part, we follow the proof of Eq. (5) to build the connection between  
 257  $h_{i^*}(t+1) - h_{i^*}(t)$  and  $H_l$  in a weighted sum relationship, with an additional constant term  $Q_i$ .  
 258 For two different student neurons close to the same teacher neuron, these weights are the same. By  
 259 studying the changes of  $\theta_{i^*}$  and  $\theta_{il}$  during this phase,  $|Q_i(t)|$  will be bounded by a small quantity.  
 260 Then we conclude the result by summing and combining the results with Eq. (5). For the second  
 261 part, based on Eq. (7), we can derive  $H_l \geq 0$  and finish the upper bound by combining the results  
 262 from the first part. For the lower bound, we derive  $h_{i^*}(t+1) - h_{i^*}(t) \geq 0$ , which implies that  $h_{i^*}$  is  
 263 monotonic increasing. Combining this with Eqs. (3) and (5), the proof is complete.

264 **Proof of Eq. (7):** Using the above analysis about  $h_{i^*}(t+1) - h_{i^*}(t) \geq 0$  and the relationship  
 265 between  $h_{i^*}$  and  $H_{\tau_i}$ , we can establish a recursive relationship between  $H(t+1)$  and  $H(t)$  as well.  
 266 Note that there is a coupling between different  $\mathbf{H}$  and interference from small quantities  $Q_i$ , so we  
 267 express the iterative formula in matrix form. To be specific, by denoting  $\mathbf{H} := \{H_l\}_{l=1}^k$  (we write it  
 268 in a matrix formulation), it admits the following recursive relationship:

$$\mathbf{H}(t+1) = \mathbf{A}\mathbf{H}(t) + \mathbf{Q}(t) \quad \text{for a certain transition matrix } \mathbf{A} \text{ and } \mathbf{Q}(t) \text{ depends on } Q_i(t).$$

269 By analyzing the eigenvalues of the transition matrix  $\mathbf{A}$ , we estimate the upper and lower bounds  
 270 of such a dynamic system. For the small quantities  $Q_i$ , we adopt the same approach used in  
 271 proving Eq. (4). Finally, we prove that  $\mathbf{H}$  converges to a small value at a linear convergence rate.

272 **Proof of Eq. (8):** The proof here is similar to Eq. (4), as it also analyzes the dynamics of  $\cos \theta_{i^*}$ .  
 273 However, the difficulty lies in that at this phase, the influence of  $\mathbf{w}$  in the gradient is no longer  
 274 negligible compared to  $\mathbf{v}$ , making the iterative relationship between angles more complex. First, by  
 275 proving

$$\frac{\|\mathbf{w}_i(t)\|}{\|\mathbf{w}_j(t)\|} = \Theta(1), \forall i, j \in [m], T_1 \leq t \leq T_2,$$

276 we are able to analyse the dynamics of  $\cos \theta_{i^*}$  (i.e.,  $I_2$  and  $I_3$  in Eq. (4)) based on two cases  
 277  $\tau_i = (\neq) \tau_j$ . First, we use some properties of trigonometric functions to decouple this relationship so  
 278 that it only involves the coupling between each student neuron and its nearest teacher neuron. Then,

279 we estimate the difference  $\sin^2 \left( \frac{\theta_{i^*}(t+1)}{2} \right) - \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right)$  for the final  $\sin \theta_{i^*}(t)$ . Unlike in phase  
 280 1, here we obtain an upper bound for the linear growth of the angle  $\theta_{i^*}$ . However, we can still prove  
 281 that within the range of  $T_2$ , the angle remains small.

282 At the end of phase 2, we can draw the following results:

283 **Corollary 2** (Final state of Phase 2). *Under the same conditions as Theorem 4, at time  $T_2$ , the*  
 284 *following statements hold with probability at least  $1 - \delta$ :*

$$\frac{\|\mathbf{v}\|}{3m\tau_i} \leq \|\mathbf{w}_i(T_2)\| \leq \frac{3\|\mathbf{v}\|}{m\tau_i}, \forall i \in [m], \quad \text{and} \quad L(\mathbf{W}(T_2)) \leq \frac{1}{2}k^2\epsilon_2^{0.05}\|\mathbf{v}\|^2. \quad (9)$$

285 **Remark:** After phase 2, the norms of each student neuron have balanced, and the loss has decreased  
 286 to a very small value. This provides the foundation for proving local convergence in phase three.

287 **Proof of first part of Eq. (9):** We use the results in Theorem 4 to prove this result. For the lower  
 288 bound, we first observe from Eq. (7) that  $H_l$  is very small at time  $T_2$ , meaning the sum of  $h$  among  
 289 student neurons near each teacher neuron is close to  $\|\mathbf{v}\|$ , i.e.,  $H_l(T_2) \leq \frac{\|\mathbf{v}\|}{3}$ . Using the balance  
 290 of them in Eq. (6), we can then establish a lower bound for  $h_{i^*}(\|\mathbf{w}_i\| \cos \theta_{i^*})$ , which further allows  
 291 us to derive a lower bound for  $\|\mathbf{w}_i\|$ . Similarly, for the upper bound, we first observe from Eq. (7)  
 292 that at time  $T_2$ , the sum of  $h$  among student neurons near each teacher neuron is close to but still

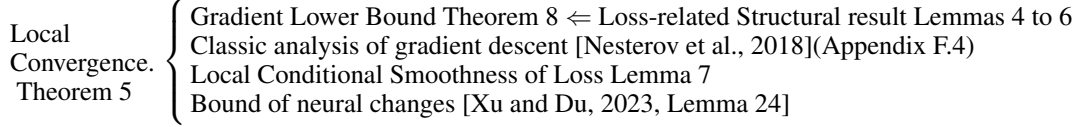


Figure 1: Proof framework of Phase 3 - Local convergence

less than  $\|\mathbf{v}\|$ . Using the balance of them in Eq. (6), we can then establish an upper bound for  $h_{i^*}(\|\mathbf{w}_i\| \cos \theta_{i^*})$ . Given that the angle between each student neuron and its nearest teacher neuron is very small (second part in Eq. (8)), we can further derive a lower bound for  $\|\mathbf{w}_i\|$ .

**Proof of second part of Eq. (9):** The key point of this proof involves introducing an auxiliary function  $g$  to help decompose the  $L$ . The loss  $L$  can be expressed in the summation of  $g$ , see Appendix B for details. First, based on the upper bound of the angle in phase 2 (second part in Eq. (8)), we know that there are two scenarios for the angle in the closed form of the loss: close to 0 and nearly orthogonal. We discuss the upper and lower bounds of auxiliary function  $g$  in these two cases. Then, according to Eq. (7), we find that at time  $T_2$ , the sum of the norms of the student neurons near each teacher neuron close to the norm of teacher neurons, i.e.  $\sum_{i=1}^m \mathbb{I}_{\tau_i=l} h_{i^*}(T_2) \geq (1 - o(1)) \|\mathbf{v}\|$ . Combining these two results, we can derive an upper bound for the loss  $L$ .

### 3.2.3 Phase 3 - Local convergence

In this section, we present the results of phase 3 - local convergence. Specifically, we show that when the loss is already small enough, the loss function converges to zero at a rate of  $\mathcal{O}(\frac{1}{T^3})$  Theorem 5. Our results build upon the previous works of Xu and Du [2023], Zhou et al. [2021], Safran et al. [2021]. The detailed derivation can be found in Appendix F.

**Theorem 5** (Local convergence). *Suppose the initial condition in Lemma 1 and Assumption 1 2 and 3 holds. If we set  $\epsilon_2 = o(\text{poly}(1))$  and  $\eta = o(1)$  in Theorem 4, then  $\forall T \in \mathbb{N}$ , the following statements hold with probability at least  $1 - \delta$ :*

$$L(\mathbf{W}(T + T_2)) \leq \frac{1}{\left( L(\mathbf{W}(T_2))^{-\frac{1}{3}} + \Omega\left(k^{-4} \|\mathbf{v}\|^{-\frac{2}{3}}\right) \eta T \right)^3}, \quad (10)$$

and

$$\frac{\|\mathbf{v}\|}{4m\tau_i} \leq \|\mathbf{w}_i(T + T_2)\| \leq \frac{4\|\mathbf{v}\|}{m\tau_i} \quad \forall i \in [m]. \quad (11)$$

**Remark:** This theorem shows that, under the condition that the loss is less than a very small value and the neurons remain balanced at the end of phase two, GD training can achieve the global minimum with a convergence rate of  $\frac{1}{T^3}$ . This result is consistent with the previous result in Xu and Du [2023] and is superior to  $\frac{1}{T}$  in Zhou et al. [2021]. Furthermore, this result also indicates that, without using regularization during training, every student neuron will implicitly converge to the directions of specific teacher neurons, and there is a balance among student neurons that converge to the direction of the same teacher neuron.

Then, we briefly introduce our proof technique.

**Proof of Eq. (10):** The proof framework of Theorem 5 is standard based on the local convergence analysis, e.g., [Zhou et al., 2021, Xu and Du, 2023], as illustrated in Fig. 1.

The key point to prove Eq. (10) is utilizing the result of classic optimization in Appendix F.4 and the lower bound of the gradient to satisfy the conditions of [Xu and Du, 2023, Lemma 24]. First, we follow [Zhou et al., 2021] to derive several lemmas related to the properties of the loss function. Based on these lemmas, we can obtain the lower bound of the gradient in terms of the loss. Then, similar to Safran et al. [2021], we deduce that when the neurons maintain a certain balance, the loss is locally smooth. This allows us to directly apply the classic optimization theory conclusion regarding the relationship between adjacent iterations of gradient descent Appendix F.4. Finally, we build the final convergence result by Xu and Du [2023, Lemma 24]. Additionally, our proof requires that the balance condition of the neurons is consistently maintained Eq. (11), which can be proven using induction and convergence results alternately.



333 Finally, by combining results from Sections 3.2.1 to 3.2.3 with the hyper-parameter selection in Ap-  
334 pendix A, we obtain the global convergence result in Theorem 2.

335 When  $k = 1$ , compared to the results of Xu and Du [2023], our paper needs stronger requirements on  
336  $\sigma$ ,  $\eta$  and time. This is due to the upper bound of the loss after phase 2 in Eq. (9) and its relationship  
337 with  $\epsilon_2$ . Due to multiple teacher neurons, the number of student neurons converging to each teacher  
338 neuron directions are different. This leads to different norms for the student neurons, which makes a  
339 looser upper bound. However, in the case of  $k = 1$ , such handling is not necessary. Therefore, our  
340 results can cover the results of Xu and Du [2023] with only minor modifications and have a better  
341 constant factor in phase 1 discussed before.

## 342 4 Numerical Validation

343 In this section, we empirically validate our theoretical results by plotting the convergence curves  
344 under the following setting: we set  $\|v\| = 5$ , data dimension  $d = 100$ , batch size of 512, and a  
345 total of 5000 batches. The total number of training samples (equivalent to the previously mentioned  
346  $T^* + T$ ) is  $2.56 \times 10^6$ . Besides, we have added a  $1/T^3$  reference line in the log-log plot for better  
347 comparison.

348 First, we selected four sets of parameters  $k$  and  $m$ :  $k = 2, m = 20$ ,  $k = 4, m = 12$ ,  $k = 4, m = 20$ ,  
349 and  $k = 4, m = 40$  with initialization variance  $\sigma = 10^{-6}$  and learning rate  $\eta = 5 \times 10^{-4}$ . The plots  
350 in Fig. 2 show the cosine of the angle and norm convergence during training (top row) and the log-log  
351 plot of the loss during training (bottom row) for different values of  $k$  and  $m$ . The results show that  
352 larger  $k$  values lead to longer  $t_1$  and  $t_2$  and slower convergence rates, while larger  $m$  values result in  
353 shorter  $t_1$  and  $t_2$  but have little effect on the convergence rate. This matches our theoretical results  
354 such that using more (student) neurons decreases the time for alignment. We admit that learning more  
355 (teacher) neurons generally requires more time but this is given under the same initialization strategy.  
356 Instead, our initialization strategy depends on  $m$  and  $k$ , leading to different learning dynamics.

357 In the second set of experiments, we selected four sets of parameters  $\sigma$  and  $\eta$ :  $\sigma = 10^{-4}, \eta =$   
358  $5 \times 10^{-4}$ ,  $\sigma = 10^{-5}, \eta = 5 \times 10^{-4}$ ,  $\sigma = 10^{-6}, \eta = 5 \times 10^{-4}$ , and  $\sigma = 10^{-6}, \eta = 10^{-3}$  with  
359  $k = 4$  and  $m = 20$ . The plots in Fig. 3 show the cosine of the angle and norm convergence during  
360 training (top row) and the log-log plot of the loss during training (bottom row) for different values of  
361  $\sigma$  and  $\eta$ . The results show that smaller stepsizes and initialization variances are crucial for stable and  
362 predictable training dynamics. Specifically, we found that larger initialization variance and stepsize  
363 can reduce the period in the first two phases, but it slows down the convergence rate in the third phase.  
364 This also suggests that empirically, smaller learning rates and initialization variances are better under  
365 this setting.

366 Regarding the timescale experiments, we divided the training dynamics into three phases for analysis.  
367 We can observe the clear “align then fit” phenomena where in phases 1 and 2, the angle aligns  
368 and the tangential grows until the norm of neurons’ weights is unchanged. In phase 3, the loss  
369 function decreases for fitting data. The phase transition from Phase 1 to 2 is not very clear in the  
370 experiments but can still be observed with a distinct difference in that Phase 2 finishes later than  
371 Phase 1. Nonetheless, we have marked the figure’s approximate endpoints of the first and second  
372 phases.

## 373 5 Related work

374 **Dynamics of gradient descent in the teacher-student setting:** Li and Yuan [2017] studied the  
375 exact-parameterized setting and proved convergence for SGD with initialization near identity. The  
376 separation between kernel methods and two-layer neural networks is further described in Li et al.  
377 [2020]. To further understand the convergence and generalization of regression tasks using non-linear  
378 networks, it is essential to thoroughly analyze the dynamics throughout gradient-based training,  
379 commonly described as “align then fit” [Maennel et al., 2018, Boursier and Flammarion, 2024] in a  
380 three-phase analysis framework. Xu and Du [2023] provide a global convergence of learning with a  
381 single ReLU neuron, where the proof for the local convergence (i.e., the third phase) is given by Zhou  
382 et al. [2021]. This analysis framework is also used in various settings, e.g., binary classification [Min  
383 et al., 2023] and matrix sensing [Xiong et al., 2024].

384 Besides, our problem can be cast as a special case of learning with multi-index model [Bietti et al.,  
385 2023] where the link function (i.e., the activation function used in this work) is unknown. However,  
386 the techniques are different and our three-phase analysis framework allows for a better understanding

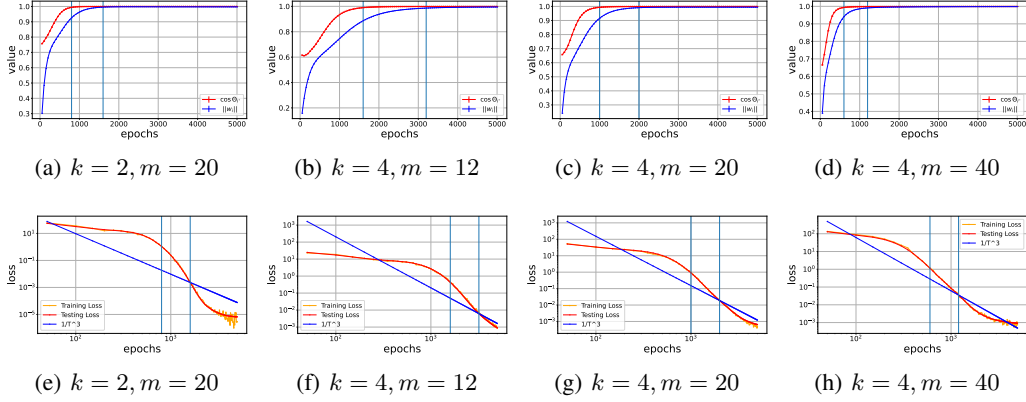


Figure 2: Convergence curves for different  $m$  and  $k$ .

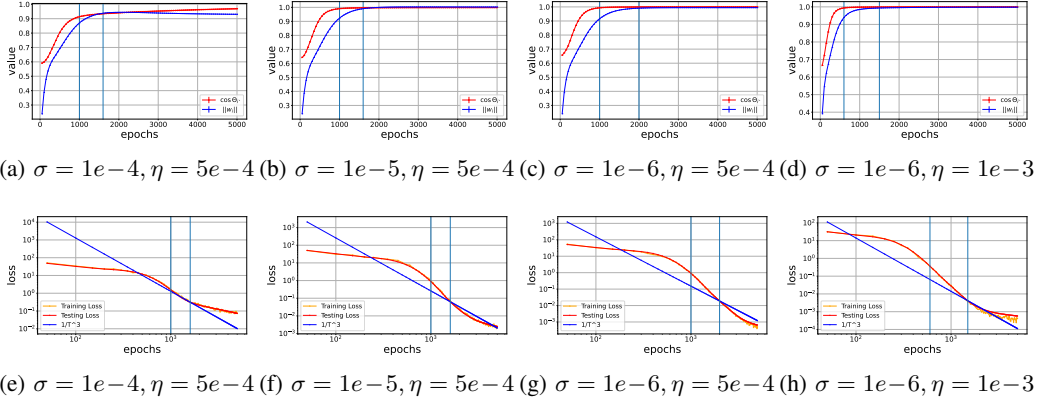


Figure 3: Convergence curves for different  $\sigma$  and  $\eta$ .

387 of global convergence. Some statistical physics studies work have explored related topics but differ  
 388 from our work [Goldt et al., 2020, Arnaboldi et al., 2023] by focusing on generalization errors without  
 389 providing convergence rates or detailed analyses of training dynamics and convergence phases.

390 **Implicit bias:** Recent studies suggest that gradient descent is implicitly biased towards a low-rank  
 391 hidden weight matrix or a sparse number of directions represented by the neurons [Safran et al.,  
 392 2022, Shevchenko et al., 2022, Chizat and Bach, 2020]. This implicit bias is often characterized by  
 393 the minimal norm interpolator, which is closely related to sparsely represented directions [Lyu and  
 394 Li, 2020]. These findings indicate that the early alignment phase enforces the alignment of weights  
 395 towards a small number of directions, even with omnidirectional initialization, leading to implicit  
 396 regularization at convergence [Boursier and Flammarion, 2023].

## 397 6 Conclusion

398 Our three-phase analysis framework provides a comprehensive analysis on global convergence, i.e.,  
 399 1) *alignment*: the angle decreases  $\theta_{i^*}(T_1) \leq \mathcal{O}(\epsilon_1)$  satisfying the balance condition but the norm  
 400 of student neuron gradually increases with  $T_1$ ; 2) *tangential growth*: the projection of the student  
 401 neurons near teacher neurons gradually increases. The angle is still small but slightly larger than  
 402 that of phase 1 due to the additional cost of handling the convergence of tangential difference; 3)  
 403 *local convergence*: the loss is close to zero and the neurons are still well-balanced thus achieving the  
 404 global convergence at the rate of  $\mathcal{O}(T^{-3})$ .

405 One potential drawback of this work is the weak recovery which simplifies the analysis. However,  
 406 without weak recovery, the analysis will be quite complex, remaining unsolved, and thus we leave it  
 407 as future work.

## 408 References

- 409 E. Abbe, E. B. Adsera, and T. Misiakiewicz. The merged-staircase property: a necessary and  
410 nearly sufficient condition for sgd learning of sparse functions on two-layer neural networks. In  
411 *Conference on Learning Theory (COLT)*, 2022.
- 412 E. Abbe, E. B. Adsera, and T. Misiakiewicz. Sgd learning on neural networks: leap complexity and  
413 saddle-to-saddle dynamics. In *Conference on Learning Theory (COLT)*, 2023.
- 414 S. Akiyama and T. Suzuki. Excess risk of two-layer reLU neural networks in teacher-student settings  
415 and its superiority to kernel methods. In *International Conference on Learning Representations*  
416 *(ICLR)*, 2023.
- 417 Z. Allen-Zhu, Y. Li, and Z. Song. A convergence theory for deep learning via over-parameterization.  
418 In *International Conference on Machine Learning (ICML)*, 2019.
- 419 L. Arnaboldi, L. Stephan, F. Krzakala, and B. Loureiro. From high-dimensional & mean-field  
420 dynamics to dimensionless odes: A unifying approach to sgd in two-layers networks, 2023. URL  
421 <https://arxiv.org/abs/2302.05882>.
- 422 S. Arora, S. S. Du, W. Hu, Z. Li, R. R. Salakhutdinov, and R. Wang. On exact computation with an  
423 infinitely wide neural net. In *Advances in neural information processing systems (NeurIPS)*, 2019.
- 424 S. Arora, Z. Li, and A. Panigrahi. Understanding gradient descent on the edge of stability in deep  
425 learning. In *International Conference on Machine Learning (ICML)*, 2022.
- 426 G. B. Arous, R. Gheissari, and A. Jagannath. Online stochastic gradient descent on non-convex losses  
427 from high-dimensional inference. *Journal of Machine Learning Research*, 22(106):1–51, 2021.
- 428 F. Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine*  
429 *Learning Research*, 2017.
- 430 A. Bietti, J. Bruna, and L. Pillaud-Vivien. On learning gaussian multi-index models with gradient  
431 flow. *arXiv:2310.19793*, 2023.
- 432 E. Boursier and N. Flammarion. Penalising the biases in norm regularisation enforces sparsity. In  
433 *Advances in neural information processing systems (NeurIPS)*, 2023.
- 434 E. Boursier and N. Flammarion. Early alignment in two-layer networks training is a two-edged sword.  
435 *arXiv:2401.10791*, 2024.
- 436 E. Boursier, L. Pillaud-Vivien, and N. Flammarion. Gradient flow dynamics of shallow relu networks  
437 for square loss and orthogonal inputs. In *Advances in neural information processing systems*  
438 *(NeurIPS)*, 2022.
- 439 S. Chen, Z. Dou, S. Goel, A. Klivans, and R. Meka. Learning narrow one-hidden-layer relu networks.  
440 In G. Neu and L. Rosasco, editors, *Proceedings of Thirty Sixth Conference on Learning Theory*,  
441 volume 195 of *Proceedings of Machine Learning Research*, pages 5580–5614. PMLR, 12–15 Jul  
442 2023. URL <https://proceedings.mlr.press/v195/chen23a.html>.
- 443 D. Chistikov, M. Englert, and R. Lazic. Learning a neuron by a shallow reLU network: Dynamics  
444 and implicit bias for correlated inputs. In *Advances in neural information processing systems*  
445 *(NeurIPS)*, 2023.
- 446 L. Chizat and F. Bach. Implicit bias of gradient descent for wide two-layer neural networks trained  
447 with the logistic loss. In *Conference on Learning Theory (COLT)*, 2020.
- 448 A. Damian, L. Pillaud-Vivien, J. Lee, and J. Bruna. Computational-statistical gaps in gaussian single-  
449 index models. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 1262–1262.  
450 PMLR, 2024.
- 451 Y. Dandi, E. Troiani, L. Arnaboldi, L. Pesce, L. Zdeborová, and F. Krzakala. The benefits of reusing  
452 batches for gradient descent in two-layer networks: Breaking the curse of information and leap  
453 exponents, 2024. URL <https://arxiv.org/abs/2402.03220>.

- 454 S. Goldt, M. S. Advani, A. M. Saxe, F. Krzakala, and L. Zdeborová. Dynamics of stochastic  
455 gradient descent for two-layer neural networks in the teacher–student setup\*. *Journal of Statistical*  
456 *Mechanics: Theory and Experiment*, 2020(12):124010, Dec. 2020. ISSN 1742-5468. doi: 10.  
457 1088/1742-5468/abc61e. URL <http://dx.doi.org/10.1088/1742-5468/abc61e>.
- 458 T. Hastie and R. Tibshirani. Generalized additive models: some applications. *Journal of the American*  
459 *Statistical Association*, 82(398):371–386, 1987.
- 460 A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural  
461 networks. In *Advances in neural information processing systems (NeurIPS)*, 2018.
- 462 Y. Li and Y. Yuan. Convergence analysis of two-layer neural networks with relu activation. In  
463 *Advances in neural information processing systems (NeurIPS)*, 2017.
- 464 Y. Li, T. Ma, and H. R. Zhang. Learning over-parametrized two-layer relu neural networks beyond  
465 ntk. *arXiv:2007.04596*, 2020.
- 466 K. Lyu and J. Li. Gradient descent maximizes the margin of homogeneous neural networks. In  
467 *International Conference on Learning Representations (ICLR)*, 2020.
- 468 H. Maennel, O. Bousquet, and S. Gelly. Gradient descent quantizes relu network features.  
469 *arXiv:1803.08367*, 2018.
- 470 H. Min, R. Vidal, and E. Mallada. Early neuron alignment in two-layer relu networks with small  
471 initialization. *arXiv:2307.12851*, 2023.
- 472 A. Mousavi-Hosseini, S. Park, M. Girotti, I. Mitliagkas, and M. A. Erdogdu. Neural networks  
473 efficiently learn low-dimensional representations with SGD. In *International Conference on*  
474 *Learning Representations (ICLR)*, 2023.
- 475 Y. Nesterov et al. *Lectures on convex optimization*. Springer, 2018.
- 476 K. Oko, Y. Song, T. Suzuki, and D. Wu. Learning sum of diverse features: computational hardness  
477 and efficient gradient-based training for ridge combinations, 2024. URL <https://arxiv.org/abs/2406.11828>.
- 479 I. Safran and O. Shamir. Spurious local minima are common in two-layer relu neural networks. In  
480 *International Conference on Machine Learning (ICML)*, 2018.
- 481 I. Safran, G. Vardi, and J. D. Lee. On the effective number of linear regions in shallow univariate relu  
482 networks: Convergence guarantees and implicit bias. *arXiv:2205.09072*, 2022.
- 483 I. M. Safran, G. Yehudai, and O. Shamir. The effects of mild over-parameterization on the optimization  
484 landscape of shallow relu neural networks. In *Conference on Learning Theory (COLT)*, 2021.
- 485 A. Shevchenko, V. Kungurtsev, and M. Mondelli. Mean-field analysis of piecewise linear solutions  
486 for wide relu networks. *Journal of Machine Learning Research*, 2022.
- 487 B. Simsek, A. Bendjeddou, W. Gerstner, and J. Brea. Should under-parameterized student networks  
488 copy or average teacher weights? In *Thirty-seventh Conference on Neural Information Processing*  
489 *Systems*, 2023. URL <https://openreview.net/forum?id=MG0mYskXN2>.
- 490 C. J. Stone. Additive regression and other nonparametric models. *The annals of Statistics*, 13(2):  
491 689–705, 1985.
- 492 S. Wojtowytsch. On the convergence of gradient descent training for two-layer relu-networks in the  
493 mean field regime. *arXiv:2005.13530*, 2020.
- 494 C. Wu, J. Luo, and J. D. Lee. No spurious local minima in a two hidden unit reLU network, 2018.
- 495 J. Wu, D. Zou, Z. Chen, V. Braverman, Q. Gu, and S. M. Kakade. Finite-sample analysis of learning  
496 high-dimensional single relu neuron. In *International Conference on Machine Learning (ICML)*,  
497 2023.

- 498 N. Xiong, L. Ding, and S. S. Du. How over-parameterization slows down gradient descent in matrix  
499 sensing: The curses of symmetry and initialization. In *International Conference on Learning*  
500 *Representations (ICLR)*, 2024.
- 501 W. Xu and S. S. Du. Over-parameterization exponentially slows down gradient descent for learning a  
502 single neuron. In *Conference on Learning Theory (COLT)*, 2023.
- 503 G. Yehudai and S. Ohad. Learning a single neuron with gradient methods. In *Conference on Learning*  
504 *Theory (COLT)*, 2020.
- 505 M. Zhou, R. Ge, and C. Jin. A local convergence theory for mildly over-parameterized two-layer  
506 neural network. In *Conference on Learning Theory (COLT)*, 2021.

507 **Appendix introduction**

508 The Appendix is organized as follows:

- 509 • In Appendix A, we discuss the selection of hyperparameters in this paper.
- 510 • In Appendix B, we provide a detailed explanation of the closed-form expression for the loss
- 511 and its gradient as mentioned in the main text.
- 512 • In Appendix C, we provide a detailed analysis of Assumption 1.
- 513 • In Appendix D, we present the main results of phase 1 along with detailed proofs.
- 514 • In Appendix E, we present the main results of phase 2 along with detailed proofs.
- 515 • In Appendix F, we present the main results of phase 3 along with detailed proofs.

516 **A Selection of hyper-parameters**

- 517 • We set  $\epsilon_2 = o(m^{-60}k^{-100}) = o(\text{poly}(1))$  in Theorem 4 as required by Theorem 5.
- 518 • We set  $\epsilon_1^2 = o(\epsilon_2^{\Theta(k)}/m) = o(\text{poly}(\epsilon_2))$  in Theorem 3 as required by Theorem 4.
- 519 • We set  $\sigma \leq \frac{\epsilon_1^{16k+2}\|\mathbf{v}\|}{10000m\sqrt{d}} = o(\frac{\text{poly}(\epsilon_1^2)\|\mathbf{v}\|}{\sqrt{d}})$  in Theorem 3 as required by Theorems 3 and 4.
- 520 • We set  $\eta = o\left(\frac{m\epsilon_1^2s_1^2}{k^2\|\mathbf{v}\|^2}\right) \leq o\left(\frac{\epsilon_1^{32k+6}}{mk}\right) = o(\text{poly}(\epsilon_1^2))$  in Theorem 3 as required by Theo-
- 521 rem 4.
- 522 • We set  $T_1 := \frac{\epsilon_1^2}{100\eta km} = \Theta\left(\frac{\epsilon_1^2}{\eta}\right)$  in Theorem 3.
- 523 • We set  $T_2 = T_1 + \frac{k}{2\eta m} \ln\left(\frac{1}{48\pi\epsilon_2}\right) = \Theta\left(\frac{1}{\eta} \ln \epsilon_2^{-1}\right) = \Omega\left(\frac{1}{\eta}\right)$  in Theorem 4.

524 **B Expression of loss  $L$  and its gradient  $\nabla L$**

525 In this section, we provide a detailed explanation of the closed-form expression for the loss and  
 526 its gradient as mentioned in the main text. The main content of this section follows [Safran and  
 527 Shamir, 2018, Section 4.1.1]. Besides, the bounded norm of  $\|\mathbf{w}_i\|$  for any  $i \in [m]$  is also given in  
 528 this subsection. We include these results here just for completeness.

529 For notational simplicity, we introduce the following auxiliary function:

$$g(\mathbf{a}, \mathbf{b}) := \mathbb{E}_{\mathbf{x}}[\sigma(\mathbf{a}^\top \mathbf{x})\sigma(\mathbf{b}^\top \mathbf{x})] = \frac{\|\mathbf{a}\|\|\mathbf{b}\|}{2\pi} \left( \sin \angle(\mathbf{a}, \mathbf{b}) + (\pi - \angle(\mathbf{a}, \mathbf{b})) \cos \angle(\mathbf{a}, \mathbf{b}) \right), \quad (12)$$

530 which implies that

- 531 • if  $\mathbf{a}$  and  $\mathbf{b}$  are orthogonal, i.e.,  $\langle \mathbf{a}, \mathbf{b} \rangle = 0$ , then  $g(\mathbf{a}, \mathbf{b}) = \frac{\|\mathbf{a}\|\|\mathbf{b}\|}{2\pi}$ .
- 532 • If  $\mathbf{a} = \mathbf{b}$ , then  $g(\mathbf{a}, \mathbf{b}) = \frac{1}{2} \|\mathbf{a}\| \|\mathbf{b}\| = \frac{1}{2} \|\mathbf{a}\|^2 = \frac{1}{2} \|\mathbf{b}\|^2$ .

533 Then we can derive that the gradient for  $g(\mathbf{a}, \mathbf{b})$  w.r.t  $\mathbf{a}$  as follows:

$$g'(\mathbf{a}, \mathbf{b}) = \frac{\partial g(\mathbf{a}, \mathbf{b})}{\partial \mathbf{a}} = \frac{1}{2\pi} \left( \|\mathbf{b}\| \sin \angle(\mathbf{a}, \mathbf{b}) \frac{\mathbf{a}}{\|\mathbf{a}\|} + (\pi - \angle(\mathbf{a}, \mathbf{b})) \mathbf{b} \right). \quad (13)$$

534 Using this auxiliary function, we can rewrite the loss function in Eq. (1) as the following form:

$$\begin{aligned} L(\mathbf{W}) &= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0,1)} \left( \sum_{i=1}^m \sigma(\mathbf{w}_i^\top \mathbf{x}) - \sum_{i=1}^k \sigma(\mathbf{v}_i^\top \mathbf{x}) \right)^2 \\ &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m g(\mathbf{w}_i, \mathbf{w}_j) + \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k g(\mathbf{v}_i, \mathbf{v}_j) - \sum_{i=1}^m \sum_{j=1}^k g(\mathbf{w}_i, \mathbf{v}_j). \end{aligned} \quad (14)$$

535 Accordingly, when  $\mathbf{w}_i \neq \mathbf{0}$ . for  $\forall i \in [n]$ , the loss function is differentiable with gradient given by:

$$\begin{aligned}
\nabla_i &:= \frac{\partial L(\mathbf{W})}{\partial \mathbf{w}_i} \\
&= \sum_{j=1, j \neq i}^m \frac{\partial g(\mathbf{w}_i, \mathbf{w}_j)}{\partial \mathbf{w}_i} + \frac{1}{2} \frac{\partial g(\mathbf{w}_i, \mathbf{w}_i)}{\partial \mathbf{w}_i} - \sum_{l=1}^k \frac{\partial g(\mathbf{w}_i, \mathbf{v}_l)}{\partial \mathbf{w}_i} \\
&= \frac{\mathbf{w}_i}{2} + \frac{1}{2\pi} \sum_{j=1, j \neq i}^m \left( \|\mathbf{w}_j\| \sin \varphi_{ij} \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|} + (\pi - \varphi_{ij}) \mathbf{w}_j \right) - \frac{1}{2\pi} \sum_{l=1}^k \left( \|\mathbf{v}_l\| \sin \theta_{il} \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|} + (\pi - \theta_{il}) \mathbf{v}_l \right) \\
&= \frac{1}{2} \sum_{j=1}^m \mathbf{w}_j - \frac{1}{2} \sum_{l=1}^k \mathbf{v}_l + \frac{1}{2\pi} \left[ \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|} \left( \sum_{j=1, j \neq i}^m \sin \varphi_{ij} \|\mathbf{w}_j\| - \sum_{l=1}^k \sin \theta_{il} \|\mathbf{v}_l\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij} \mathbf{w}_j + \sum_{l=1}^k \theta_{il} \mathbf{v}_l \right].
\end{aligned} \tag{15}$$

536 The bounded norm  $\|\mathbf{w}_i\|$  at initialization can be given as below.

537 **Lemma 1** (Adapted from Lemma 3 in [Xu and Du, 2023]). *Let  $s_1 := \frac{1}{2}\sigma\sqrt{d}$ .  $s_2 := 2\sigma\sqrt{d}$ , if*  
538  *$d = \Omega(\log(m/\delta))$ , with probability at least  $1 - \delta$ , the following properties hold at the initialization:*

$$s_1 \leq \|\mathbf{w}_i(0)\| \leq s_2, \quad \forall i \in [m].$$

539 **Remark:** This is a standard fact in high-dimensional statistics, and on this basis, our result only  
540 involves this randomness. In the rest of the analysis in this paper is deterministic.

## 541 C Detailed analysis for Assumption 1

542 Here we prove the following lemma:

543 **Lemma 2.** *When  $d = \Omega(\frac{\log(mk/\delta)}{\zeta^2})$  with  $\zeta = o(\text{poly}(m^{-k^2}, k^{-k^2}))$ , then with probability at least*  
544  *$1 - \delta$ , the following property hold at the initialization:*

$$\frac{\pi}{2} - \zeta \leq \theta_{ij}(0) \leq \frac{\pi}{2} + \zeta. \quad \forall i \in [m], j \in [k].$$

545 *Proof.* According to Lemma 1, we have for  $\forall i \in [m], j \in [k]$ , we have:

$$|\langle \mathbf{w}_i(0), \bar{\mathbf{v}}_j \rangle| \leq \frac{\zeta}{4} \sigma \sqrt{d} \wedge \|\mathbf{w}_i(0)\| \geq \frac{1}{2} \sigma \sqrt{d} \Rightarrow |\cos \theta_{ij}(0)| \leq \frac{\zeta}{2} \Rightarrow \frac{\pi}{2} - \zeta \leq \theta_{ij}(0) \leq \frac{\pi}{2} + \zeta.$$

546 By concentration inequality of Gaussian, we have:

$$\mathbb{P} \left( |\langle \mathbf{w}_i(0), \bar{\mathbf{v}}_j \rangle| \geq \frac{\zeta}{4} \sigma \sqrt{d} \right) \leq 2 \exp \left( - \frac{(\frac{\zeta}{4} \sigma \sqrt{d})^2}{2\sigma^2} \right) \leq \frac{\delta}{3mk}.$$

547 Then:

$$\mathbb{P} \left( \theta_{ij}(0) \geq \frac{\pi}{2} + \zeta \vee \theta_{ij}(0) \leq \frac{\pi}{2} - \zeta \forall j \in [k] \right) \leq \frac{\delta}{3mk} * k + \frac{\delta}{3m} = \frac{2\delta}{3m}.$$

548 Applying the union bound for  $\forall i \in [m]$ , which finishes the proof.  $\square$

## 549 D Global Convergence: Phase 1 (Alignment)

550 In Phase 1, we are interested in the dynamics of  $\theta_{i^*}$  as well as the angle difference between the student  
551 neuron and its closest teacher neuron. The theorem we prove below is a combination of Theorem 3  
552 and Corollary 1 from the main text.

553 **Theorem 6** (Phase 1: Alignment). Assume  $d = \Omega(\log(m/\delta))$  with  $\delta \in (0, 1)$ , for any  $\epsilon_1 > 0$ , under  
554 Assumption 1 with  $10k\zeta \leq \epsilon_1^2 = o(\zeta_i^3)$  and Assumptions 2, 3 such that  $\sigma \leq \frac{\epsilon_1^{16k+2}\|\mathbf{v}\|}{10000m\sqrt{d}}$  in our random  
555 Gaussian initialization, and the stepsize satisfies  $\eta \leq \frac{\sigma\sqrt{d}\epsilon_1^2}{100k^2\|\mathbf{v}\|}$ , then by setting  $T_1 := \frac{\epsilon_1^2}{100\eta km}$ , for  
556  $0 \leq t \leq T_1$ , the following statements hold with probability at least  $1 - \delta$ :

$$s_1 \leq \|\mathbf{w}_i(t)\| \leq s_2 + 2k\eta\|\mathbf{v}\|t, \quad \forall i \in [m], \quad \text{with } s_1 := \frac{1}{2}\sigma\sqrt{d}, \quad s_2 := 2\sigma\sqrt{d}, \quad (16)$$

557 and

$$\sin^2\left(\frac{\theta_{i^*}(t)}{2}\right) - \epsilon_1^2 \leq \left(1 + \frac{\eta k\|\mathbf{v}\|t}{s_2}\right)^{-\frac{1}{8k}} \left(\sin^2\left(\frac{\theta_{i^*}(0)}{2}\right) - \epsilon_1^2\right), \quad \forall i \in [m]. \quad (17)$$

558 After Phase 1, we have:

$$\theta_{i^*}(T_1) \leq 4\epsilon_1, \quad \forall i \in [m]. \quad (18)$$

559 and

$$h_{i^*}(T_1) \leq 2h_{j^*}(T_1), \quad \forall i, j \in [m]. \quad (19)$$

560 *Proof.* The proof is given by induction. We firstly prove Eqs. (16) and (17) and then Eqs. (18)  
561 and (19).

562 At the initialization time  $t = 0$ , Eq. (16) and Eq. (17) directly hold according to Lemma 1. Note  
563 that the probability in this work only relates to the random initialization as given by Lemma 1. For  
564 description simplicity, we do not include this probability during the derivation but just mention it in  
565 our theorem.

566 Before proving Eqs. (16) and (17), we first analyse the learning dynamics of  $\theta_{i^*}$ . For any  $\forall i \in [m]$   
567 and  $0 < t < T_1$ , according to the inductive hypothesis, we have:

$$\sin^2\left(\frac{\theta_{i^*}(t)}{2}\right) \leq \max\left\{\sin^2\left(\frac{\theta_{i^*}(0)}{2}\right), \epsilon_1^2\right\} = \sin^2\left(\frac{\theta_{i^*}(0)}{2}\right) = \sin^2\left(\frac{\pi}{4} - \frac{\zeta_i}{2}\right),$$

568 where the right part of the above inequality is given by the following fact with Assumption 1:

$$\sin^2\left(\frac{\theta_{i^*}(0)}{2}\right) = \sin^2\left(\frac{\pi}{4} - \frac{\zeta_i}{2}\right) = \frac{1}{2} + \frac{1}{2}\sin(\zeta_i) = \Theta(1) \geq \epsilon_1^2 = o(\zeta_i^3) = o(1),$$

569 which means:

$$\theta_{i^*}(t) \leq \frac{\pi}{2} - \zeta_i. \quad (20)$$

570 Then we assume Eqs. (16) and (17) hold for any  $0 < t < T_1$  to prove Eqs. (16) and (17) for  $t + 1$ .

571 **Proof of right part of Eq. (16):**

572 According to the inductive hypothesis and  $s_2 := 2\sigma\sqrt{d}$  in Lemma 1, we have:

$$\|\mathbf{w}_i(t)\| \leq s_2 + 2k\eta\|\mathbf{v}\|T_1 \leq \frac{\epsilon_1^{16k+2}\|\mathbf{v}\|}{50m} + \frac{\epsilon_1^2\|\mathbf{v}\|}{50m} \leq \frac{\epsilon_1^2\|\mathbf{v}\|}{48m} = o\left(\frac{\|\mathbf{v}\|}{m}\right) \leq \frac{\|\mathbf{v}\|}{3m}, \quad \forall i \in [m]. \quad (21)$$

573 That means the teacher neuron's norm controls all of the student neurons' norm at  $t \in [0, T_1]$ . Then  
574 by triangle inequality and Eq. (21), the gradient norm can be upper bounded by



$$\begin{aligned}
& \|\nabla_i(t)\| \\
& \leq \left\| \frac{1}{2} \sum_{j=1}^m \mathbf{w}_j(t) \right\| + \left\| \frac{1}{2} \sum_{l=1}^k \mathbf{v}_l \right\| \\
& + \left\| \frac{1}{2\pi} \left[ \frac{\mathbf{w}_i(t)}{\|\mathbf{w}_i(t)\|} \left( \sum_{j=1, j \neq i}^m \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij}(t) \mathbf{w}_j(t) + \sum_{l=1}^k \theta_{il}(t) \mathbf{v}_l(t) \right] \right\| \\
& \leq \frac{m}{2} \times \frac{\|\mathbf{v}\|}{3m} + \frac{k}{2} \|\mathbf{v}\| + \frac{1}{2\pi} \left( m \times \frac{\|\mathbf{v}\|}{3m} + k \|\mathbf{v}\| + m\pi \times \frac{\|\mathbf{v}\|}{3m} + k\pi \|\mathbf{v}\| \right) \\
& < 2k \|\mathbf{v}\|, \quad \forall i \in [m].
\end{aligned} \tag{22}$$

575 One can see that, the gradient norm is upper bounded by all of the teacher neuron's norm. Accordingly,  
576 based on the gradient iteration, by the above results, we have:

$$\|\mathbf{w}_i(t+1)\| = \|\mathbf{w}_i(t) - \eta \nabla_i(t)\| \leq \|\mathbf{w}_i(t)\| + \|\eta \nabla_i(t)\| \leq s_2 + 2k\eta \|\mathbf{v}\| (t+1), \quad \forall i \in [m], \tag{23}$$

577 which concludes the proof.

578 **Proof of left part of Eq. (16):**

579 Here we need to prove the lower bound, we have:

$$\|\mathbf{w}_i(t+1)\| \geq \|\mathbf{w}_i(t)\| \geq s_1, \quad \forall i \in [m].$$

580 According to the gradient iteration:

$$\|\mathbf{w}_i(t+1)\|^2 - \|\mathbf{w}_i(t)\|^2 = \|\mathbf{w}_i(t) - \eta \nabla_i(t)\|^2 - \|\mathbf{w}_i(t)\|^2 = -2\eta \langle \mathbf{w}_i(t), \nabla_i(t) \rangle + \eta^2 \|\nabla_i(t)\|^2, \quad \forall i \in [m],$$

581 we only need to prove  $\forall i \in [m], \langle \mathbf{w}_i(t), \nabla_i(t) \rangle \leq 0$ . To be specific, we split  $\langle \mathbf{w}_i(t), \nabla_i(t) \rangle$  into two  
582 parts:

$$\begin{aligned}
& \langle \mathbf{w}_i(t), \nabla_i(t) \rangle \\
&= \left\langle \mathbf{w}_i(t), \frac{1}{2} \sum_{j=1}^m \mathbf{w}_j(t) - \frac{1}{2} \sum_{l=1}^k \mathbf{v}_l \right\rangle \\
&+ \left\langle \mathbf{w}_i(t), \frac{1}{2\pi} \left[ \frac{\mathbf{w}_i(t)}{\|\mathbf{w}_i(t)\|} \left( \sum_{j=1, j \neq i}^m \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij}(t) \mathbf{w}_j(t) + \sum_{l=1}^k \theta_{il}(t) \mathbf{v}_l \right] \right\rangle \\
&= \frac{1}{2} \sum_{j=1}^m \langle \mathbf{w}_i(t), \mathbf{w}_j(t) \rangle - \frac{1}{2} \sum_{l=1}^k \langle \mathbf{w}_i(t), \mathbf{v}_l \rangle \\
&+ \frac{1}{2\pi} \|\mathbf{w}_i(t)\| \sum_{j=1, j \neq i}^m \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \frac{1}{2\pi} \|\mathbf{w}_i(t)\| \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| \\
&- \frac{1}{2\pi} \sum_{j=1, j \neq i}^m \varphi_{ij}(t) \langle \mathbf{w}_i(t), \mathbf{w}_j(t) \rangle + \frac{1}{2\pi} \sum_{l=1}^k \theta_{il}(t) \langle \mathbf{w}_i(t), \mathbf{v}_l \rangle \\
&= \frac{1}{2} \|\mathbf{w}_i(t)\| \left( \sum_{j=1}^m \|\mathbf{w}_j(t)\| \cos \varphi_{ij}(t) - \sum_{l=1}^k \|\mathbf{v}\| \cos \theta_{il}(t) \right) \\
&+ \frac{1}{\pi} \sum_{j=1, j \neq i}^m \|\mathbf{w}_j(t)\| \sin \varphi_{ij}(t) - \frac{1}{\pi} \sum_{l=1}^k \|\mathbf{v}\| \sin \theta_{il}(t) \\
&- \frac{1}{\pi} \sum_{j=1, j \neq i}^m \|\mathbf{w}_j(t)\| \cos \varphi_{ij}(t) \varphi_{ij}(t) + \frac{1}{\pi} \sum_{l=1}^k \|\mathbf{v}\| \cos \theta_{il}(t) \theta_{il}(t) \\
&= \frac{1}{2\pi} \|\mathbf{w}_i(t)\| \|\mathbf{v}\| \sum_{l=1}^k \underbrace{\left( -\pi \cos \theta_{il}(t) - \sin \theta_{il}(t) + \cos \theta_{il}(t) \theta_{il}(t) \right)}_{I_1} \\
&+ \frac{1}{2\pi} \|\mathbf{w}_i(t)\| \sum_{j=1}^m \|\mathbf{w}_j(t)\| \underbrace{\left( \pi \cos \varphi_{ij}(t) + \sin \varphi_{ij}(t) - \cos \varphi_{ij}(t) \varphi_{ij}(t) \right)}_{\tilde{I}_1},
\end{aligned}$$

583 where the last equality holds by including the additional term related to  $\varphi_{ii} = 0$  for any  $i \in [m]$ .

584 One hand, for  $I_1$ , by Eq. (17),  $I_1$  is a monotonically increase function of  $\theta_{il}(t)$  on the interval  $[0, \pi]$ .

585 Then by Eq. (20), we have  $\theta_{il}(t) \leq \theta_{i^*}(t) + \frac{\pi}{2} \leq \pi - \zeta_i$ , which implies that:

$$\begin{aligned}
I_1 &= -\pi \cos \theta_{il}(t) - \sin \theta_{il}(t) + \cos \theta_{il}(t) \theta_{il}(t) \\
&\leq -\pi \cos(\pi - \zeta_i) - \sin(\pi - \zeta_i) + \cos(\pi - \zeta_i)(\pi - \zeta_i) \\
&= \zeta_i \cos(\zeta_i) - \sin(\zeta_i) \\
&\leq -\frac{\zeta_i^3}{4},
\end{aligned}$$

586 where the last inequality holds by the fact that  $\zeta_i \cos(\zeta_i) - \sin(\zeta_i) \leq -\frac{\zeta_i^3}{4}$  is always true on the  
587 interval  $[0, \frac{\pi}{2}]$ .

588 On the other hand, to estimate  $\tilde{I}_1$ , recall  $\|\mathbf{w}_i(t)\| \leq \frac{\epsilon_1^2 \|\mathbf{v}\|}{48m} = o(1)$  in Eq. (21) and the fact  $|\tilde{I}_1| \leq \pi$ ,  
589 we have:

$$\frac{1}{2\pi} \|\mathbf{w}_i(t)\| \sum_{j=1}^m \|\mathbf{w}_j(t)\| \tilde{I}_1 \leq \frac{1}{2\pi} \|\mathbf{w}_i(t)\| \sum_{j=1}^m \|\mathbf{w}_j(t)\| |\tilde{I}_1| \leq \frac{1}{96} \|\mathbf{w}_i(t)\| \|\mathbf{v}\| \epsilon_1^2.$$

590 Accordingly, combining the above derivation over  $I_1$  and  $\tilde{I}_1$ , we have:

$$\langle \mathbf{w}_i(t), \nabla_i(t) \rangle \leq \frac{1}{2\pi} \|\mathbf{w}_i(t)\| \|\mathbf{v}\| \left( -\frac{k\zeta_i^3}{4} + \Theta(\epsilon_1^2) \right) \leq 0,$$

591 due to  $\epsilon_1^2 = o(\zeta_i^3)$ , we conclude that  $\|\mathbf{w}_i(t+1)\| \geq \|\mathbf{w}_i(t)\| \geq s_1$  and finish the proof for Eq. (16).

592 **Proof of Eq. (17):**

593 We analyze the learning dynamics of  $\cos \theta_{i^*}$  by splitting it into two parts (first-order term and the  
594 second-order term) as follows:

$$\begin{aligned} & \cos \theta_{i^*}(t+1) - \cos \theta_{i^*}(t) \\ &= \frac{\langle \mathbf{w}_i(t+1), \mathbf{v}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\| \|\mathbf{v}\|} - \frac{\langle \mathbf{w}_i(t), \mathbf{v}_{\tau_i} \rangle}{\|\mathbf{w}_i(t)\| \|\mathbf{v}\|} \\ &= \frac{\|\mathbf{w}_i(t)\| \langle \mathbf{w}_i(t+1), \mathbf{v}_{\tau_i} \rangle - \|\mathbf{w}_i(t+1)\| \langle \mathbf{w}_i(t), \mathbf{v}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\| \|\mathbf{w}_i(t)\| \|\mathbf{v}\|} \\ &= \frac{\|\mathbf{w}_i(t)\| \langle \mathbf{w}_i(t) - \eta \nabla_i(t), \mathbf{v}_{\tau_i} \rangle - \|\mathbf{w}_i(t+1)\| \langle \mathbf{w}_i(t), \mathbf{v}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\| \|\mathbf{w}_i(t)\| \|\mathbf{v}\|} \\ &= \frac{(\|\mathbf{w}_i(t)\| - \|\mathbf{w}_i(t+1)\|) \langle \mathbf{w}_i(t), \mathbf{v}_{\tau_i} \rangle - \|\mathbf{w}_i(t)\| \langle \eta \nabla_i(t), \mathbf{v}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\| \|\mathbf{w}_i(t)\| \|\mathbf{v}\|} \\ &= \frac{\left( \frac{\|\mathbf{w}_i(t)\|^2 - \|\mathbf{w}_i(t+1)\|^2}{\|\mathbf{w}_i(t)\| + \|\mathbf{w}_i(t+1)\|} \right) \langle \mathbf{w}_i(t), \mathbf{v}_{\tau_i} \rangle - \|\mathbf{w}_i(t)\| \langle \eta \nabla_i(t), \mathbf{v}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\| \|\mathbf{w}_i(t)\| \|\mathbf{v}\|} \tag{24} \\ &= \frac{\left( \frac{2\eta \langle \mathbf{w}_i(t), \nabla_i(t) \rangle - \eta^2 \|\nabla_i(t)\|^2}{\|\mathbf{w}_i(t)\| + \|\mathbf{w}_i(t+1)\|} \right) \langle \mathbf{w}_i(t), \mathbf{v}_{\tau_i} \rangle - \eta \|\mathbf{w}_i(t)\| \langle \nabla_i(t), \mathbf{v}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\| \|\mathbf{w}_i(t)\| \|\mathbf{v}\|} \\ &= \underbrace{\frac{\eta}{\|\mathbf{w}_i(t+1)\|} \langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \nabla_i(t) \rangle}_{I_2} \\ &+ \underbrace{\frac{\eta \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\|} \left( \frac{\langle \bar{\mathbf{w}}_i(t), \nabla_i(t) \rangle (\|\mathbf{w}_i(t)\| - \|\mathbf{w}_i(t+1)\|) - \eta \|\nabla_i(t)\|^2}{\|\mathbf{w}_i(t+1)\| + \|\mathbf{w}_i(t)\|} \right)}_{I_3}. \end{aligned}$$

595 One can see that we need to estimate the respective two parts  $I_2$  and  $I_3$ . For term  $I_2$ , note that  
596  $\langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \mathbf{w}_i(t) \rangle = 0$ , then we have:

$$\begin{aligned}
I_2 &= \frac{\eta}{\|\mathbf{w}_i(t+1)\|} \langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \nabla_i(t) \rangle \\
&= \frac{\eta}{\|\mathbf{w}_i(t+1)\|} \left( \left\langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \frac{1}{2} \sum_{j=1}^m \mathbf{w}_j(t) - \frac{1}{2} \sum_{l=1}^k \mathbf{v}_l \right\rangle \right. \\
&\quad \left. + \left\langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \frac{1}{2\pi} \left[ - \sum_{j=1, j \neq i}^m \varphi_{ij}(t) \mathbf{w}_j(t) + \sum_{l=1}^k \theta_{il}(t) \mathbf{v}_l \right] \right\rangle \right) \\
&= \frac{\eta}{\|\mathbf{w}_i(t+1)\|} \left( \frac{1}{2} \left\langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \sum_{j=1}^m \mathbf{w}_j(t) \right\rangle - \frac{1}{2} \left\langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \sum_{l=1}^k \mathbf{v}_l \right\rangle \right. \\
&\quad \left. - \frac{1}{2\pi} \left\langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \sum_{j=1, j \neq i}^m \varphi_{ij}(t) \mathbf{w}_j(t) \right\rangle + \frac{1}{2\pi} \left\langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \sum_{l=1}^k \theta_{il}(t) \mathbf{v}_l \right\rangle \right) \\
&= \frac{\eta}{2\pi \|\mathbf{w}_i(t+1)\|} \left( \left\langle \cos \theta_{i^*}(t) \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \sum_{j=1}^m (\pi - \varphi_{ij}(t)) \mathbf{w}_j(t) \right\rangle - \left\langle \cos \theta_{i^*}(t) \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \sum_{l=1}^k (\pi - \theta_{il}(t)) \mathbf{v}_l \right\rangle \right) \\
&= \frac{\eta}{2\pi \|\mathbf{w}_i(t+1)\|} \left( \left\langle \cos \theta_{i^*}(t) \bar{\mathbf{w}}_i(t), \sum_{j=1}^m (\pi - \varphi_{ij}(t)) \mathbf{w}_j(t) \right\rangle - \left\langle \bar{\mathbf{v}}_{\tau_i}, \sum_{j=1}^m (\pi - \varphi_{ij}(t)) \mathbf{w}_j(t) \right\rangle \right. \\
&\quad \left. - \left\langle \cos \theta_{i^*}(t) \bar{\mathbf{w}}_i(t), \sum_{l=1}^k (\pi - \theta_{il}(t)) \mathbf{v}_l \right\rangle + \left\langle \bar{\mathbf{v}}_{\tau_i}, \sum_{l=1}^k (\pi - \theta_{il}(t)) \mathbf{v}_l \right\rangle \right) \\
&= \frac{\eta}{2\pi \|\mathbf{w}_i(t+1)\|} \left( \sum_{j=1}^m (\|\mathbf{w}_j(t)\| (\pi - \varphi_{ij}(t)) \cos \theta_{i^*}(t) \cos \varphi_{ij}(t)) - \sum_{j=1}^m (\|\mathbf{w}_j(t)\| (\pi - \varphi_{ij}(t)) \cos \theta_{j\tau_i}(t)) \right. \\
&\quad \left. - \sum_{l=1, l \neq \tau_i}^k (\|\mathbf{v}_l\| (\pi - \theta_{il}(t)) \cos \theta_{i^*}(t) \cos \theta_{il}(t)) + \|\mathbf{v}\| \sin^2 \theta_{i^*}(t) (\pi - \theta_{i^*}(t)) \right) \\
&\geq \frac{\eta \|\mathbf{v}\|}{2\pi \|\mathbf{w}_i(t+1)\|} \left( \sin^2 \theta_{i^*}(t) (\pi - \theta_{i^*}(t)) - \sum_{l=1, l \neq \tau_i}^k ((\pi - \theta_{il}(t)) \cos \theta_{i^*}(t) \cos \theta_{il}(t)) - \frac{\pi}{12} \epsilon_1^2 \right) \quad [\text{Eq. (21)}] \\
&\geq \frac{\eta \|\mathbf{v}\|}{2\pi \|\mathbf{w}_i(t+1)\|} \left( \frac{\pi}{2} \sin^2 \theta_{i^*}(t) - 2k\pi\zeta - \frac{\pi}{12} \epsilon_1^2 \right) \\
&\geq \frac{\eta \|\mathbf{v}\|}{4 \|\mathbf{w}_i(t+1)\|} \left( \sin^2 \theta_{i^*}(t) - \frac{17}{30} \epsilon_1^2 \right), \tag{25}
\end{aligned}$$

597 which builds the connection between  $I_2$  and  $\sin^2 \theta_{i^*}(t)$ . For term  $I_3$ :

$$\begin{aligned}
I_3 &= \frac{\eta \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\|} \left( \frac{\langle \bar{\mathbf{w}}_i(t), \nabla_i(t) \rangle (\|\mathbf{w}_i(t)\| - \|\mathbf{w}_i(t+1)\|) - \eta \|\nabla_i(t)\|^2}{\|\mathbf{w}_i(t+1)\| + \|\mathbf{w}_i(t)\|} \right) \\
&\geq -\frac{\eta}{\|\mathbf{w}_i(t+1)\|} \left( \frac{\|\nabla_i(t)\| \|\eta \nabla_i(t)\| + \eta \|\nabla_i(t)\|^2}{\|\mathbf{w}_i(t+1)\| + \|\mathbf{w}_i(t)\|} \right) \quad [\text{using Eq. (23)}] \\
&\geq -\frac{\eta}{\|\mathbf{w}_i(t+1)\|} \left( \frac{\|\nabla_i(t)\| \|\eta \nabla_i(t)\| + \eta \|\nabla_i(t)\|^2}{2s_1} \right) \quad [\text{using Eq. (16)}] \tag{26} \\
&= -\frac{\eta^2 \|\nabla_i(t)\|^2}{s_1 \|\mathbf{w}_i(t+1)\|} \\
&\geq -\frac{4k^2 \eta^2 \|\mathbf{v}\|^2}{s_1 \|\mathbf{w}_i(t+1)\|}. \quad [\text{using Eq. (22)}]
\end{aligned}$$

598 Take Eq. (25) and Eq. (26) into Eq. (24), we have:

$$\begin{aligned}
\cos \theta_{i^*}(t+1) - \cos \theta_{i^*}(t) &= I_2 + I_3 \\
&\geq \frac{\eta \|\mathbf{v}\|}{4 \|\mathbf{w}_i(t+1)\|} \left( \sin^2 \theta_{i^*}(t) - \frac{17}{30} \epsilon_1^2 - \frac{16k^2 \eta \|\mathbf{v}\|}{s_1} \right) \\
&\geq \frac{\eta \|\mathbf{v}\|}{4 \|\mathbf{w}_i(t+1)\|} \left( \sin^2 \theta_{i^*}(t) - \frac{133}{150} \epsilon_1^2 \right) \\
&\geq \frac{\eta \|\mathbf{v}\|}{4 \|\mathbf{w}_i(t+1)\|} \left( \sin^2 \theta_{i^*}(t) - \epsilon_1^2 \right).
\end{aligned}$$

599 Accordingly, we transform the dynamics analysis on  $\theta_{i^*}$  from cos to sin, which allows for estimat-  
600 ing Eq. (17) as below. Recall  $\cos 2x = 1 - 2 \sin^2 x$ , the above inequality implies:

$$\begin{aligned}
&\sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) - \sin^2 \left( \frac{\theta_{i^*}(t+1)}{2} \right) \\
&= \frac{\cos \theta_{i^*}(t+1) - \cos \theta_{i^*}(t)}{2} \\
&\geq \frac{\eta \|\mathbf{v}\|}{8 \|\mathbf{w}_i(t+1)\|} \left( \sin^2 \theta_{i^*}(t) - \epsilon_1^2 \right) \tag{27} \\
&= \frac{\eta \|\mathbf{v}\|}{8 \|\mathbf{w}_i(t+1)\|} \left( 4 \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) \cos^2 \left( \frac{\theta_{i^*}(t)}{2} \right) - \epsilon_1^2 \right) \\
&\geq \frac{\eta \|\mathbf{v}\|}{4 \|\mathbf{w}_i(t+1)\|} \left( \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) - \epsilon_1^2 \right),
\end{aligned}$$

601 which implies:

$$\begin{aligned}
\sin^2 \left( \frac{\theta_{i^*}(t+1)}{2} \right) - \epsilon_1^2 &\leq \left( 1 - \frac{\eta \|\mathbf{v}\|}{4 \|\mathbf{w}_i(t+1)\|} \right) \left( \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) - \epsilon_1^2 \right) \\
&\leq \left( 1 - \frac{\eta \|\mathbf{v}\|}{4(s_2 + 2\eta k \|\mathbf{v}\| (t+1))} \right) \left( \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) - \epsilon_1^2 \right) \quad [\text{using Eq. (16)}] \\
&\leq \prod_{u=1}^{t+1} \left( 1 - \frac{\eta \|\mathbf{v}\|}{4(s_2 + 2\eta k \|\mathbf{v}\| u)} \right) \left( \sin^2 \left( \frac{\theta_{i^*}(0)}{2} \right) - \epsilon_1^2 \right) \\
&\leq \exp \left( \int_{u=1}^{t+2} -\frac{\eta \|\mathbf{v}\|}{4(s_2 + 2\eta k \|\mathbf{v}\| u)} du \right) \left( \sin^2 \left( \frac{\theta_{i^*}(0)}{2} \right) - \epsilon_1^2 \right) \quad [\text{using } 1 - x \leq e^{-x}] \\
&= \exp \left( -\frac{1}{8k} \ln \left( \frac{s_2 + 2\eta k \|\mathbf{v}\| (t+2)}{s_2 + 2\eta k \|\mathbf{v}\|} \right) \right) \left( \sin^2 \left( \frac{\theta_{i^*}(0)}{2} \right) - \epsilon_1^2 \right) \\
&\leq \left( 1 + \frac{\eta k \|\mathbf{v}\| (t+1)}{s_2} \right)^{-\frac{1}{8k}} \left( \sin^2 \left( \frac{\theta_{i^*}(0)}{2} \right) - \epsilon_1^2 \right).
\end{aligned}$$

602 Accordingly, we finish the proof of Eq. (17).

603 **Proof of Eq. (18):**

604 Let  $t_0 := \frac{T}{50} \in \mathbb{N}$ , for any  $t \in [t_0, T_1]$ , using Eq. (17) and definitions of  $s_2, \sigma$ , we have:

$$\begin{aligned}
\sin^2\left(\frac{\theta_{i^*}(t)}{2}\right) - \epsilon_1^2 &\leq \left(1 + \frac{\eta k \|\mathbf{v}\| t}{s_2}\right)^{-\frac{1}{8k}} \left(\sin^2\left(\frac{\theta_{i^*}(0)}{2}\right) - \epsilon_1^2\right) \\
&\leq \left(1 + \frac{\eta k \|\mathbf{v}\| t}{s_2}\right)^{-\frac{1}{8k}} \\
&\leq \left(\frac{\eta k \|\mathbf{v}\| t_0}{s_2}\right)^{-\frac{1}{8k}} \\
&= \left(\frac{\eta k \|\mathbf{v}\| T_1}{100\sigma\sqrt{d}}\right)^{-\frac{1}{8k}} \\
&\leq \left(\frac{\|\mathbf{v}\| \epsilon_1^2}{10000m\sigma\sqrt{d}}\right)^{-\frac{1}{8k}} \\
&\leq \epsilon_1^2.
\end{aligned}$$

605 That means:  $\sin^2\left(\frac{\theta_{i^*}(t)}{2}\right) \leq 2\epsilon_1^2$ . So  $\forall t \in [T_1/50, T_1]$  and  $\forall i \in [m]$ , we have  $\theta_{i^*}(t) \leq 4\epsilon_1$ .

606 Consequently, each student neuron has aligned to a teacher neuron by the end of phase 1.

607 **Proof of Eq. (19):** For any  $t \in [T_1/50, T_1]$ , we study the dynamics of  $h_{i^*}$  (i.e., the inner product  
608 between the projection of gradient and teacher neuron) admitting the following formulation:

$$\begin{aligned}
&h_{i^*}(t+1) - h_{i^*}(t) \\
&= \langle \mathbf{w}_i(t+1), \bar{\mathbf{v}}_{\tau_i} \rangle - \langle \mathbf{w}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \\
&= -\eta \langle \nabla_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \\
&= -\frac{\eta}{2} \left\langle \sum_{j=1}^m \mathbf{w}_j(t) - \mathbf{v}_{\tau_i}, \bar{\mathbf{v}}_{\tau_i} \right\rangle \\
&\quad - \frac{\eta}{2\pi} \left\langle \frac{\mathbf{w}_i(t)}{\|\mathbf{w}_i(t)\|} \left( \sum_{j=1, j \neq i}^m \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij}(t) \mathbf{w}_j(t) + \theta_{i^*}(t) \mathbf{v}_{\tau_i}, \bar{\mathbf{v}}_{\tau_i} \right\rangle \\
&= \frac{\eta}{2} \left( \|\mathbf{v}\| - \sum_{j=1}^m h_{j\tau_i}(t) \right) \\
&\quad - \frac{\eta}{2\pi} \left( \cos \theta_{i^*}(t) \left( \sum_{j=1, j \neq i}^m \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij}(t) h_{j\tau_i}(t) + \theta_{i^*}(t) \|\mathbf{v}\| \right).
\end{aligned} \tag{28}$$

609 To analyse this dynamics, we need to study the  $\sin \theta_{il}(t)$  at first. According to Assumption 2, we  
610 have:

$$\frac{\pi}{2} - \theta_{i^*}(t) \leq \theta_{il}(t) \leq \frac{\pi}{2} + \theta_{i^*}(t), \quad \forall i \in [m], \tau_i \neq l \in [k].$$

611 So we have:

$$-\theta_{i^*}(t) \leq \frac{\pi}{2} - \theta_{il}(t) \leq \theta_{i^*}(t), \quad \forall i \in [m], \tau_i \neq l \in [k].$$

612 That is:

$$1 \geq \sin \theta_{il}(t) = \cos\left(\frac{\pi}{2} - \theta_{il}(t)\right) \geq \cos \theta_{i^*}(t) \geq \cos(4\epsilon_1) \geq 1 - 8\epsilon_1^2, \quad \forall i \in [m], \tau_i \neq l \in [k].$$

613 Then taking it back to Eq. (28), we have:

$$\begin{aligned}
h_{i^*}(t+1) - h_{i^*}(t) &\leq \frac{\eta}{2} \|\mathbf{v}\| - \frac{\eta}{2\pi} \left( \cos \theta_{i^*}(t) \left( - \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij}(t) h_{j\tau_i}(t) \right) \\
&\leq \frac{\eta}{2} \|\mathbf{v}\| + \frac{\eta}{2\pi} \left( \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| + \sum_{j=1, j \neq i}^m \varphi_{ij}(t) \|\mathbf{w}_j(t)\| \right) \text{ [using } \cos \theta_{i^*}(t) \leq 1] \\
&\leq \frac{\eta}{2} \|\mathbf{v}\| + \frac{\eta}{2\pi} \left( k \|\mathbf{v}\| + (m-1)\pi \frac{\epsilon_1^2 \|\mathbf{v}\|}{48m} \right) \text{ [using Eq. (21) and } \varphi_{ij} < \pi] \\
&\leq \frac{k + \pi - 0.5}{2\pi} \eta \|\mathbf{v}\|, \quad \forall i \in [m].
\end{aligned}$$

614 Similarly, we can derive that:

$$h_{i^*}(t+1) - h_{i^*}(t) \geq \frac{k + \pi - 1.5}{2\pi} \eta \|\mathbf{v}\|, \quad \forall i \in [m].$$

615 Then, we accumulate over the time:

$$\frac{49(k + \pi - 1.5)}{100\pi} \eta T_1 \|\mathbf{v}\| \leq h_{i^*}(T_1) - h_{i^*}\left(\frac{T_1}{50}\right) \leq \frac{49(k + \pi - 0.5)}{100\pi} \eta T_1 \|\mathbf{v}\|, \quad \forall i \in [m]. \quad (29)$$

616 The remaining thing left is to bound  $h_{i^*}\left(\frac{T_1}{50}\right)$ :

$$\left| h_{i^*}\left(\frac{T_1}{50}\right) \right| \leq \left\| \mathbf{w}_i\left(\frac{T_1}{50}\right) \right\| \leq s_2 + 2k\eta \|\mathbf{v}\| \frac{T_1}{50} \leq \frac{k}{20} \eta T_1 \|\mathbf{v}\|, \quad \forall i \in [m]. \quad (30)$$

617 Combine Eqs. (29) and (30), we have:

$$\frac{49k + 49\pi - 5\pi k - 73.5}{100\pi} \eta T_1 \|\mathbf{v}\| \leq h_{i^*}(T_1) \leq \frac{49k + 49\pi + 5\pi k - 24.5}{100\pi} \eta T_1 \|\mathbf{v}\|, \quad \forall i \in [m]. \quad (31)$$

618 Hence we finish the proof of Eq. (19).  $\square$

## 619 E Global Convergence: Phase 2 (Behaviors on the tangential growth)

620 In Phase 2, we are interested in the dynamics of  $h_i^*$  as well as the tangential difference between the  
621 student neuron and its closest teacher neuron.

### 622 E.1 Global Convergence: Phase 2 (Tangential growth process)

623 In this section, we will restate and prove Theorem 4.

624 **Theorem 7** (Phase 2: Tangential Growth, restate version of Theorem 4). *Assume  $d = \Omega(\log(m/\delta))$   
625 with  $\delta \in (0, 1)$ , for any  $\epsilon_1 > 0, \epsilon_2 > 0$ , under Assumption 1 with  $10k\zeta \leq \epsilon_1^2 = o(\zeta_i^3) =$   
626  $o(\epsilon_2^{\Theta(k)}/m), \epsilon_2 = o(1)$ , Assumptions 2, 3 such that  $\sigma \leq \frac{\epsilon_1^{16k+2} \|\mathbf{v}\|}{10000m\sqrt{d}}$  in our random Gaussian initial-  
627 ization, and the stepsize satisfies  $\eta = o\left(\frac{m\epsilon_1^2 s_1^2}{k^2 \|\mathbf{v}\|^2}\right) \leq \frac{\sigma \sqrt{d} \epsilon_1^2}{100k^2 \|\mathbf{v}\|}$ , then by setting  $T_1 := \frac{\epsilon_1^2}{100\eta km}$  and  
628  $T_2 = T_1 + \frac{k}{2\eta m} \ln\left(\frac{1}{48\pi\epsilon_2}\right)$ , then  $\forall T_1 \leq t \leq T_2$ , we define  $H_l(t) := \|\mathbf{v}\| - \sum_{i=1}^m \mathbb{I}_{\tau_i=l} h_{i^*}(t)$  for  
629  $l \in [k]$ , the following statements hold with probability at least  $1 - \delta$ :*

$$h_{i^*}(t) \leq 2h_{j^*}(t), \forall i, j \in [m] \text{ and } \tau_i = \tau_j. \quad (32)$$

$$\left(1 - \frac{\eta m}{9k}\right)^{t-T_1} \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\| \geq H_l(t) \geq \frac{2}{3} \|\mathbf{v}\| \left(1 - \frac{3\eta m}{2k}\right)^{t-T_1} - 8\pi\epsilon_2 \|\mathbf{v}\| \geq 24\pi\epsilon_2 \|\mathbf{v}\|, \forall l \in [k]. \quad (33)$$

$$\frac{2\|\mathbf{v}\|}{m\tau_i} \geq h_{i^*}(t) \geq \frac{s_1}{2}, \forall i \in [m], \quad (34)$$

630 *and*

$$\theta_{i^*}(t) \leq \epsilon_2, \forall i \in [m]. \quad (35)$$

631 *Proof.* We use induction to prove this theorem.

632 First, for  $t = T_1$ , according to Eq. (19) and Eq. (18), we have Eq. (32) and Eq. (35) hold directly.

633 For Eq. (33), by Eq. (21), we have:

$$\|\mathbf{v}\| \geq H_l(T_1) = \|\mathbf{v}\| - \sum_{i=1}^m \mathbb{I}_{\tau_i=l} h_{i^*}(t) \geq \frac{2}{3} \|\mathbf{v}\|, \forall l \in [k]. \quad (36)$$

634 For Eq. (34), for the left part, by Eq. (21) we have:

$$h_{i^*}(T_1) \leq \|\mathbf{w}_i(T_1)\| \leq \frac{2\|\mathbf{v}\|}{m} \leq \frac{2\|\mathbf{v}\|}{m_j}, \forall i \in [m],$$

635 *and for the right part, by Eq. (18) and Lemma 1, we have:*

$$h_{i^*}(T_1) = \|\mathbf{w}_i(T_1)\| \cos \theta_{i^*}(T_1) \geq (1 - 8\epsilon_1^2) \|\mathbf{w}_i(T_1)\| \geq \frac{s_1}{2}, \forall i \in [m].$$

636 Next step, we assume Eqs. (32) to (35) hold for  $T_1, T_1 + 1, \dots, t$  for any  $T_1 < t < T_2$ , and then  
637 prove Eqs. (32) to (35) for  $t + 1$ .

638 **Proof of Eq. (32):**

639 By Eq. (28), for any  $i \in [m]$ , we decompose the tangential difference  $h_{i^*}(t + 1) - h_{i^*}(t)$  as below:



$$\begin{aligned}
& h_{i^*}(t+1) - h_{i^*}(t) \\
&= \frac{\eta}{2} \left( \|\mathbf{v}\| - \sum_{j=1}^m h_{j\tau_i}(t) \right) \\
&\quad - \frac{\eta}{2\pi} \left( \cos \theta_{i^*}(t) \left( \sum_{j=1, j \neq i}^m \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij}(t) h_{j\tau_i}(t) + \theta_{i^*}(t) \|\mathbf{v}\| \right) \\
&= \frac{\eta}{2} \left( H_{\tau_i}(t) - \sum_{j=1}^m \mathbb{I}_{\tau_j \neq \tau_i} h_{j\tau_i}(t) \right) \\
&\quad - \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \sum_{l=1, l \neq \tau_i}^k \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=l} \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sin \theta_{il}(t) \|\mathbf{v}\| \right) \\
&\quad - \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sin \theta_{i^*}(t) \|\mathbf{v}\| \right) \\
&\quad + \frac{\eta}{2\pi} \left( \sum_{j=1, j \neq i}^m \varphi_{ij}(t) h_{j\tau_i}(t) + \theta_{i^*}(t) \|\mathbf{v}\| \right) \\
&= \frac{\eta}{2} \left( H_{\tau_i}(t) - \sum_{j=1}^m \mathbb{I}_{\tau_j \neq \tau_i} h_{j\tau_i}(t) \right) + \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \sum_{l=1, l \neq \tau_i}^k \left( \sin \theta_{il}(t) H_l(t) \right) \\
&\quad - \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \sum_{l=1, l \neq \tau_i}^k \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=l} \|\mathbf{w}_j(t)\| \left[ \sin \varphi_{ij}(t) - \cos \theta_{jl}(t) \sin \theta_{il}(t) \right] \right) \\
&\quad - \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sin \theta_{i^*}(t) \|\mathbf{v}\| \right) \\
&\quad + \frac{\eta}{2\pi} \left( \sum_{j=1, j \neq i}^m \varphi_{ij}(t) h_{j\tau_i}(t) + \theta_{i^*}(t) \|\mathbf{v}\| \right) \\
&= \frac{\eta}{2} H_{\tau_i}(t) + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_i}^k H_l(t) \\
&\quad - \frac{\eta}{2} \sum_{j=1}^m \mathbb{I}_{\tau_j \neq \tau_i} h_{j\tau_i}(t) + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_i}^k \left( \left[ \cos \theta_{i^*}(t) \sin \theta_{il}(t) - 1 \right] H_l(t) \right) \\
&\quad - \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \sum_{l=1, l \neq \tau_i}^k \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=l} \|\mathbf{w}_j(t)\| \left[ \sin \varphi_{ij}(t) - \cos \theta_{jl}(t) \sin \theta_{il}(t) \right] \right) \\
&\quad - \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sin \theta_{i^*}(t) \|\mathbf{v}\| \right) \\
&\quad + \frac{\eta}{2\pi} \left( \sum_{j=1}^m \mathbb{I}_{\tau_i \neq \tau_j} \varphi_{ij}(t) h_{j\tau_i}(t) + \sum_{j=1}^m \mathbb{I}_{\tau_i=\tau_j} \varphi_{ij}(t) h_{j\tau_i}(t) + \theta_{i^*}(t) \|\mathbf{v}\| \right) \\
&:= \frac{\eta}{2} H_{\tau_i}(t) + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_i}^k H_l(t) + Q_i(t),
\end{aligned} \tag{37}$$

640 where the  $Q_i(t)$  is defined as:

$$\begin{aligned}
Q_i(t) &:= -\frac{\eta}{2} \sum_{j=1}^m \mathbb{I}_{\tau_j \neq \tau_i} h_{j\tau_i}(t) + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_i}^k \left( \left[ \cos \theta_{i^*}(t) \sin \theta_{il}(t) - 1 \right] H_l(t) \right) \\
&\quad - \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \sum_{l=1, l \neq \tau_i}^k \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=l} \|\mathbf{w}_j(t)\| \left[ \sin \varphi_{ij}(t) - \cos \theta_{jl}(t) \sin \theta_{il}(t) \right] \right) \\
&\quad - \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sin \theta_{i^*}(t) \|\mathbf{v}\| \right) \\
&\quad + \frac{\eta}{2\pi} \left( \sum_{j=1}^m \mathbb{I}_{\tau_i \neq \tau_j} \varphi_{ij}(t) h_{j\tau_i}(t) + \sum_{j=1}^m \mathbb{I}_{\tau_i=\tau_j} \varphi_{ij}(t) h_{j\tau_i}(t) + \theta_{i^*}(t) \|\mathbf{v}\| \right).
\end{aligned}$$

641 To bound  $Q_i$ , we need to estimate  $\varphi_{ij}$  and  $\theta_{il}$  at first. By Eq. (35) and Assumption 2, we have that  
642 for  $\tau_j = l$  and  $\tau_i \neq l$ :

$$\frac{\pi}{2} - 2\epsilon_2 \leq \frac{\pi}{2} - \theta_{i^*}(t) - \theta_{j^*}(t) \leq \varphi_{ij}(t) \leq \frac{\pi}{2} + \theta_{i^*}(t) + \theta_{j^*}(t) \leq \frac{\pi}{2} + 2\epsilon_2.$$

643 And for a similar reason, we have:

$$\frac{\pi}{2} - \epsilon_2 \leq \theta_{il}(t) \leq \frac{\pi}{2} + \epsilon_2, \quad \text{and} \quad -\epsilon_2 \leq \theta_{jl}(t) \leq \epsilon_2,$$

644 which implies that for a sufficient small  $\epsilon_2$ :

$$\begin{aligned}
\sin \varphi_{ij}(t) - \cos \theta_{jl}(t) \sin \theta_{il}(t) &\leq |\sin \varphi_{ij}(t) - 1| + |1 - \cos \theta_{jl}(t) \sin \theta_{il}(t)| \\
&= \left( 1 - \cos \left( \frac{\pi}{2} - \varphi_{ij}(t) \right) \right) + \left( 1 - \cos \theta_{jl}(t) \cos \left( \frac{\pi}{2} - \theta_{il}(t) \right) \right) \\
&\cong (1 - \cos 2\epsilon_2) + (1 - \cos^2 \epsilon_2) \\
&\leq 2\epsilon_2^2 + \epsilon_2^2 \\
&= 3\epsilon_2^2.
\end{aligned}$$

645 Then using this result as well as Eqs. (34) and (35) to bound  $|Q_i(t)|$ , for  $\forall i \in [m]$ , we have:

$$\begin{aligned}
|Q_i(t)| &\leq \frac{\eta}{2} \sum_{j=1}^m \mathbb{I}_{\tau_j \neq \tau_i} \frac{2 \|\mathbf{v}\| \sin \theta_{j^*}(t)}{m_{\tau_j} \cos \theta_{j^*}(t)} + \frac{\eta}{2\pi} (k-1) \sin^2 \theta_{i^*}(t) \|\mathbf{v}\| \\
&\quad + \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \sum_{l=1, l \neq \tau_i}^k \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=l} \frac{2 \|\mathbf{v}\|}{m_{\tau_j} \cos \theta_{j^*}(t)} 3\epsilon_2^2 \right) \\
&\quad + \frac{\eta}{2\pi} \cos \theta_{i^*}(t) \sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} \frac{2 \|\mathbf{v}\| \sin \varphi_{ij}(t)}{m_{\tau_i} \cos \theta_{j^*}(t)} + \frac{\eta}{2\pi} \frac{\sin 2\theta_{i^*}(t)}{2} \|\mathbf{v}\| \\
&\quad + \frac{\eta}{2\pi} \sum_{j=1}^m \mathbb{I}_{\tau_i \neq \tau_j} \left( \frac{\pi}{2} + 2\epsilon_2 \right) \frac{2 \|\mathbf{v}\| \sin \theta_{j^*}(t)}{m_{\tau_j} \cos \theta_{j^*}(t)} + \frac{\eta}{2\pi} \sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} 2\epsilon_2 \frac{2 \|\mathbf{v}\|}{m_{\tau_j}} + \frac{\eta}{2\pi} \theta_{i^*}(t) \|\mathbf{v}\| \\
&\leq \eta \sum_{j=1}^m \mathbb{I}_{\tau_j \neq \tau_i} \frac{\|\mathbf{v}\| \epsilon_2 (1 + \epsilon_2^2)}{m_{\tau_j}} + \frac{\eta}{2\pi} k \epsilon_2^2 \|\mathbf{v}\| \\
&\quad + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_i}^k \left( \sum_{j=1}^m \mathbb{I}_{\tau_j=l} \frac{2 \|\mathbf{v}\| (1 + \epsilon_2^2)}{m_{\tau_j}} 3\epsilon_2^2 \right) \\
&\quad + \frac{\eta}{2\pi} \sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} \frac{2 \|\mathbf{v}\| 2\epsilon_2 (1 + \epsilon_2^2)}{m_{\tau_i}} + \frac{\eta}{2\pi} \frac{2\epsilon_2}{2} \|\mathbf{v}\| \\
&\quad + \frac{\eta}{2\pi} \sum_{j=1}^m \mathbb{I}_{\tau_i \neq \tau_j} \left( \frac{\pi}{2} + 2\epsilon_2 \right) \frac{2 \|\mathbf{v}\| \epsilon_2 (1 + \epsilon_2^2)}{m_{\tau_j}} + \frac{\eta}{2\pi} \sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} 2\epsilon_2 \frac{2 \|\mathbf{v}\|}{m_{\tau_j}} + \frac{\eta}{2\pi} \epsilon_2 \|\mathbf{v}\| \\
&\leq 1.1\eta k \epsilon_2 \|\mathbf{v}\| + \eta k \epsilon_2^2 \|\mathbf{v}\| + 2\eta k \epsilon_2^2 \|\mathbf{v}\| + 0.7\eta \epsilon_2 \|\mathbf{v}\| + 0.2\eta \epsilon_2 \|\mathbf{v}\| + 0.6\eta k \epsilon_2 \|\mathbf{v}\| + 0.7\eta \epsilon_2 \|\mathbf{v}\| + 0.2\eta \epsilon_2 \|\mathbf{v}\| \\
&\leq 4\eta k \epsilon_2 \|\mathbf{v}\| \\
&\leq \frac{1}{3} \left( \frac{\eta}{2} H_{\tau_i}(t) + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_i}^k H_l(t) \right),
\end{aligned} \tag{38}$$

646 where the last inequality use Eq. (33).

647 Then  $\forall i, j \in [m]$  and  $\tau_i = \tau_j$ , we have:

$$\begin{aligned}
h_{i^*}(t+1) &= h_{i^*}(t) + \frac{\eta}{2} H_{\tau_i}(t) + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_i}^k H_l(t) + Q_i(t) \\
&\leq 2h_{j^*}(t) + 2 \left( \frac{\eta}{2} H_{\tau_j}(t) + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_j}^k H_l(t) + Q_j(t) \right) \\
&\leq 2h_{j^*}(t+1),
\end{aligned}$$

648 which finishes the proof of Eq. (32).

649 **Proof of Eq. (33):**

650 Then we derive the dynamics of  $H_l(t)$ , for any  $l \in [k]$ , we have:

$$\begin{aligned}
H_l(t+1) &= H_l(t) - \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \left( h_{i^*}(t+1) - h_{i^*}(t) \right) \\
&= H_l(t) - \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \left( \frac{\eta}{2} H_{\tau_i}(t) + \frac{\eta}{2\pi} \sum_{j=1, j \neq l}^k H_j(t) + Q_i(t) \right) \\
&= \left( 1 - \frac{m_l \eta}{2} \right) H_l(t) - \frac{m_l \eta}{2\pi} \sum_{j=1, j \neq \tau_i}^k H_j(t) + \sum_{i=1}^m \mathbb{I}_{\tau_i=l} Q_i(t).
\end{aligned}$$

651 For ease of description, we write the recursive iteration in a matrix form

$$\mathbf{H}(t+1) = \left( \mathbf{I} - \frac{\eta}{2\pi} \mathbf{Diag}(m)(\mathbf{1}\mathbf{1}^\top + (\pi-1)\mathbf{I}) \right) \mathbf{H}(t) + \mathbf{Q}(t).$$

652 by defining the following quantities

$$\mathbf{H}(t) := [H_1(t), H_2(t), \dots, H_k(t)]^\top \in \mathbb{R}^k,$$

$$\mathbf{Diag}(m) := \mathbf{Diag}(m_1, m_2, \dots, m_k) \in \mathbb{R}^{k \times k},$$

$$\mathbf{Q}(t) := [\sum_{i=1}^m \mathbb{I}_{\tau_i=1} Q_i(t), \sum_{i=1}^m \mathbb{I}_{\tau_i=2} Q_i(t), \dots, \sum_{i=1}^m \mathbb{I}_{\tau_i=k} Q_i(t)]^\top \in \mathbb{R}^k.$$

653 In the next, we aim to derive the upper and lower bound of  $\mathbf{H}(t+1)$ . Denote  $\mathbf{A} :=$   
654  $[\frac{8\pi k \epsilon_2 \|\mathbf{v}\|}{\pi+k-1}, \frac{8\pi k \epsilon_2 \|\mathbf{v}\|}{\pi+k-1}, \dots, \frac{8\pi k \epsilon_2 \|\mathbf{v}\|}{\pi+k-1}]^\top \in \mathbb{R}^k$ , according to Eq. (38) and Assumption 3, we have:

$$\begin{aligned}
\mathbf{H}(t+1) - \mathbf{A} &\preceq \left( \mathbf{I} - \frac{\eta}{2\pi} \mathbf{Diag}(m)(\mathbf{1}\mathbf{1}^\top + (\pi-1)\mathbf{I}) \right) \mathbf{H}(t) + 4\eta k \epsilon_2 \|\mathbf{v}\| \mathbf{Diag}(m)\mathbf{1} - \mathbf{A}. \\
&= \left( \mathbf{I} - \frac{\eta}{2\pi} \mathbf{Diag}(m)(\mathbf{1}\mathbf{1}^\top + (\pi-1)\mathbf{I}) \right) (\mathbf{H}(t) - \mathbf{A}) \\
&\preceq \left( \mathbf{I} - \frac{\eta}{2\pi} \frac{m}{3k} (\pi-1)\mathbf{I} \right) \mathbf{H}(t) \\
&\preceq \left( 1 - \frac{\eta m (\pi-1)}{6\pi k} \right) \mathbf{H}(t).
\end{aligned}$$

655 Here  $\preceq$  means that all elements of the previous vector are smaller than the following vector. Then for  
656  $l \in [k]$ , we have:

$$\begin{aligned}
H_l(t+1) &\leq \left( 1 - \frac{\eta m (\pi-1)}{6\pi k} \right)^{t+1-T_1} H_l(T_1) + \frac{8\pi k \epsilon_2 \|\mathbf{v}\|}{\pi+k-1} \\
&\leq \left( 1 - \frac{\eta m (\pi-1)}{6\pi k} \right)^{t+1-T_1} \|\mathbf{v}\| + \frac{8\pi k \epsilon_2 \|\mathbf{v}\|}{\pi+k-1} \\
&\leq \left( 1 - \frac{\eta m}{9k} \right)^{t+1-T_1} \|\mathbf{v}\| + 8\pi \epsilon_2 \|\mathbf{v}\|.
\end{aligned}$$

657 Similarly, we have

$$\begin{aligned}
\mathbf{H}(t+1) + \mathbf{A} &\succeq \left( \mathbf{I} - \frac{\eta}{2\pi} \frac{3m}{k} (\mathbf{1}\mathbf{1}^\top + (\pi-1)\mathbf{I}) \right) (\mathbf{H}(t) + \mathbf{A}) \\
&\succeq \left( \mathbf{I} - \frac{3\eta m}{2\pi k} (\mathbf{1}\mathbf{1}^\top + (\pi-1)\mathbf{I}) \right) \mathbf{H}(t).
\end{aligned}$$

658 Here  $\succ$  means that all elements of the previous vector are greater than the following vector. The  
659 eigenvalues of matrix  $\mathbf{I} - \frac{3\eta m}{2\pi k}(\mathbf{1}\mathbf{1}^\top + (\pi - 1)\mathbf{I})$  is calculated to be one  $1 - \frac{3\eta m(k + \pi - 1)}{2\pi k}$  and the  
660 rest  $k - 1$  are  $1 - \frac{3\eta m(\pi - 1)}{2\pi k}$ . Then according to Eq. (36), for  $l \in [k]$ , we have:

$$\begin{aligned} H_l(t+1) &\geq \frac{2}{3} \|\mathbf{v}\| \left( 1 - \frac{3\eta m(k + \pi - 1 + (k-1)(\pi-1))}{2\pi k^2} \right)^{t+1-T_1} - 8\pi\epsilon_2 \|\mathbf{v}\| \\ &= \frac{2}{3} \|\mathbf{v}\| \left( 1 - \frac{3\eta m}{2k} \right)^{t+1-T_1} - 8\pi\epsilon_2 \|\mathbf{v}\|. \end{aligned}$$

661 Based on the above results, for  $l \in [k]$ , we have:

$$\left( 1 - \frac{\eta m}{9k} \right)^{t+1-T_1} \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\| \geq H_l(t) \geq \frac{2}{3} \|\mathbf{v}\| \left( 1 - \frac{3\eta m}{2k} \right)^{t+1-T_1} - 8\pi\epsilon_2 \|\mathbf{v}\|. \quad (39)$$

662 Due to  $\eta m \ll 1$ , we have  $(1-x) \geq \exp(-1.5x)$  with  $x := \eta m$ . Using this fact, for any  $t \leq T_2$ , the  
663 last inequality can be further lower bounded by:

$$\begin{aligned} &\frac{2}{3} \|\mathbf{v}\| \left( 1 - \frac{3\eta m}{2k} \right)^{t-T_1} - 8\pi\epsilon_2 \\ &\geq \frac{2}{3} \|\mathbf{v}\| \exp\left( -\frac{2\eta m}{k} \frac{k}{2\eta m} \ln\left( \frac{1}{48\pi\epsilon_2} \right) \right) - 8\pi\epsilon_2 \|\mathbf{v}\| \\ &= 24\pi\epsilon_2 \|\mathbf{v}\|. \end{aligned}$$

664 **Proof of Eq. (34):**

665 To prove the left part, by Eq. (33), we have:  $H_l(t+1) = \|\mathbf{v}\| - \sum_{i=1}^m \mathbb{I}_{\tau_i=l} h_{i^*}(t+1) \geq 0$ . Then  
666 we have:

$$\|\mathbf{v}\| \geq \sum_{i=1}^m \mathbb{I}_{\tau_i=l} h_{i^*}(t+1) \geq \frac{m_{\tau_i}}{2} h_{i^*}(t+1), \quad \forall i \in [m].$$

667 For the right part, we have:

$$h_{i^*}(t+1) - h_{i^*}(t) = \frac{\eta}{2} H_{\tau_i}(t) + \frac{\eta}{2\pi} \sum_{l=1, l \neq \tau_i}^k H_l(t) + Q_i(t) \geq \frac{\eta k}{2\pi} 24\pi\epsilon_2 \|\mathbf{v}\| - 4\eta k \epsilon_2 \|\mathbf{v}\| \geq 0.$$

668 So we have  $h_{i^*}(t+1) \geq h_{i^*}(t) \geq h_{i^*}(T_1) \geq \frac{s_1}{2}$ .

669 **Proof of Eq. (35):**

670 First, we prove that for  $\forall i, j \in [m]$ ,  $T_1 \leq t \leq T_2$ , we have  $\frac{\|\mathbf{w}_i(t)\|}{\|\mathbf{w}_j(t)\|} = \Theta(1)$ .

671 When  $t = T_1$ , according to Eq. (31) we have:

$$\frac{1}{2} \leq \frac{h_{i^*}(T_1)}{h_{j^*}(T_1)} \leq 2, \quad \forall i, j \in [m],$$

672 which implies:

$$\frac{\|\mathbf{w}_i(T_1)\|}{\|\mathbf{w}_j(T_1)\|} = \frac{h_{i^*}(T_1) \cos \theta_{i^*}(T_1)}{h_{j^*}(T_1) \cos \theta_{j^*}(T_1)} = \Theta(1), \quad \forall i, j \in [m]. \quad (40)$$

673 Then by defining  $t_s = \frac{9k \ln(2)}{\eta m} + T_1$ , when  $T_1 \leq t \leq t_s$ , according to Eq. (33), for any  $l \in [k]$ , we  
674 have:

$$\begin{aligned}
H_l(t) &\geq \frac{2}{3} \|\mathbf{v}\| \left(1 - \frac{3\eta m}{2k}\right)^{t-T_1} - 8\pi\epsilon_2 \|\mathbf{v}\| \\
&\geq \frac{2}{3} \|\mathbf{v}\| \exp\left(-\frac{2\eta m}{k} \frac{9k \ln(2)}{\eta m}\right) - 8\pi\epsilon_2 \|\mathbf{v}\| \\
&\geq \frac{2}{3} \left(\frac{1}{2}\right)^{18} \|\mathbf{v}\| - 8\pi\epsilon_2 \|\mathbf{v}\|.
\end{aligned}$$

675 So for  $\forall l_1, l_2 \in [k]$ , we have  $\frac{H_{l_1}(t)}{H_{l_2}(t)} = \Theta(1)$ .

676 Then for  $\forall i, j \in [m]$ , according to Eq. (37), for  $T_1 \leq t_0 < t$ , we have  $\frac{h_{i^*}(t_0+1)-h_{i^*}(t_0)}{h_{j^*}(t_0+1)-h_{j^*}(t_0)} = \Theta(1)$ .

677 Then consider Eq. (31), we have  $\frac{h_{i^*}(t)}{h_{j^*}(t)} = \Theta(1)$ .

678 That means for  $\forall i, j \in [m]$ , when  $T_1 \leq t \leq t_s$ , we have:

$$\frac{\|\mathbf{w}_i(t)\|}{\|\mathbf{w}_j(t)\|} = \frac{h_{i^*}(t) \cos \theta_{i^*}(t)}{h_{j^*}(t) \cos \theta_{j^*}(t)} = \Theta(1). \quad (41)$$

679 When  $t_s \leq t \leq T_2$ , according to Eq. (33), for  $\forall l \in [k]$ , we have:

$$\begin{aligned}
H_l(t) &\leq \left(1 - \frac{\eta m}{9k}\right)^{t-T_1} \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\| \\
&\leq \exp\left(-\frac{\eta m}{9k} \frac{9k \ln(2)}{\eta m}\right) \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\| \quad [\text{using } (1-x) \leq \exp(-x), \forall x \geq 0] \\
&= \frac{1}{2} \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\|.
\end{aligned}$$

680 Then we have:

$$\sum_{i=1}^m \mathbb{I}_{\tau_i=l} h_{i^*}(t) = \|\mathbf{v}\| - H_l(t) \geq \frac{1}{2} \|\mathbf{v}\| - 8\pi\epsilon_2 \|\mathbf{v}\| \geq \frac{1}{3} \|\mathbf{v}\|.$$

681 Then for  $\forall i, j \in [m]$ , we have:

$$\begin{aligned}
h_{i^*}(t) &\geq \frac{\sum_{l=i}^m \mathbb{I}_{\tau_l=i} h_{l^*}(t)}{2m_{\tau_i}} \\
&\geq \frac{\|\mathbf{v}\|}{6m_{\tau_i}} \\
&\geq \frac{\sum_{l=i}^m \mathbb{I}_{\tau_l=j} h_{l^*}(t)}{6m_{\tau_i}} \\
&\geq \frac{m_{\tau_j} h_{j^*}(t)}{12m_{\tau_i}} \\
&\geq \frac{h_{j^*}(t)}{108}.
\end{aligned}$$

682 That means for  $\forall i, j \in [m]$ , when  $t_s \leq t \leq T_2$ , we have:

$$\frac{\|\mathbf{w}_i(t)\|}{\|\mathbf{w}_j(t)\|} = \frac{h_{i^*}(t) \cos \theta_{i^*}(t)}{h_{j^*}(t) \cos \theta_{j^*}(t)} = \Theta(1). \quad (42)$$

683 So combine Eqs. (40) to (42), for  $\forall i, j \in [m]$ , when  $T_1 \leq t \leq T_2$ , we have:

$$\frac{\|\mathbf{w}_i(t)\|}{\|\mathbf{w}_j(t)\|} = \frac{h_{i^*}(t) \cos \theta_{i^*}(t)}{h_{j^*}(t) \cos \theta_{j^*}(t)} = \Theta(1). \quad (43)$$

684 Then, we analyze the change in angle, recall the dynamics of  $\cos \theta_{i^*}$  in Eq. (24) is given by:

$$\cos \theta_{i^*}(t+1) - \cos \theta_{i^*}(t) =: I_2 + I_3.$$

685 For  $I_2$ , we have:

$$\begin{aligned} I_2 &= \frac{\eta}{\|\mathbf{w}_i(t+1)\|} \langle \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle \bar{\mathbf{w}}_i(t) - \bar{\mathbf{v}}_{\tau_i}, \nabla_i(t) \rangle \\ &= \frac{\eta}{2\pi \|\mathbf{w}_i(t+1)\|} \left( \sum_{j=1}^m (\|\mathbf{w}_j(t)\| (\pi - \varphi_{ij}(t)) (\cos \theta_{i^*}(t) \cos \varphi_{ij}(t) - \cos \theta_{j\tau_i}(t))) \right. \\ &\quad \left. - \sum_{l=1, l \neq \tau_i}^k (\|\mathbf{v}\| (\pi - \theta_{il}(t)) \cos \theta_{i^*}(t) \cos \theta_{il}(t)) + \|\mathbf{v}\| \sin^2 \theta_{i^*}(t) (\pi - \theta_{i^*}(t)) \right). \end{aligned}$$

686 To bound  $I_2$ , we need handle  $\cos \theta_{i^*}(t) \cos \varphi_{ij}(t) - \cos \theta_{j\tau_i}(t)$  at first. For  $\tau_i \neq \tau_j$ , without loss of  
687 generality, we assume that:  $\bar{\mathbf{v}}_{\tau_i} = [1, 0, 0, \dots, 0]^\top \in \mathbb{R}^d$  and  $\bar{\mathbf{v}}_{\tau_j} = [0, 1, 0, 0, \dots, 0]^\top \in \mathbb{R}^d$ . Let  
688  $\bar{\mathbf{w}}_i = [w_{i1}, w_{i2}, \dots, w_{id}]^\top \in \mathbb{R}^d$  and  $\bar{\mathbf{w}}_j = [w_{j1}, w_{j2}, \dots, w_{jd}]^\top \in \mathbb{R}^d$ , then we have:

$$\begin{aligned} &\cos \theta_{i^*}(t) \cos \varphi_{ij}(t) - \cos \theta_{j\tau_i}(t) \\ &= \langle \bar{\mathbf{w}}_i, \bar{\mathbf{v}}_{\tau_i} \rangle \langle \bar{\mathbf{w}}_i, \bar{\mathbf{w}}_j \rangle - \langle \bar{\mathbf{w}}_j, \bar{\mathbf{v}}_{\tau_i} \rangle \\ &= w_{i1}(t) \sum_{l=1}^d w_{il}(t) w_{jl}(t) - w_{j1}(t) \\ &= w_{i1}(t) \left( w_{i1}(t) w_{j1}(t) + w_{i2}(t) w_{j2}(t) + \sum_{l=3}^d w_{il}(t) w_{jl}(t) \right) - w_{j1}(t) \\ &= w_{i1}(t) \left( w_{i2}(t) w_{j2}(t) + \sum_{l=3}^d w_{il}(t) w_{jl}(t) \right) - \sin^2 \theta_{i^*}(t) w_{j1}(t) \\ &\geq - \left| w_{i2}(t) w_{j2}(t) + \sum_{l=3}^d w_{il}(t) w_{jl}(t) \right| - \sin^2 \theta_{i^*}(t) |w_{j1}(t)| \\ &\geq - |w_{i2}(t) w_{j2}(t)| - \left| \sum_{l=3}^d w_{il}(t) w_{jl}(t) \right| - \sin^2 \theta_{i^*}(t) |w_{j1}(t)| \\ &\geq - |w_{i2}(t) w_{j2}(t)| - \left| \left( \sum_{l=3}^d w_{il}(t)^2 \right)^{\frac{1}{2}} \left( \sum_{l=3}^d w_{jl}(t)^2 \right)^{\frac{1}{2}} \right| - \sin^2 \theta_{i^*}(t) |w_{j1}(t)| \quad [\text{Cauchy-Schwarz inequality}] \\ &\geq - |w_{i2}(t) w_{j2}(t)| - \sin \theta_{i^*}(t) \sin \theta_{j^*}(t) - \sin^2 \theta_{i^*}(t) |w_{j1}(t)| \\ &\geq - \sin \theta_{i^*}(t) \sin \theta_{j^*}(t) - 2\zeta. \end{aligned}$$

689 For  $\tau_i = \tau_j$ , without loss of generality, we assume that:  $\bar{\mathbf{v}}_{\tau_i} = \bar{\mathbf{v}}_{\tau_j} = [1, 0, 0, \dots, 0]^\top \in \mathbb{R}^d$ . Then,  
690 we let  $\bar{\mathbf{w}}_i = [w_{i1}, w_{i2}, \dots, w_{id}]^\top \in \mathbb{R}^d$  and  $\bar{\mathbf{w}}_j = [w_{j1}, w_{j2}, \dots, w_{jd}]^\top \in \mathbb{R}^d$ . Then we have:

$$\begin{aligned}
& \cos \theta_{i^*}(t) \cos \varphi_{ij}(t) - \cos \theta_{j\tau_i}(t) \\
&= \langle \bar{\mathbf{w}}_i, \bar{\mathbf{v}}_{\tau_i} \rangle \langle \bar{\mathbf{w}}_i, \bar{\mathbf{w}}_j \rangle - \langle \bar{\mathbf{w}}_j, \bar{\mathbf{v}}_{\tau_i} \rangle \\
&= w_{i1}(t) \sum_{l=1}^d w_{il}(t) w_{jl}(t) - w_{j1}(t) \\
&= w_{i1}(t) \left( w_{i1}(t) w_{j1}(t) + \sum_{l=2}^d w_{il}(t) w_{jl}(t) \right) - w_{j1}(t) \\
&= w_{i1}(t) \left( \sum_{l=2}^d w_{il}(t) w_{jl}(t) \right) - \sin^2 \theta_{i^*}(t) w_{j1}(t) \\
&= \cos \theta_{i^*}(t) \left( \sum_{l=2}^d w_{il}(t) w_{jl}(t) \right) - \sin^2 \theta_{i^*}(t) \cos \theta_{j\tau_i}(t) \\
&\geq - \left( \sum_{l=2}^d w_{il}(t)^2 \right)^{\frac{1}{2}} \left( \sum_{l=2}^d w_{jl}(t)^2 \right)^{\frac{1}{2}} - \sin^2 \theta_{i^*}(t) \quad [\text{Cauchy-Schwarz inequality}] \\
&= - \sin \theta_{i^*}(t) \sin \theta_{j^*}(t) - \sin^2 \theta_{i^*}(t).
\end{aligned}$$

691 Then we have:

$$\begin{aligned}
I_2 &= \frac{\eta}{2\pi \|\mathbf{w}_i(t+1)\|} \left( \sum_{j=1}^m (\|\mathbf{w}_j(t)\| (\pi - \varphi_{ij}(t)) (\cos \theta_{i^*}(t) \cos \varphi_{ij}(t) - \cos \theta_{j\tau_i}(t))) \right. \\
&\quad \left. - \sum_{l=1, l \neq \tau_i}^k (\|\mathbf{v}\| (\pi - \theta_{il}(t)) \cos \theta_{i^*}(t) \cos \theta_{il}(t) + \|\mathbf{v}\| \sin^2 \theta_{i^*}(t) (\pi - \theta_{i^*}(t))) \right) \\
&\geq - \frac{\eta}{\|\mathbf{w}_i(t+1)\|} \left( \sum_{j=1}^m [\|\mathbf{w}_j(t)\| (2\zeta + \sin \theta_{i^*}(t) \sin \theta_{j^*}(t))] + \sum_{j=1}^m \mathbb{I}_{\tau_j = \tau_i} (\|\mathbf{w}_j(t)\| \sin^2 \theta_{i^*}(t)) \right. \\
&\quad \left. + (k-1) \|\mathbf{v}\| \pi \zeta - (\pi - \theta_{i^*}(t)) \|\mathbf{v}\| \sin^2 \theta_{i^*}(t) \right) \\
&\geq -C^* \eta \sin \theta_{i^*}(t) \sum_{j=1}^m \sin \theta_{j^*}(t) - \frac{6k\eta\zeta \|\mathbf{v}\|}{\|\mathbf{w}_i(t+1)\|} \quad [\text{Eq. (43)}] \\
&\geq -C^* \eta \sin \theta_{i^*}(t) \sum_{j=1}^m \sin \theta_{j^*}(t) - \frac{12k\eta\zeta \|\mathbf{v}\|}{s_1} \quad [\text{Eq. (34)}].
\end{aligned} \tag{44}$$

692 In the next, we aim to bound  $I_3$ , which requires the estimation of the gradient. Similar to Eq. (22),  
693 we have:



$$\begin{aligned}
& \|\nabla_i(t)\| \\
& \leq \left\| \frac{1}{2} \sum_{j=1}^m \mathbf{w}_j(t) \right\| + \left\| \frac{1}{2} \sum_{l=1}^k \mathbf{v}_l \right\| \\
& + \left\| \frac{1}{2\pi} \left[ \frac{\mathbf{w}_i(t)}{\|\mathbf{w}_i(t)\|} \left( \sum_{j=1, j \neq i}^m \sin \varphi_{ij}(t) \|\mathbf{w}_j(t)\| - \sum_{l=1}^k \sin \theta_{il}(t) \|\mathbf{v}\| \right) - \sum_{j=1, j \neq i}^m \varphi_{ij}(t) \mathbf{w}_j(t) + \sum_{l=1}^k \theta_{il}(t) \mathbf{v}_l(t) \right] \right\| \\
& \leq \frac{m}{2} \times \frac{9k \|\mathbf{v}\|}{m} + \frac{k}{2} \|\mathbf{v}\| + \frac{1}{2\pi} \left( m \times \frac{9k \|\mathbf{v}\|}{m} + k \|\mathbf{v}\| + m\pi \times \frac{9k \|\mathbf{v}\|}{m} + k\pi \|\mathbf{v}\| \right) \\
& < 15k \|\mathbf{v}\|.
\end{aligned} \tag{45}$$

694 Combining with this result, we can derive the lower bound for  $I_3$ :

$$\begin{aligned}
I_3 & = \frac{\eta \langle \bar{\mathbf{w}}_i(t), \bar{\mathbf{v}}_{\tau_i} \rangle}{\|\mathbf{w}_i(t+1)\|} \left( \frac{\langle \bar{\mathbf{w}}_i(t), \nabla_i(t) \rangle (\|\mathbf{w}_i(t)\| - \|\mathbf{w}_i(t+1)\|) - \eta \|\nabla_i(t)\|^2}{\|\mathbf{w}_i(t+1)\| + \|\mathbf{w}_i(t)\|} \right) \\
& \geq -\frac{\eta}{\|\mathbf{w}_i(t+1)\|} \left( \frac{\|\nabla_i(t)\| \|\eta \nabla_i(t)\| + \eta \|\nabla_i(t)\|^2}{s_1} \right) \\
& = -\frac{4\eta^2 \|\nabla_i(t)\|^2}{s_1^2} \\
& \geq -\frac{900k^2 \eta^2 \|\mathbf{v}\|^2}{s_1^2}.
\end{aligned} \tag{46}$$

695 Subsequently, we need to estimate the difference  $\sin^2 \left( \frac{\theta_{i^*}(t+1)}{2} \right) - \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right)$  for our final  
696 estimation for  $\sin \theta_{i^*}$ . Hence, similar to Eq. (27), combining Eq. (44), for  $\forall i \in [m]$ , we have:

$$\begin{aligned}
& \sin^2 \left( \frac{\theta_{i^*}(t+1)}{2} \right) - \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) \\
& = -\frac{1}{2} \left( \cos \theta_{i^*}(t+1) - \cos \theta_{i^*}(t) \right) \\
& \leq -\frac{1}{2} \left( -C^* \eta \sin \theta_{i^*}(t) \sum_{j=1}^m \sin \theta_{j^*}(t) - \frac{12k\eta\zeta \|\mathbf{v}\|}{s_1} - \frac{900k^2 \eta^2 \|\mathbf{v}\|^2}{s_1^2} \right) \quad [\text{using Eq. (44) and Eq. (46)}] \\
& \leq 2C^* \eta \sin \left( \frac{\theta_{i^*}(t)}{2} \right) \sum_{j=1}^m \sin \left( \frac{\theta_{j^*}(t)}{2} \right) + \frac{6k\eta\zeta \|\mathbf{v}\|}{s_1} + \frac{450k^2 \eta^2 \|\mathbf{v}\|^2}{s_1^2}.
\end{aligned}$$

697 Summing over all student neurons yields:

$$\begin{aligned}
& \sum_{i=1}^m \sin^2 \left( \frac{\theta_{i^*}(t+1)}{2} \right) - \sum_{i=1}^m \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) \\
& \leq \sum_{i=1}^m \left[ 2C^* \eta \sin \left( \frac{\theta_{i^*}(t)}{2} \right) \sum_{j=1}^m \sin \left( \frac{\theta_{j^*}(t)}{2} \right) + \frac{6k\zeta\eta \|\mathbf{v}\|}{s_1} + \frac{450k^2 \eta^2 \|\mathbf{v}\|^2}{s_1^2} \right] \\
& = 2C^* \eta \sum_{i=1}^m \sin \left( \frac{\theta_{i^*}(t)}{2} \right) \sum_{j=1}^m \sin \left( \frac{\theta_{j^*}(t)}{2} \right) + \frac{6km\zeta\eta \|\mathbf{v}\|}{s_1} + \frac{450k^2 m \eta^2 \|\mathbf{v}\|^2}{s_1^2} \\
& \leq 2C^* \eta m \sum_{i=1}^m \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) + \frac{6km\zeta\eta \|\mathbf{v}\|}{s_1} + \frac{450k^2 m \eta^2 \|\mathbf{v}\|^2}{s_1^2}. \quad [\text{using AM-GM inequality}]
\end{aligned} \tag{47}$$

698 Then we have:

$$\begin{aligned}
& \sum_{i=1}^m \sin^2 \left( \frac{\theta_{i^*}(t+1)}{2} \right) \\
& \leq \sum_{i=1}^m \sin^2 \left( \frac{\theta_{i^*}(t+1)}{2} \right) + \frac{3k\zeta \|\mathbf{v}\|}{C^* s_1} + \frac{225k^2 \eta \|\mathbf{v}\|^2}{C^* s_1^2} \\
& \leq (1 + 2C^* \eta m) \left( \sum_{i=1}^m \sin^2 \left( \frac{\theta_{i^*}(t)}{2} \right) + \frac{3k\zeta \|\mathbf{v}\|}{C^* s_1} + \frac{225k^2 \eta \|\mathbf{v}\|^2}{C^* s_1^2} \right) \quad [\text{Eq. (47)}] \\
& \leq (1 + 2C^* \eta m)^{t+1-T_1} \left( \sum_{i=1}^m \sin^2 \left( \frac{\theta_{i^*}(T_1)}{2} \right) + \frac{3k\zeta \|\mathbf{v}\|}{C^* s_1} + \frac{225k^2 \eta \|\mathbf{v}\|^2}{C^* s_1^2} \right) \\
& \leq (1 + 2C^* \eta m)^{t+1-T_1} 4m\epsilon_1^2 \quad [\text{by Assumption 1, choosing } \zeta = o\left(\frac{m\epsilon_1^2 s_1}{k \|\mathbf{v}\|}\right)] \\
& \leq \exp \left( 2C^* \eta m \frac{k}{2\eta m} \ln \left( \frac{1}{48\pi\epsilon_2} \right) \right) 4m\epsilon_1^2 \quad [\text{using } 1+x \leq \exp(x)] \\
& \leq \frac{4m\epsilon_1^2}{(48\pi\epsilon_2)^{C^* k}} \\
& \leq \frac{\epsilon_2^2}{16},
\end{aligned}$$

699 where the last inequality needs  $\epsilon_1^2 \leq \frac{(48\pi\epsilon_2)^{C^* k} \epsilon_2^2}{64m}$ .

700 Finally we finish the proof for Eq. (35), i.e.,

$$\theta_{i^*}(t+1) \leq \epsilon_2, \quad \forall i \in [m].$$

701 which finishes the proof.

702

□

## 703 E.2 Global Convergence: Phase 2 (Final state)

704 Here we prove the bounds on the student neurons and the loss function at the end of phase 2.

705 **Lemma 3** (Final state of Phase 2, restate version of Corollary 2). *Under the same conditions*  
706 *as Theorem 7, at time  $T_2$ , we have the following statements hold with probability at least  $1 - \delta$ :*

$$\frac{\|\mathbf{v}\|}{3m_{\tau_i}} \leq \|\mathbf{w}_i(T_2)\| \leq \frac{3\|\mathbf{v}\|}{m_{\tau_i}}, \quad \forall i \in [m],$$

707 and

$$L(\mathbf{W}(T_2)) \leq \frac{1}{2} k^2 \epsilon_2^{0.05} \|\mathbf{v}\|^2.$$

708 *Proof.* Firstly we derive the bound for the  $\|\mathbf{w}_i(T_2)\|$ . By Eq. (33) in Theorem 7, for any  $l \in [k]$ , we  
709 have:

$$\begin{aligned}
H_i(T_2) &\leq \left(1 - \frac{\eta m}{9k}\right)^{T_2 - T_1} \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\| \\
&\leq \exp\left(-\frac{\eta m}{9k}(T_2 - T_1)\right) \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\| \\
&= (48\pi\epsilon_2)^{\frac{1}{18}} \|\mathbf{v}\| + 8\pi\epsilon_2 \|\mathbf{v}\| \\
&\leq (49\pi\epsilon_2)^{\frac{1}{18}} \|\mathbf{v}\| \\
&\leq \frac{1}{3} \|\mathbf{v}\|,
\end{aligned}$$

710 which implies

$$\frac{2}{3} \|\mathbf{v}\| \leq \|\mathbf{v}\| - H_{\tau_i}(T_2) = \sum_{j=1}^m \mathbb{I}_{\tau_j = \tau_i} h_{j^*}(T_2) \leq 2m_{\tau_i} h_{i^*}(T_2), \quad \forall i \in [m].$$

711 So we have the lower bound  $\|\mathbf{w}_i(T_2)\| \geq h_{i^*}(T_2) \geq \frac{\|\mathbf{v}\|}{3m_{\tau_i}}$ . For the upper bound, for any  $i \in [m]$ ,  
712 we have  $H_{\tau_i}(T_2) \geq 0$ :

$$\|\mathbf{v}\| \geq \sum_{j=1}^m \mathbb{I}_{\tau_j = \tau_i} h_{j^*}(T_2) \geq \frac{1}{2} m_{\tau_i} h_{i^*}(T_2) = \frac{1}{2} m_{\tau_i} \|\mathbf{w}_i(T_2)\| \cos \theta_{i^*}(T_2) \geq \frac{1}{3} m_{\tau_i} \|\mathbf{w}_i(T_2)\|,$$

713 which implies  $\|\mathbf{w}_i(T_2)\| \leq \frac{3\|\mathbf{v}\|}{m_{\tau_i}}$  and the following estimation which is used for estimating the loss.

714 To be specific, for any  $l \in [k]$ , we have:

$$\sum_{i=1}^m \mathbb{I}_{\tau_i = l} \|\mathbf{w}_i(T_2)\| = \sum_{i=1}^m \mathbb{I}_{\tau_i = l} \frac{h_{i^*}(T_2)}{\cos \theta_{i^*}(T_2)} \leq (1 + \epsilon_2^2) \sum_{i=1}^m \mathbb{I}_{\tau_i = l} h_{i^*}(T_2) \leq (1 + \epsilon_2^2) \|\mathbf{v}\|,$$

715 and

$$\sum_{i=1}^m \mathbb{I}_{\tau_i = l} \|\mathbf{w}_i(T_2)\| \geq \sum_{i=1}^m \mathbb{I}_{\tau_i = l} h_{i^*}(T_2) \geq (1 - (49\pi\epsilon_2)^{\frac{1}{18}}) \|\mathbf{v}\| \geq (1 - \epsilon_2^{0.05}) \|\mathbf{v}\|.$$

716 Combine the lower and upper bound, we have:

$$(1 - \epsilon_2^{0.05}) \|\mathbf{v}\| \leq \sum_{i=1}^m \mathbb{I}_{\tau_i = l} \|\mathbf{w}_i(T_2)\| \leq (1 + \epsilon_2^2) \|\mathbf{v}\|. \quad (48)$$

717 Before we bound the loss, we need to analyze  $g(\mathbf{a}, \mathbf{b})$  defined in Eq. (12). If  $\angle(\mathbf{a}, \mathbf{b}) \leq 2\epsilon_2$  we have:

$$\frac{\pi - 2\epsilon_2}{2\pi} \|\mathbf{a}\| \|\mathbf{b}\| \leq g(\mathbf{a}, \mathbf{b}) = \frac{\|\mathbf{a}\| \|\mathbf{b}\|}{2\pi} \left( \sin \angle(\mathbf{a}, \mathbf{b}) + (\pi - \angle(\mathbf{a}, \mathbf{b})) \cos \angle(\mathbf{a}, \mathbf{b}) \right) \leq \frac{1}{2} \|\mathbf{a}\| \|\mathbf{b}\|, \quad (49)$$

718 Besides, if  $-2\epsilon_2 \leq \frac{\pi}{2} - \angle(\mathbf{a}, \mathbf{b}) \leq 2\epsilon_2$ , we have:

$$\frac{1 - 4\epsilon_2}{2\pi} \|\mathbf{a}\| \|\mathbf{b}\| \leq g(\mathbf{a}, \mathbf{b}) \leq \frac{1 + 4\epsilon_2}{2\pi} \|\mathbf{a}\| \|\mathbf{b}\|. \quad (50)$$

719 According to Eq. (35) in Theorem 7, then when  $\tau_i = \tau_j$ , we have  $\varphi_{ij} \leq 2\epsilon_2$  and when  $\tau_i \neq \tau_j$ , we  
720 have  $-2\epsilon_2 \leq \frac{\pi}{2} - \varphi_{ij} \leq 2\epsilon_2$ .

721 Then, according to Eqs. (48) to (50), we have:

$$\begin{aligned}
L(\mathbf{W}(T_2)) &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m g(\mathbf{w}_i(T_2), \mathbf{w}_j(T_2)) + \frac{1}{2} \sum_{i=1}^k \sum_{j=1}^k g(\mathbf{v}_i, \mathbf{v}_j) - \sum_{i=1}^m \sum_{j=1}^k g(\mathbf{w}_i(T_2), \mathbf{v}_j) \\
&\leq \frac{1}{2} \sum_{l=1}^k \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \sum_{j=1}^m \mathbb{I}_{\tau_j=l} \frac{1}{2} \|\mathbf{w}_i(T_2)\| \|\mathbf{w}_j(T_2)\| + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \mathbb{I}_{\tau_i \neq \tau_j} \frac{1+4\epsilon_2}{2\pi} \|\mathbf{w}_i(T_2)\| \|\mathbf{w}_j(T_2)\| \\
&\quad + \frac{k}{2} \frac{\|\mathbf{v}\|^2}{2} + \frac{k(k-1)}{2} \frac{\|\mathbf{v}\|^2}{2\pi} \\
&\quad - \sum_{l=1}^k \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \frac{\pi-2\epsilon_2}{2\pi} \|\mathbf{w}_i(T_2)\| \|\mathbf{v}_l\| - \sum_{l=1}^k \mathbb{I}_{\tau_i \neq l} \frac{1-4\epsilon_2}{2\pi} \|\mathbf{w}_i(T_2)\| \|\mathbf{v}_l\| \\
&= \frac{1}{2} \sum_{l=1}^k \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \sum_{j=1}^m \mathbb{I}_{\tau_j=l} \frac{1}{2} \|\mathbf{w}_i(T_2)\| \|\mathbf{w}_j(T_2)\| + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \mathbb{I}_{\tau_i \neq \tau_j} \frac{1}{2\pi} \|\mathbf{w}_i(T_2)\| \|\mathbf{w}_j(T_2)\| \\
&\quad + \frac{k}{4} \frac{\|\mathbf{v}\|^2}{4} + \frac{k(k-1)}{4\pi} \frac{\|\mathbf{v}\|^2}{4\pi} - \sum_{l=1}^k \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \frac{1}{2} \|\mathbf{w}_i(T_2)\| \|\mathbf{v}_l\| - \sum_{l=1}^k \sum_{i=1}^m \mathbb{I}_{\tau_i \neq l} \frac{1}{2\pi} \|\mathbf{w}_i(T_2)\| \|\mathbf{v}_l\| \\
&\quad + \sum_{i=1}^m \sum_{j=1}^m \mathbb{I}_{\tau_i \neq \tau_j} \frac{\epsilon_2}{\pi} \|\mathbf{w}_i(T_2)\| \|\mathbf{w}_j(T_2)\| + \sum_{l=1}^k \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \frac{\epsilon_2}{\pi} \|\mathbf{w}_i(T_2)\| \|\mathbf{v}_l\| \\
&\quad + \sum_{l=1}^k \sum_{i=1}^m \mathbb{I}_{\tau_i \neq l} \frac{2\epsilon_2}{\pi} \|\mathbf{w}_i(T_2)\| \|\mathbf{v}_l\| \quad [\text{Eqs. (49) and (50)}] \\
&\leq \frac{k(1+\epsilon_2^2)^2 \|\mathbf{v}\|^2}{4} + \frac{k(k-1)(1+\epsilon_2^2)^2 \|\mathbf{v}\|^2}{4\pi} + \frac{k}{4} \frac{\|\mathbf{v}\|^2}{4} + \frac{k(k-1)}{4\pi} \frac{\|\mathbf{v}\|^2}{4\pi} \\
&\quad - \frac{k(1-\epsilon_2^{0.05}) \|\mathbf{v}\|^2}{2} - \frac{k(k-1)(1-\epsilon_2^{0.05}) \|\mathbf{v}\|^2}{2\pi} \\
&\quad + \frac{k(k-1)(1+\epsilon_2^2)^2 \epsilon_2 \|\mathbf{v}\|^2}{\pi} + \frac{k(1+\epsilon_2^2) \epsilon_2 \|\mathbf{v}\|^2}{\pi} + \frac{2k(k-1)(1+\epsilon_2^2) \epsilon_2 \|\mathbf{v}\|^2}{\pi} \quad [\text{Eq. (48)}] \\
&\leq \frac{1}{2} k^2 \epsilon_2^{0.05} \|\mathbf{v}\|^2,
\end{aligned} \tag{51}$$

722 which concludes the proof.  $\square$

## 723 F Global Convergence: Phase 3 (local convergence)

724 In phase 3, we focus on the local convergence of the network when the loss function has an upper  
725 bound. First, we introduce some structural lemmas related to the loss function of neural network.

### 726 F.1 Structural Lemmas

727 **Lemma 4.** We define that  $\mathbf{w}_i^* := \frac{h_{i^*}}{\sum_{j=1}^m \mathbb{I}_{\tau_j=\tau_i} h_{j^*}} \mathbf{v}_{\tau_i}$ , and  $\theta_{\max} := \max_{i \in [m]} \theta_{i^*}$ , then we have:

$$\sum_{i=1}^m \left\langle \frac{\partial}{\partial \mathbf{w}_i} L(\mathbf{W}), \mathbf{w}_i - \mathbf{w}_i^* \right\rangle \geq 2L(\mathbf{W}) - \mathcal{O}(k\theta_{\max}^2 \sum_{l=1}^k \|\mathbf{r}_l\| \|\mathbf{v}\|).$$

728 *Proof.* First, we decomposes the residual function  $R(\mathbf{x})$  into two terms:

$$\begin{aligned}
R(\mathbf{x}) &:= \sum_{i=1}^m \sigma(\mathbf{w}_i^\top \mathbf{x}) - \sum_{l=1}^k \sigma(\mathbf{v}_l^\top \mathbf{x}) \\
&= \sum_{i=1}^m (\mathbf{w}_i^\top \mathbf{x}) \sigma'(\mathbf{w}_i^\top \mathbf{x}) - \sum_{l=1}^k (\mathbf{v}_l^\top \mathbf{x}) \sigma'(\mathbf{v}_l^\top \mathbf{x}) \quad [\text{using ReLU property: } \sigma(x) = x \sigma'(x)] \\
&= \sum_{i=1}^m (\mathbf{w}_i^\top \mathbf{x}) \sigma'(\mathbf{w}_i^\top \mathbf{x}) - \sum_{l=1}^k \left( \left( \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \mathbf{w}_i - \mathbf{r}_l \right)^\top \mathbf{x} \right) \sigma'(\mathbf{v}_l^\top \mathbf{x}) \quad [\text{using definition of } \mathbf{r}_l] \\
&= \sum_{i=1}^m (\mathbf{w}_i^\top \mathbf{x}) \sigma'(\mathbf{w}_i^\top \mathbf{x}) - \sum_{l=1}^k \left( \left( \sum_{i=1}^m \mathbb{I}_{\tau_i=l} \mathbf{w}_i \right)^\top \mathbf{x} \right) \sigma'(\mathbf{v}_l^\top \mathbf{x}) + \sum_{l=1}^k (\mathbf{r}_l^\top \mathbf{x}) \sigma'(\mathbf{v}_l^\top \mathbf{x}) \\
&:= \underbrace{\sum_{l=1}^k \sum_{i=1}^m \mathbb{I}_{\tau_i=l} (\mathbf{w}_i^\top \mathbf{x}) \left( \sigma'(\mathbf{w}_i^\top \mathbf{x}) - \sigma'(\mathbf{v}_l^\top \mathbf{x}) \right)}_{R_1(\mathbf{x})} + \underbrace{\sum_{l=1}^k (\mathbf{r}_l^\top \mathbf{x}) \sigma'(\mathbf{v}_l^\top \mathbf{x})}_{R_2(\mathbf{x})}.
\end{aligned}$$

729 Then we can derive the lower bound for  $\sum_{i=1}^m \left\langle \frac{\partial}{\partial \mathbf{w}_i} L(\mathbf{W}), \mathbf{w}_i - \mathbf{w}_i^* \right\rangle$  that:

$$\begin{aligned}
&\sum_{i=1}^m \left\langle \frac{\partial}{\partial \mathbf{w}_i} L(\mathbf{W}), \mathbf{w}_i - \mathbf{w}_i^* \right\rangle \\
&= \sum_{i=1}^m \mathbb{E}_{\mathbf{x}} \left( R(\mathbf{x}) \sigma'(\mathbf{w}_i^\top \mathbf{x}) \mathbf{x}^\top (\mathbf{w}_i - \mathbf{w}_i^*) \right) \\
&= \mathbb{E}_{\mathbf{x}} \left[ R(\mathbf{x}) \sum_{i=1}^m \left( \sigma(\mathbf{w}_i^\top \mathbf{x}) - \sigma'(\mathbf{w}_i^\top \mathbf{x}) \mathbf{x}^\top \mathbf{w}_i^* \right) \right] \\
&= 2L(\mathbf{W}) + \mathbb{E}_{\mathbf{x}} \left[ R(\mathbf{x}) \left( \sum_{i=1}^m \left( \sigma(\mathbf{w}_i^\top \mathbf{x}) - \sigma'(\mathbf{w}_i^\top \mathbf{x}) \mathbf{x}^\top \mathbf{w}_i^* \right) - R(\mathbf{x}) \right) \right] \quad [\text{using } L(\mathbf{W}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} R(\mathbf{x})^2] \\
&= 2L(\mathbf{W}) + \mathbb{E}_{\mathbf{x}} \left[ R(\mathbf{x}) \left( \sum_{l=1}^k \sigma(\mathbf{v}_l^\top \mathbf{x}) - \sum_{i=1}^m \left( \sigma'(\mathbf{w}_i^\top \mathbf{x}) \mathbf{x}^\top \mathbf{w}_i^* \right) \right) \right] \quad [\text{using definition of } R(\mathbf{x})] \\
&= 2L(\mathbf{W}) + \mathbb{E}_{\mathbf{x}} \left[ R(\mathbf{x}) \left( \sum_{i=1}^m \left( \sigma'(\mathbf{x}^\top \mathbf{w}_i^*) \mathbf{x}^\top \mathbf{w}_i^* \right) - \sum_{i=1}^m \left( \sigma'(\mathbf{w}_i^\top \mathbf{x}) \mathbf{x}^\top \mathbf{w}_i^* \right) \right) \right] \quad [\text{using definition of } R(\mathbf{w}_i^*)] \\
&= 2L(\mathbf{W}) + \mathbb{E}_{\mathbf{x}} \left[ R(\mathbf{x}) \sum_{i=1}^m (\mathbf{x}^\top \mathbf{w}_i^*) \left( \sigma'(\mathbf{x}^\top \mathbf{w}_i^*) - \sigma'(\mathbf{w}_i^\top \mathbf{x}) \right) \right] \\
&:= 2L(\mathbf{W}) + \underbrace{\mathbb{E}_{\mathbf{x}} \left[ R_1(\mathbf{x}) \sum_{i=1}^m (\mathbf{x}^\top \mathbf{w}_i^*) \left( \sigma'(\mathbf{x}^\top \mathbf{w}_i^*) - \sigma'(\mathbf{w}_i^\top \mathbf{x}) \right) \right]}_{I_4} \\
&\quad + \underbrace{\mathbb{E}_{\mathbf{x}} \left[ R_2(\mathbf{x}) \sum_{i=1}^m (\mathbf{x}^\top \mathbf{w}_i^*) \left( \sigma'(\mathbf{x}^\top \mathbf{w}_i^*) - \sigma'(\mathbf{w}_i^\top \mathbf{x}) \right) \right]}_{I_5}.
\end{aligned}$$

730 For term  $I_4$ , note that for  $\forall i \in [m]$ , when  $\mathbf{w}_i^\top \mathbf{x} \geq 0$ , we have  $\sigma'(\mathbf{w}_i^\top \mathbf{x}) = 1$ , which means  
731  $\sigma'(\mathbf{w}_i^\top \mathbf{x}) - \sigma'(\mathbf{v}_l^\top \mathbf{x}) \geq 0$ . Then we have  $R_1(\mathbf{x}) \geq 0$ . Similar, we have  $\sum_{i=1}^m (\mathbf{x}^\top \mathbf{w}_i^*) \left( \sigma'(\mathbf{x}^\top \mathbf{w}_i^*) - \right.$   
732  $\left. \sigma'(\mathbf{w}_i^\top \mathbf{x}) \right) \geq 0$ . So we have  $I_4 \geq 0$ .

733 For term  $I_5$ , we have:

$$\begin{aligned}
I_5 &= \mathbb{E}_{\mathbf{x}} \sum_{l=1}^k (\mathbf{r}_l^\top \mathbf{x}) \sigma'(\mathbf{v}_l^\top \mathbf{x}) \sum_{i=1}^m (\mathbf{x}^\top \mathbf{w}_i^*) \left( \sigma'(\mathbf{x}^\top \mathbf{w}_i^*) - \sigma'(\mathbf{w}_i^\top \mathbf{x}) \right) \\
&= \sum_{l=1}^k \sum_{i=1}^m \mathbb{E}_{\mathbf{x}} (\mathbf{r}_l^\top \mathbf{x}) \sigma'(\mathbf{v}_l^\top \mathbf{x}) (\mathbf{x}^\top \mathbf{w}_i^*) \left( \sigma'(\mathbf{x}^\top \mathbf{w}_i^*) - \sigma'(\mathbf{w}_i^\top \mathbf{x}) \right) \\
&\geq - \sum_{l=1}^k \sum_{i=1}^m \mathcal{O}(\|\mathbf{r}_l\| \theta_{i^*}^2 \|\mathbf{w}_i^*\|),
\end{aligned}$$

734 where the last inequality is from the proof of Xu and Du [2023, Lemma 8].

735 Thus we have:

$$\begin{aligned}
\sum_{i=1}^m \left\langle \frac{\partial}{\partial \mathbf{w}_i} L(\mathbf{W}), \mathbf{w}_i - \mathbf{w}_i^* \right\rangle &= 2L(\mathbf{W}) + I_4 + I_5 \\
&\geq 2L(\mathbf{W}) - \sum_{l=1}^k \sum_{i=1}^m \mathcal{O}(\|\mathbf{r}_l\| \theta_{i^*}^2 \|\mathbf{w}_i^*\|) \\
&\geq 2L(\mathbf{W}) - \mathcal{O}(k \theta_{\max}^2 \sum_{l=1}^k \|\mathbf{r}_l\| \|\mathbf{v}\|),
\end{aligned}$$

736 which finishes the proof.  $\square$

737 **Lemma 5** (Bounds of  $\theta_{i^*}$  and  $\|\mathbf{r}\|$ ). *Given that  $\frac{\|\mathbf{v}\|}{3m\tau_i} \leq \|\mathbf{w}_i\| \leq \frac{3\|\mathbf{v}\|}{m\tau_i}$  and  $L(\mathbf{W}) = o(\|\mathbf{v}\|^2 k^{10})$ ,*  
738 *then we have:*

$$\|\mathbf{r}_l\| \leq \mathcal{O}(k^{\frac{11}{4}} \|\mathbf{v}\|^{\frac{1}{4}} L^{\frac{3}{8}}(\mathbf{W})), \quad \forall l \in [l]. \quad (52)$$

$$\|\mathbf{v}\|^2 \theta_{i^*}^3 = \Theta(k^3 L(\mathbf{W})), \quad \forall i \in [m]. \quad (53)$$

739 *Proof.* The proof technique here heavily depends on [Zhou et al., 2021], so we simplify our proof  
740 here. To be specific, using the same proof method as [Zhou et al., 2021, Lemma C.6], we have:

$$\sum_{i=1}^m \|\mathbf{w}_i\|^2 \theta_{i^*}^2 = \mathcal{O}(L^{\frac{1}{2}}(\mathbf{W})).$$

741 Similarly, following [Zhou et al., 2021, Lemma 12], we have:

$$\mathbb{E}_{\mathbf{x}} R_1(\mathbf{x})^2 = \mathcal{O}\left(k^{\frac{5}{2}} \|\mathbf{v}\|^{\frac{1}{2}} L^{\frac{3}{4}}(\mathbf{W})\right)$$

742 Based on Zhou et al. [2021, Lemma 11], we can derive that:

$$\mathbb{E}_{\mathbf{x}} R_2(\mathbf{x})^2 = \Omega\left(\frac{\|\mathbf{r}_l\|^2}{k^3}\right), \quad \forall l \in [k].$$

743 Combine the previous results, for any  $l \in [k]$ , the upper bound of  $\|\mathbf{r}_l\|$  is:

$$\begin{aligned}
\frac{\|\mathbf{r}_l\|}{k^{\frac{3}{2}}} &= \mathcal{O}(\mathbb{E}_{\mathbf{x}} R_2(\mathbf{x})) \\
&\leq \mathcal{O}(\mathbb{E}_{\mathbf{x}} R(\mathbf{x}) + \mathbb{E}_{\mathbf{x}} R_1(\mathbf{x})) \\
&= \mathcal{O}\left(L^{\frac{1}{2}}(\mathbf{W}) + k^{\frac{5}{4}} \|\mathbf{v}\|^{\frac{1}{4}} L^{\frac{3}{8}}(\mathbf{W})\right) \\
&\leq \mathcal{O}\left(k^{\frac{5}{4}} \|\mathbf{v}\|^{\frac{1}{4}} L^{\frac{3}{8}}(\mathbf{W})\right) \quad [\text{using } L(\mathbf{W}) = \mathcal{O}(k^{10} \|\mathbf{v}\|^2)].
\end{aligned}$$

744 Accordingly, we finish the proof of Eq. (52). Based on this, using the same proof method as Zhou  
745 et al. [2021, Lemma 9], we can directly obtain Eq. (53).  $\square$

746 **Lemma 6** (Bound of  $\|\mathbf{w}_i - \mathbf{w}_i^*\|$ ). *Given that  $\frac{\|\mathbf{v}\|}{3m\tau_i} \leq \|\mathbf{w}_i\| \leq \frac{3\|\mathbf{v}\|}{m\tau_i}$  and  $L(\mathbf{W}) = o(\frac{\|\mathbf{v}\|^2}{k^{\frac{22}{3}}})$ , then*  
747 *for  $\forall i \in [m]$ , we have:*

$$\|\mathbf{w}_i - \mathbf{w}_i^*\| \leq \mathcal{O}\left(\frac{k^{\frac{2}{3}} m^{\frac{2}{3}} L^{\frac{1}{3}}(\mathbf{W})}{\|\mathbf{v}\|^{\frac{2}{3}}}\right) \|\mathbf{w}_i\|.$$

748 *Proof.* By Lemma 5, we have  $\theta_{i^*} = \mathcal{O}\left(\frac{kL^{\frac{1}{3}}(\mathbf{W})}{\|\mathbf{v}\|^{\frac{2}{3}}}\right)$  and  $|H_l| = |\langle \mathbf{r}_l, \bar{\mathbf{v}}_l \rangle| = \|\mathbf{r}_l\| \leq$   
749  $\mathcal{O}(k^{\frac{11}{4}} \|\mathbf{v}\|^{\frac{1}{4}} L^{\frac{3}{8}}(\mathbf{W})) = o(\|\mathbf{v}\|)$ . Then we have:

$$\begin{aligned}
\|\mathbf{w}_i - \mathbf{w}_i^*\| &\leq \|\mathbf{w}_i - h_{i^*} \bar{\mathbf{v}}_{\tau_i}\| + \|h_{i^*} \bar{\mathbf{v}}_{\tau_i} - \mathbf{w}_i^*\| \\
&= \|\mathbf{w}_i - h_{i^*} \bar{\mathbf{v}}_{\tau_i}\| + \left| h_{i^*} \left(1 - \frac{\|\mathbf{v}\|}{\sum_{j=1}^m \mathbb{I}_{\tau_j = \tau_i} h_{i^*}}\right) \right| \\
&= \|\mathbf{w}_i\| \sin \theta_{i^*} + \frac{h_{i^*} |H_l|}{\|\mathbf{v}\| - |H_l|} \quad [\text{using definition of } H_l] \\
&\leq \|\mathbf{w}_i\| \mathcal{O}\left(\frac{kL^{\frac{1}{3}}(\mathbf{W})}{\|\mathbf{v}\|^{\frac{2}{3}}}\right) + \frac{\|\mathbf{w}_i\| \mathcal{O}(k^{\frac{11}{4}} \|\mathbf{v}\|^{\frac{1}{4}} L^{\frac{3}{8}}(\mathbf{W}))}{\|\mathbf{v}\|} \\
&\leq \mathcal{O}\left(\frac{kL^{\frac{1}{3}}(\mathbf{W})}{\|\mathbf{v}\|^{\frac{2}{3}}}\right) \|\mathbf{w}_i\| + \mathcal{O}\left(\frac{k^{\frac{11}{4}} L^{\frac{3}{8}}(\mathbf{W})}{\|\mathbf{v}\|^{\frac{3}{4}}}\right) \|\mathbf{w}_i\| \\
&\leq \mathcal{O}\left(\frac{kL^{\frac{1}{3}}(\mathbf{W})}{\|\mathbf{v}\|^{\frac{2}{3}}}\right) \|\mathbf{w}_i\| \quad [\text{using } L(\mathbf{W}) = \mathcal{O}(\frac{\|\mathbf{v}\|^2}{k^{\frac{7}{2}}})].
\end{aligned}$$

750  $\square$

## 751 F.2 Gradient Lower Bound

752 In this subsection, we use the structural lemmas in Appendix F.1 to derive the local gradient lower  
753 bound.

754 **Theorem 8.** *Given that  $\frac{\|\mathbf{v}\|}{3m\tau_i} \leq \|\mathbf{w}_i\| \leq \frac{3\|\mathbf{v}\|}{m\tau_i}$  for  $\forall i \in [m]$  and  $L(\mathbf{W}) = o(\frac{\|\mathbf{v}\|^2}{k^{162}})$ , then we have:*

$$\left\| \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} \right\| \geq \Omega\left(\frac{L^{\frac{2}{3}}(\mathbf{W})}{k^2 \|\mathbf{v}\|^{\frac{1}{3}}}\right).$$

755

756 *Proof.* According to Lemmas 4 and 5, we have:

$$\begin{aligned}
\sum_{i=1}^m \left\langle \frac{\partial}{\partial \mathbf{w}_i} L(\mathbf{W}), \mathbf{w}_i - \mathbf{w}_i^* \right\rangle &\geq 2L(\mathbf{W}) - \mathcal{O}(k\theta_{\max}^2 \sum_{l=1}^k \|\mathbf{r}_l\| \|\mathbf{v}\|) \\
&\geq 2L(\mathbf{W}) - \mathcal{O}\left(\left(\frac{kL^{\frac{1}{3}}(\mathbf{W})}{\|\mathbf{v}\|^{\frac{2}{3}}}\right)^2 k^2 (k^{\frac{11}{4}} \|\mathbf{v}\|^{\frac{1}{4}} L^{\frac{3}{8}}(\mathbf{W})) \|\mathbf{v}\|\right) \\
&\geq 2L(\mathbf{W}) - \mathcal{O}\left(\frac{L^{\frac{25}{24}}(\mathbf{W}) k^{\frac{27}{4}}}{\|\mathbf{v}\|^{\frac{1}{12}}}\right) \\
&\geq L(\mathbf{W}) \quad [\text{using } L(\mathbf{W}) = \mathcal{O}(\frac{\|\mathbf{v}\|^2}{k^{162}})].
\end{aligned}$$

757 Then according to Lemma 6, we have:

$$\begin{aligned}
L(\mathbf{W}) &\leq \sum_{i=1}^m \left\langle \frac{\partial}{\partial \mathbf{w}_i} L(\mathbf{W}), \mathbf{w}_i - \mathbf{w}_i^* \right\rangle \\
&\leq \sum_{i=1}^m \left\| \frac{\partial}{\partial \mathbf{w}_i} L(\mathbf{W}) \right\| \|\mathbf{w}_i - \mathbf{w}_i^*\| \\
&\leq \left\| \frac{\partial}{\partial \mathbf{W}} L(\mathbf{W}) \right\| \mathcal{O}\left(\frac{kL^{\frac{1}{3}}(\mathbf{W})}{\|\mathbf{v}\|^{\frac{2}{3}}}\right) \sum_{i=1}^m \|\mathbf{w}_i\| \\
&= \mathcal{O}\left(k^2 L^{\frac{1}{3}}(\mathbf{W}) \|\mathbf{v}\|^{\frac{1}{3}}\right) \left\| \frac{\partial L(\mathbf{W})}{\partial \mathbf{W}} \right\|,
\end{aligned}$$

758 which concludes the proof.

759

□

### 760 F.3 Local Conditional Smoothness of Loss

761 In this subsection, we deal with the non-smoothness of  $L$ . We will prove the smoothness of  $L$  when  
762 the student neuron has upper and lower bounds

763 **Lemma 7** (Local Conditional Smoothness of  $L$ ). *Given that  $\frac{\|\mathbf{v}\|}{5m\tau_i} \leq \|\mathbf{w}_i\| \leq \frac{5\|\mathbf{v}\|}{m\tau_i}$  for any  $i \in [m]$ ,  
764 define the Hessian matrix of  $L$  as  $\mathbf{\Lambda} = \frac{\partial^2 L(\mathbf{W})}{\partial \mathbf{W}^2}$ , then we have  $\|\mathbf{\Lambda}\|_2 \leq \mathcal{O}(m^2)$ .*

765 *Proof.* According to Safran et al. [2021], we have that  $L$  is twice differentiable and the closed-form  
766 expression of Hessian  $\mathbf{\Lambda} = \frac{\partial^2 L(\mathbf{W})}{\partial \mathbf{W}^2} \in \mathbb{R}^{md \times md}$  can be write as:

$$\mathbf{\Lambda} = \begin{pmatrix} \mathbf{\Lambda}_{1,1} & \cdots & \mathbf{\Lambda}_{1,m} \\ \vdots & \ddots & \vdots \\ \mathbf{\Lambda}_{m,1} & \cdots & \mathbf{\Lambda}_{1,m} \end{pmatrix},$$

767 where  $\mathbf{\Lambda}_{i,j} \in \mathbb{R}^{d \times d}$ ,  $\forall i, j \in [m]$ , we will discuss below.

768 For diagonal elements:

$$\mathbf{\Lambda}_{i,i} = \frac{1}{2} \mathbf{I} + \sum_{j=1, j \neq i}^m \mathbf{\Lambda}_1(\mathbf{w}_i, \mathbf{w}_j) - \sum_{l=1}^k \mathbf{\Lambda}_1(\mathbf{w}_i, \mathbf{v}_l), \quad \forall i \in [m],$$

769 and by defining  $\mathbf{n}_{\mathbf{w}, \mathbf{v}} = \bar{\mathbf{v}} - \cos \angle(\mathbf{w}, \mathbf{v}) \bar{\mathbf{w}}$ ,  $\mathbf{\Lambda}_1$  can be rewritten as:

$$\mathbf{\Lambda}_1(\mathbf{w}, \mathbf{v}) = \frac{\sin \angle(\mathbf{w}, \mathbf{v}) \|\mathbf{v}\|}{2\pi \|\mathbf{w}\|} \left( \mathbf{I} - \bar{\mathbf{w}} \bar{\mathbf{w}}^\top + \bar{\mathbf{n}}_{\mathbf{w}, \mathbf{v}} \bar{\mathbf{n}}_{\mathbf{w}, \mathbf{v}}^\top \right).$$



770 We can bound that

$$\|\mathbf{\Lambda}_1(\mathbf{w}, \mathbf{v})\| \leq \frac{\|\mathbf{v}\|}{\|\mathbf{w}\|}.$$

771 Then we have:

$$\begin{aligned} \|\mathbf{\Lambda}_{i,i}\| &\leq \left\| \frac{1}{2} \mathbf{I} \right\| + \sum_{j=1, j \neq i}^m \|\mathbf{\Lambda}_1(\mathbf{w}_i, \mathbf{w}_j)\| + \sum_{l=1}^k \|\mathbf{\Lambda}_1(\mathbf{w}_i, \mathbf{v}_l)\| \\ &= \mathcal{O}(1) + m\mathcal{O}(1) + k\mathcal{O}\left(\frac{m}{k}\right) \\ &= \mathcal{O}(m), \quad \forall i \in [m]. \end{aligned}$$

772 And non-diagonal elements satisfy that:

$$\mathbf{\Lambda}_{i,j} = \frac{1}{2\pi} \left( (\pi - \angle(\mathbf{w}_i, \mathbf{w}_j)) \mathbf{I} + \bar{\mathbf{n}}_{\mathbf{w}_i, \mathbf{w}_j} \bar{\mathbf{w}}_j^\top + \bar{\mathbf{n}}_{\mathbf{w}_j, \mathbf{w}_i} \bar{\mathbf{w}}_i^\top \right), \quad \forall i, j \in [m], \text{ and } i \neq j.$$

773 So we have:

$$\|\mathbf{\Lambda}_{i,j}\| \leq \frac{1}{2\pi} (\pi + 1 + 1) \leq 1, \quad \forall i, j \in [m], \text{ and } i \neq j.$$

774 Combining the above results, we have:

$$\|\mathbf{\Lambda}\| \leq \sum_{i=1}^m \sum_{j=1}^m \|\mathbf{\Lambda}_{i,j}\| \leq m(m-1) + m\mathcal{O}(m) = \mathcal{O}(m^2).$$

775

□

#### 776 F.4 Generalization Error Bound

777 In this subsection, we prove the final convergence result, which is also the generalization error bound.

778 **Theorem 9.** *Suppose the initial condition in Lemma 1 and Assumption 1 2 and 3 holds. If we set*  
 779  $\epsilon_2 = o(m^{-60}k^{-100})$  *and*  $\eta = o(\frac{1}{m})$  *in Theorem 7, then*  $\forall T \in \mathbb{N}$ , *we have the following statements*  
 780 *hold with probability at least*  $1 - \delta$ :

$$L(\mathbf{W}(T + T_2)) \leq \frac{1}{\left( L(\mathbf{W}(T_2))^{-\frac{1}{3}} + \Omega\left(k^{-4} \|\mathbf{v}\|^{-\frac{2}{3}}\right) \eta T \right)^3}, \quad (54)$$

781 and

$$\frac{\|\mathbf{v}\|}{4m_{\tau_i}} \leq \|\mathbf{w}_i(T + T_2)\| \leq \frac{4\|\mathbf{v}\|}{m_{\tau_i}} \quad \forall i \in [m]. \quad (55)$$

782 *Proof.* We prove Eqs. (54) and (55) together inductively.

783 For  $T = 0$ , Eq. (54) directly hold and by Lemma 3 we have Eq. (55) holds.

784 Then we assume Eqs. (54) and (55) hold for  $0, 1, \dots, t$  for any  $0 < t < T_1$  to prove Eqs. (54) and (55) for  $t + 1$ .

786 **Proof of Eq. (54):**

787 For  $\forall i \in [m]$ , similar to Eq. (45), we have  $\|\nabla_i(t)\| = \mathcal{O}(k \|\mathbf{v}\|)$ . Then for  $\forall t \in [0, 1]$ , we have:

$$\|\mathbf{w}_i(t) - \iota\eta\nabla_i(t)\| \geq \|\mathbf{w}_i(t)\| - \eta\|\nabla_i(t)\| \geq \frac{\|\mathbf{v}\|}{4m_{\tau_i}} - \eta\mathcal{O}(k\|\mathbf{v}\|) \geq \frac{\|\mathbf{v}\|}{5m_{\tau_i}},$$

788 and

$$\|\mathbf{w}_i(t) - \iota\eta\nabla_i(t)\| \leq \|\mathbf{w}_i(t)\| + \eta\|\nabla_i(t)\| \leq \frac{4\|\mathbf{v}\|}{m_{\tau_i}} + \eta\mathcal{O}(k\|\mathbf{v}\|) \leq \frac{5\|\mathbf{v}\|}{m_{\tau_i}}.$$

789 Then, we can use Lemma 7 for  $\mathbf{W}(t) - \iota\eta\nabla_{\mathbf{W}}(t)$  in the following proof.

790 For  $T_2 \leq t \leq T + T_2 - 1$ , according to the classic analysis of gradient descent in Nesterov et al.  
791 [2018], we have:

$$\begin{aligned} L(\mathbf{W}(t+1)) &= L(\mathbf{W}(t)) + \langle \nabla_{\mathbf{W}}(t), -\eta\nabla_{\mathbf{W}}(t) \rangle \\ &\quad + \int_{\iota=0}^1 (1-\iota)(-\eta\nabla_{\mathbf{W}}(t))^\top \frac{\partial^2 L}{\partial \mathbf{W}^2}(\mathbf{W}(t) - \iota\eta\nabla_{\mathbf{W}}(t))(-\eta\nabla_{\mathbf{W}}(t))d\iota \\ &\leq L(\mathbf{W}(t)) - \eta\|\nabla_{\mathbf{W}}(t)\|^2 + \int_{\iota=0}^1 (1-\iota)\eta^2\|\nabla_{\mathbf{W}}(t)\|^2\mathcal{O}(m^2)d\iota \quad [\text{Lemma 7}]. \end{aligned}$$

792 Then we have:

$$\begin{aligned} L(\mathbf{W}(t)) - L(\mathbf{W}(t+1)) &\geq \eta\|\nabla_{\mathbf{W}}(t)\|^2 - \int_{\iota=0}^1 (1-\iota)\eta^2\|\nabla_{\mathbf{W}}(t)\|^2\mathcal{O}(m^2)d\iota \\ &= \eta\|\nabla_{\mathbf{W}}(t)\|^2 - \frac{1}{2}\eta^2\|\nabla_{\mathbf{W}}(t)\|^2\mathcal{O}(m^2) \\ &\geq \frac{1}{2}\eta\|\nabla_{\mathbf{W}}(t)\|^2 \\ &\geq \Omega\left(\frac{\eta L^{\frac{4}{3}}(\mathbf{W}(t))}{k^4\|\mathbf{v}\|^{\frac{2}{3}}}\right). \end{aligned}$$

793 According to Xu and Du [2023, Lemma 24], let  $C_s = \Omega\left(k^{-4}\|\mathbf{v}\|^{-\frac{2}{3}}\right)$ , then we have:

$$L(\mathbf{W}(T + T_2)) \leq \frac{1}{\left(L^{-\frac{1}{3}}(\mathbf{W}(T_2)) + \Omega\left(k^{-4}\|\mathbf{v}\|^{-\frac{2}{3}}\right)\eta T\right)^3}.$$

794 **Proof of Eq. (55):** According to Lemma 3, we have  $L(\mathbf{W}(T_2)) \leq \frac{1}{2}k^2\epsilon_2^{0.05}\|\mathbf{v}\|^2 = o\left(\frac{\|\mathbf{v}\|^2}{m^3k^3}\right)$ .

795 Then for  $\forall i \in [m]$ , according to [Xu and Du, 2023, Lemma 24], we have:

$$\begin{aligned} \|\mathbf{w}_i(T + T_2)\| &\geq \|\mathbf{w}_i(T_2)\| - \sum_{t=0}^{T-1} \eta\|\nabla_{\mathbf{W}}(t + T_2)\| \\ &\geq \frac{\|\mathbf{v}\|}{3m_{\tau_i}} - 8C_s^{-\frac{1}{2}}o\left(\frac{\|\mathbf{v}\|^2}{m^3k^3}\right)^{\frac{1}{3}} \\ &\geq \frac{\|\mathbf{v}\|}{3m_{\tau_i}} - 8\mathcal{O}\left(k^{-4}\|\mathbf{v}\|^{-\frac{2}{3}}\right)^{-\frac{1}{2}}o\left(\frac{\|\mathbf{v}\|^2}{m^3k^3}\right)^{\frac{1}{3}} \\ &\geq \frac{\|\mathbf{v}\|}{3m_{\tau_i}} - o\left(\frac{k\|\mathbf{v}\|}{m}\right) \\ &\geq \frac{\|\mathbf{v}\|}{4m_{\tau_i}}, \end{aligned}$$

796 and

$$\begin{aligned}\|\mathbf{w}_i(T + T_2)\| &\leq \|\mathbf{w}_i(T_2)\| + \sum_{t=0}^{T-1} \eta \|\nabla \mathbf{w}(t + T_2)\| \\ &\leq \frac{3\|\mathbf{v}\|}{m_{\tau_i}} + 8C_s^{-\frac{1}{2}} o\left(\frac{\|\mathbf{v}\|^2}{m^3 k^3}\right)^{\frac{1}{3}} \\ &\leq \frac{3\|\mathbf{v}\|}{m_{\tau_i}} + 8\mathcal{O}\left(k^{-4}\|\mathbf{v}\|^{-\frac{2}{3}}\right)^{-\frac{1}{2}} o\left(\frac{\|\mathbf{v}\|^2}{m^3 k^3}\right)^{\frac{1}{3}} \\ &\leq \frac{3\|\mathbf{v}\|}{m_{\tau_i}} + o\left(\frac{k\|\mathbf{v}\|}{m}\right) \\ &\leq \frac{4\|\mathbf{v}\|}{m_{\tau_i}},\end{aligned}$$

797 which finishes the proof.

798

□