
Guided by the Experts: Provable Feature Learning Dynamic of Soft-Routed Mixture-of-Experts

Fangshuo Liao¹

¹Computer Science Department, Rice University

Anastasios Kyrillidis^{1,2}

²Ken Kennedy Institute, Rice University

Abstract

Mixture-of-Experts (MoE) architectures have emerged as a cornerstone of modern AI systems. In particular, MoEs route inputs dynamically to specialized experts whose outputs are aggregated through weighted summation. Despite their widespread application, theoretical understanding of MoE training dynamics remains limited to either separate expert-router optimization or only top-1 routing scenarios with carefully constructed datasets. This paper advances MoE theory by providing convergence guarantees for joint training of soft-routed MoE models with non-linear routers and experts in a student-teacher framework. We prove that, with moderate over-parameterization, the student network undergoes a feature learning phase, where the router’s learning process is “guided” by the experts, that recovers the teacher’s parameters. Moreover, we show that a post-training pruning can effectively eliminate redundant neurons, followed by a provably convergent fine-tuning process that reaches global optimality. To our knowledge, our analysis is the first to bring novel insights in understanding the optimization landscape of the MoE architecture.

1 INTRODUCTION

Mixture-of-Experts (MoE) architectures have become a fundamental building block in modern artificial intelligence systems, enabling significant advances in model capacity without corresponding increases in computational costs (Shazeer et al., 2017; Fedus et al., 2022).

Proceedings of the 29th International Conference on Artificial Intelligence and Statistics (AISTATS) 2026, Tangier, Morocco. PMLR: Volume 300. Copyright 2026 by the author(s).

At its core, a MoE system treats a complicated task as a combination of multiple simpler tasks, which can be handled efficiently by smaller models. In concept, the conditional routing of the input to different sub-modules of the MoE allows each sub-module to “specialize” in its own domain, leading to an effective decoupling of the overall task complexity.¹ This approach has achieved remarkable success in large language models (LLMs) (Fedus et al., 2022), computer vision (Riquelme et al., 2021), and multi-modal agentic systems (Mustafa et al., 2022), where conditional computation provides an efficient way to scale model capacity.

The fundamental structure of an MoE layer consists of a set of expert networks (the “experts”) and a gating network that determines the contribution of each expert to the final output. While simple, this architecture presents theoretical challenges, particularly regarding the joint optimization of both components. The gating function, typically implemented using a softmax activation, introduces non-convexity that makes the analysis difficult. This is further complicated by the interplay between expert specialization and router assignments: experts must specialize in certain inputs, while the router must correctly identify the correct combination of experts appropriate for each input.

Despite the widespread adoption of MoE architectures in practice, such a theoretical understanding of their optimization dynamic remains limited; see Related Works section. Existing theoretical work has focused on either simplified linear models or has analyzed the experts and gating networks separately; e.g. Li et al. (2025) and Kawata et al. (2025) study the setting where the experts are trained first with the router parameter fixed, followed by a fine-tuning stage of the router itself. While such a setting simplifies the analysis by decoupling the updates of the router and expert parameters, it deviates from the more beneficial sce-

¹In this paper we consider a “classic” MoE instead of the MoE used to increase parameterization without increasing computational costs such as Fedus et al. (2022).

nario, where a joint optimization is applied to handle intricate task combinations (Kong et al., 2025; Zhang et al., 2025a).

A prior work (Chen et al., 2022) studies the joint optimization of the expert and router parameters in a top-1 routed MoE on patched input data, thus reducing the interference between the learning process of the experts in each gradient step. Although top-1 routing has been a popular approach (Fedus et al., 2022), most state-of-the-art language models such as Mixtral 8x7B (Jiang et al., 2024), DeepSeek-V3 (DeepSeek-AI et al., 2025), and Qwen 3 (Yang et al., 2025) uses a top- K routing with $K > 1$, leading to gaps between theory and practice.

In general, it remains open to study the optimization dynamics of MoEs with more than one activated experts, where experts and router are jointly trained. This gap is increasingly significant as MoE architectures become fundamental components in state-of-the-art AI systems. Analysis of these dynamics would not only inform architectural improvements but also provide formal guarantees about model behavior and performance. Moreover, such an analysis could help us understand better Agentic AI systems (Hu et al., 2025; Zhang et al., 2025b,c), where component orchestration mirrors MoE routing mechanisms (Bhatt et al., 2025). These systems must dynamically select appropriate specialized modules (tools, APIs, or reasoning components) based on input context—functionally analogous to expert selection in MoE architectures.

Contributions. Given the difficulty of the task, we focus on the learning of a MoE model with one-layer sigmoid router and non-linear experts over the mean-square-error (MSE) loss in a teacher-student set-up on high dimensional Gaussian input. In particular, we show that, with moderate over-parameterization and under the gradient flow training, the student MoE enjoys a near-perfect recovery of the feature from the teacher model’s in a sequential order in $\mathcal{O}(\sqrt{d})$ time, where d is the dimensionality of the input. Moreover, after the feature learning stage, a greedy pruning can be applied to remove the unused experts. Lastly, the post-pruning fine-tuning of the student model converges to zero loss.

Notations. Without further specification, we use regular lower-case letters (e.g. a) to denote scalars, bold-face lower-case letters (e.g. \mathbf{a}) to denote vectors, and bold-face capital letters (e.g. \mathbf{A}) to denote matrices. We use $\mathcal{N}(\mu, \sigma^2)$ to denote the the Gaussian distribution with mean μ and (co)variance σ^2 . For a function $f(x)$, we use $f'(x)$, $f''(x)$ and $f'''(x)$ to denote its first three order derivatives, and $f^{(a)}(x)$ to denote is arbitray a th order derivative. We use

$\text{poly}(x_1, \dots, x_n)$ to denote the polynomial dependency in terms of x_1, \dots, x_n .

2 RELATED WORKS

Theory of Mixture-of-Experts. From an optimization perspective, Chen et al. (2022) studies the convergence rate of top-1 MoE with CNN experts on patched input data. Chowdhury et al. (2023) studies the patch-level routing under both the setting with a separately trained expert and router, and the setting of pre-trained experts. Chowdhury et al. (2024) shows the pruning effectiveness after fine-tuning a pre-trained MoE model. Kawata et al. (2025) considers the training both a top-1 MoE and a ReLU routed MoE, but under a four-stage training algorithm. Fruytier et al. (2025) studies the convergence of the Expectation-Maximization algorithm for learning MoE. Li et al. (2025) studies the optimization of MoE in a continual learning set-up. From the perspective of sample complexity, Nguyen et al. (2024a,c, 2025) studies the sample complexity of correctly identifying experts for softmax MoE under both the logistic loss and the MSE loss. Nguyen et al. (2024b) shows that sigmoid gated MoEs enjoy a better sample complexity compared with softmax gated MoEs. Other works (Kratsios et al., 2024; Wang and E, 2025) studies MoE under operator learning, and the expressive power of MoEs, respectively. Following the expressivity line of work, Boix-Adsera and Rigollet (2025) studies how the granularity of the experts affects the expressive power of MoEs.

Feature Learning of Neural Networks. As its name suggested, the feature learning framework explores the ability of the neural network to learn intrinsic features of the dataset, which is an ability not present in the traditional Neural Tangent Kernel framework (Jacot et al., 2020; Du et al., 2019). In particular Shi et al. (2022, 2023) studies the hidden-neuron evolution during training, and Damian et al. (2022); Mousavi-Hosseini et al. (2023a) investigates how gradient-based learning discovers the intrinsic low-dimensional subspace of data. Along this line of work Ba et al. (2023); Mousavi-Hosseini et al. (2023b) studies the learning with data sampled from distribution with a spiked covariance matrix. A related line of work studies the feature learning dynamics of two-layer ReLU networks, including the characterization of gradient flow dynamics under orthogonal inputs (Boursier et al., 2022), the phase diagram at infinite width (Luo et al., 2021), quantization of learned features under gradient descent (Maennel et al., 2018), early alignment phenomena and their implications for robustness and generalization (Boursier and Flammarion, 2025; Min et al., 2023), as well as implicit biases of gradient descent leading to feature averaging (Li et al., 2024).

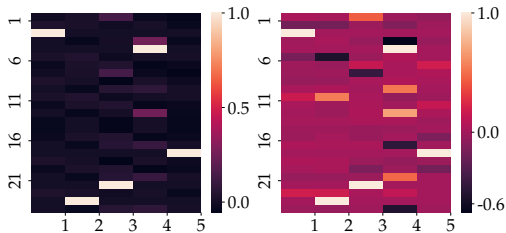


Figure 1: Training MoE in (1) on (3) with $m^* = 5, m = 25$, and $d = 1000$ with online batch SGD simulating GF on the population loss. Left: alignment values of the router parameters $\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j^*$. Right: alignment values of the expert parameters $\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j^*$.

Recently, a popular line of work studies the learning of Gaussian single/multi-index models (Bietti et al., 2022; Lee et al., 2024; Ren et al., 2025; Bietti et al., 2023; Ba et al., 2023; Şimşek et al., 2025). Noticeably, this line of work adopts the Hermite expansion of the nonlinear function to transform the loss objective into a form similar to the tensor decomposition (Ge et al., 2017). In terms of the proof technique, our work is similar to Ren et al. (2025) by utilizing the sharp phase transition that occurs from the high information exponent of the activation function.

3 PRELIMINARY AND SET-UP

Student model. In this paper, we consider the training of a normalized MoE with m experts. In particular, given inputs $\mathbf{x} \in \mathbb{R}^d$, we study the setting of a one-layer router with parameter \mathbf{V} , given by $\pi(\bar{\mathbf{V}}\mathbf{x})$. Here $\pi(\cdot)$ is an entry-wise sigmoid function, and $\bar{\mathbf{V}}$ denote the row-wise normalized version of \mathbf{V} . In short, we have

$$\pi(\bar{\mathbf{V}}\mathbf{x})_i := \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) = \pi\left(\frac{\mathbf{v}_i^\top \mathbf{x}}{\|\mathbf{v}_i\|_2}\right); \forall i \in [m]$$

where \mathbf{v}_i is the i th row of \mathbf{V} . We consider each expert as a one-layer non-linear function with parameter \mathbf{w}_i given by $\sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) = \sigma\left(\frac{\mathbf{w}_i^\top \mathbf{x}}{\|\mathbf{w}_i\|_2}\right)$. In this paper, we choose $\sigma(a) = a^3 - 3a$ to be the third-order Hermite polynomial. Letting $\boldsymbol{\theta} = \{(\mathbf{v}_i, \mathbf{w}_i)\}_{i=1}^m$, then the student model is given by

$$f(\boldsymbol{\theta}, \mathbf{x}) := \sum_{i=1}^m \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) \quad (1)$$

Our choice of the sigmoid router is motivated by Nguyen et al. (2024b) which shows that sigmoid routing is more sample efficient than softmax routing. Moreover, using Hermite polynomial as the activation function has been a popular approach in previous study of the feature learning mechanism of neural

networks (Arous et al., 2025). Lastly, the choice of normalizing the weights is also a popular choice in prior works (Wang et al., 2020; Ren et al., 2025).

Data and Teacher Model. The teacher model f^* we consider has the same structure as in the student model, but with m^* experts. In addition, we assume that the parameter of the teacher model’s parameters $\bar{\mathbf{v}}_1^*, \dots, \bar{\mathbf{v}}_{m^*}^*, \bar{\mathbf{w}}_1^*, \dots, \bar{\mathbf{w}}_{m^*}^*$ satisfy the following assumption:

Assumption 1 (Teacher Orthonormality). *The teacher model’s parameters $\bar{\mathbf{v}}_1^*, \dots, \bar{\mathbf{v}}_{m^*}^*, \bar{\mathbf{w}}_1^*, \dots, \bar{\mathbf{w}}_{m^*}^*$ form an orthonormal list.*

We consider the input data that come from a standard Gaussian distribution $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, and labels are generated by the teacher model according to

$$y = f^*(\mathbf{x}) := \sum_{i=1}^{m^*} \pi(\bar{\mathbf{v}}_i^{*\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^{*\top} \mathbf{x}) \quad (2)$$

Intuitively, this set-up implies that the input space \mathbb{R}^d is softly partitioned by the teacher’s router. The goal of the student model is to learn both the features $\bar{\mathbf{v}}_i^*$ s that gives a correct partitions, and the features $\bar{\mathbf{w}}_i^*$ s that leads to the effective specialization of experts.

We consider training the student model $f(\boldsymbol{\theta}, \mathbf{x})$ on the population mean-squared error (MSE) loss $\mathcal{L}(\boldsymbol{\theta})$ using gradient flow $\frac{d}{dt}\boldsymbol{\theta}(t) = -\nabla\mathcal{L}(\boldsymbol{\theta}(t))$ over the data distribution defined by the teacher model. To be more specific, the MSE has the form

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{2} \mathbb{E}_{\mathbf{x}, y} \left[(f(\boldsymbol{\theta}, \mathbf{x}) - y)^2 \right] \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)} \left[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}))^2 \right] \end{aligned} \quad (3)$$

We assume that the parameters of the student model are initialized according to the following assumption:

Assumption 2 (Initialization). *The expert weights are initialized as $\mathbf{w}_i(0) \sim \mathcal{N}(\mathbf{0}, d^{-1}\mathbf{I}_d)$. The router weights are initialized as $\mathbf{v}_i(0) = \hat{\mathbf{v}}_i(0) - \hat{\mathbf{v}}_i(0)^\top \bar{\mathbf{w}}_i(0) \cdot \bar{\mathbf{w}}_i(0)$, where $\hat{\mathbf{v}}_i(0) \sim \mathcal{N}(\mathbf{0}, d^{-1}\mathbf{I}_d)$, so as to decouple the router and expert weights at initialization.²*

Lastly, we make the following assumption on the sigmoid function, which is numerically checked in Appendix E.

Assumption 3. *Let $z_1, z_2 \sim \mathcal{N}(0, 1)$ with arbitrary covariance $\text{Cov}(z_1, z_2) \in [-1, 1]$, it holds that $\mathbb{E}_{z_1, z_2} [\pi'(z_1)\pi^{(3)}(z_2)] \leq 0$*

We empirically verify that this setting allows that each of the m^* experts and corresponding router parameter

²The behavior that $\mathbf{v}_i(t)^\top \mathbf{w}_i(t) = 0$ does not hold throughout gradient flow training.

in the teacher model can be recovered by one and only one expert and corresponding router in the student model. According to Figure 1, for each $\bar{\mathbf{v}}_j^*, \bar{\mathbf{w}}_j^*$ in the teacher (x-axis), there is one and only one expert and router $\bar{\mathbf{v}}_j, \bar{\mathbf{w}}_j$ that converges to $\bar{\mathbf{v}}_j^*, \bar{\mathbf{w}}_j^*$ (lighter color indicates that $\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j^*$ and $\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j^*$ are closer to one.)

4 MAIN RESULT: SEQUENTIAL FEATURE LEARNING

In this section, we present the main result of the feature learning phase. Recall the set-up of the student and teacher MoE models in (1) and (2). An ideal feature learning result would be that, for each router-expert pair $(\bar{\mathbf{v}}_i^*, \bar{\mathbf{w}}_i^*)$ in the teacher model, there is an exclusive router-expert pair $(\bar{\mathbf{v}}_i, \bar{\mathbf{w}}_i)$ that converges to it. The theorem below states that the matching between the router-expert pair from the teacher model and the router-expert pair from the student model happens in a sequential order.

Theorem 1. *Consider training the MoE model $f(\boldsymbol{\theta}, \mathbf{x})$ in (1) with respect to a teacher model given by (2) using the gradient flow on the population MSE loss in (3). Let $\delta_{\mathbb{P}} \in (0, 1/7)$ be given. If $m \geq \Omega\left(m^* \log \frac{m^*}{\delta_{\mathbb{P}}}\right)$ and $d \geq \text{poly}(m, \delta_{\mathbb{P}}^{-1})$, then there exists an injective mapping $\mathcal{I} : [m^*] \rightarrow [m]$ and time steps $0 \leq T_1 \leq \dots \leq T_{m^*} \leq T^* \leq \mathcal{O}\left(\sqrt{d}\right)$ such that for all $\ell \in [m^* - 1]$ and $t \in [T_\ell, T_{\ell+1})$, we have that*

- (Recovered expert-router pairs) $\bar{\mathbf{v}}_{\mathcal{I}(i)}(t)^\top \bar{\mathbf{v}}_i^* \geq 0.9$ and $\bar{\mathbf{w}}_{\mathcal{I}(i)}(t)^\top \bar{\mathbf{w}}_i^* \geq 0.9$ for all $i \leq \ell$.
- (Unrecovered expert-router pairs) For all $i > \ell$, $\max\{|\bar{\mathbf{v}}_{\mathcal{I}(i)}(t)^\top \bar{\mathbf{v}}_i^*|, |\bar{\mathbf{w}}_{\mathcal{I}(i)}(t)^\top \bar{\mathbf{w}}_i^*|\} \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$

Moreover, for $T^* \leq t \leq T^* + \mathcal{O}\left(\frac{\delta_{\mathbb{P}}\sqrt{d}}{m^2}\right)$, we have that

- (Learned features) $\bar{\mathbf{v}}_{\mathcal{I}(i)}(t)^\top \bar{\mathbf{v}}_i^*, \bar{\mathbf{w}}_{\mathcal{I}(i)}(T)^\top \bar{\mathbf{w}}_i^* \geq 1 - \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$ for all $i \in [m^*]$
- (Unused expert-router pairs) for all $i_1 \in [m] \setminus \mathcal{I}([m^*]), i_2 \in [m], i_1 \neq i_2$ and $j \in [m^*]$, the following quantities

$$\begin{aligned} & |\bar{\mathbf{v}}_{i_1}(t)^\top \bar{\mathbf{v}}_{i_2}(t)|, |\bar{\mathbf{v}}_{i_1}(t)^\top \bar{\mathbf{w}}_{i_1}(t)|, |\bar{\mathbf{w}}_{i_1}(t)^\top \bar{\mathbf{w}}_{i_2}(t)| \\ & |\bar{\mathbf{v}}_{i_1}(t)^\top \bar{\mathbf{w}}_{i_2}(t)|, |\bar{\mathbf{v}}_{i_2}(t)^\top \bar{\mathbf{w}}_{i_1}(t)| \\ & |\bar{\mathbf{v}}_{i_1}(t)^\top \bar{\mathbf{v}}_j^*|, |\bar{\mathbf{v}}_{i_1}(t)^\top \bar{\mathbf{w}}_j^*|, |\bar{\mathbf{w}}_{i_1}(t)^\top \bar{\mathbf{w}}_j^*|, |\bar{\mathbf{w}}_{i_1}(t)^\top \bar{\mathbf{v}}_j^*| \end{aligned}$$

are all upper bounded by $\mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$

In particular, Theorem 1 states the result that the MoE training in our set-up undergoes a sequential feature learning phase. As requirements of the theorem, we need m to be as large as $\Omega\left(m^* \log \frac{m^*}{\delta_{\mathbb{P}}}\right)$ to ensure that, at initialization, at least one of the router-expert pair from the student model have a good enough alignment for each router-expert pair in the teacher. Also we require d to be polynomially large in terms of m and $\delta_{\mathbb{P}}^{-1}$ to control the interference between the convergence of each router-expert pair, as well as between the convergence of the router parameter and the expert parameter. Due to the polynomial scaling of d in terms of m and $\delta_{\mathbb{P}}^{-1}$, quantities $\mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ and $\mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$ are in general small, and $\mathcal{O}\left(\frac{\delta_{\mathbb{P}}\sqrt{d}}{m^2}\right)$ is large.

In the set-up of the theorem, the mapping \mathcal{I} establishes the correct matching between the expert-router pair in the teacher model and the router-expert pair in the student model. Ideally, we expect $\bar{\mathbf{v}}_{\mathcal{I}(i)}$ and $\bar{\mathbf{w}}_{\mathcal{I}(i)}$ to converge to $\bar{\mathbf{v}}_i^*$ and $\bar{\mathbf{w}}_i^*$. Two key points of the theorem are outlined below:

Sequential weak recovery. The first part of the theorem states that such convergence happens in a sequential order. By its set-up, T_ℓ denotes the time where the first ℓ pairs of $(\bar{\mathbf{v}}_{\mathcal{I}(i)}, \bar{\mathbf{w}}_{\mathcal{I}(i)})$ for $i \leq \ell$ just achieved a weak convergence to $\bar{\mathbf{v}}_i^*$ and $\bar{\mathbf{w}}_i^*$ by achieving an inner product with of at least 0.9. In the mean time, before t reaches $T_{\ell+1}$, all the remaining pairs $(\bar{\mathbf{v}}_{\mathcal{I}(i)}, \bar{\mathbf{w}}_{\mathcal{I}(i)})$ for $i > \ell$ still stays within a small alignment value with their reference.

Order of Recovery. We remark that the order in which the student's expert-router pairs converge to the teacher's parameters is determined by the alignment values $\mathbf{w}_i(0)^\top \bar{\mathbf{w}}_j^*$ at initialization, as can be seen from the proof sketch in Section 4.1. Intuitively, a student expert i that achieves a larger initial alignment with a teacher expert j will converge to $\bar{\mathbf{w}}_j^*$ earlier than a student expert with smaller initial alignment, provided that $\bar{\mathbf{w}}_i$ has not converged to another $\bar{\mathbf{w}}_j^*$, and no other $\bar{\mathbf{w}}_{i'}$ has converged to $\bar{\mathbf{w}}_j^*$ yet. While one might expect the norms $\|\mathbf{w}_i(0)\|_2$ to play a role in governing this ordering, their influence is in fact negligible: by concentration, the norms of the randomly initialized weight vectors are tightly concentrated around their expectation, and thus do not meaningfully differentiate the convergence order across experts. In contrast, as shown in Lemma 2, the inner products $\mathbf{w}_i(0)^\top \bar{\mathbf{w}}_j^*$ vary across pairs and serve as the primary signal that breaks the symmetry among student experts, ultimately determining which student expert specializes to which teacher expert and at what rate.

Near-perfect recovery. The second part of the the-

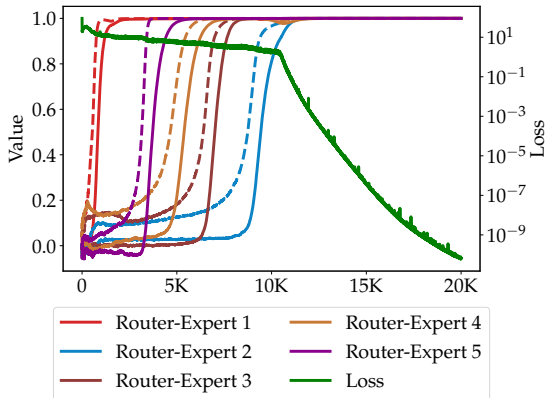


Figure 2: Dynamics of the routers’ and experts’ alignment value with the teacher’s parameter under the same set-up as Figure 1. The green curve denotes the loss value. Except for the green curve, dashed line and solid line of the same color denotes a pair of router and expert alignment value.

orem shows that for any time t that exceeds some $T^* \leq \mathcal{O}(\sqrt{d})$ but stays under $T^* + \mathcal{O}(\frac{\delta_P \sqrt{d}}{m^2})$, the learned features have converged to inner product values of at least $1 - \mathcal{O}(\frac{m^7}{\delta_P d^{\frac{3}{2}}})$. In the meantime, all the router and expert parameters in the student model that did not converge to any teacher’s parameter must stay nearly orthogonal both to the teacher model’s parameter and to each other. In Section 5, we will utilize this property to prove the theoretical guarantee of pruning these unused experts in the student model.

4.1 Guided by the Experts: A Proof Sketch of Theorem 1

In this section, we will discuss the difficulties and techniques arises in the proof of Theorem 1.

Hermite expansion of the loss and gradient. The starting point of our proof relies on the Hermite expansion of non-linear functions to study its property with Gaussian inputs. Let $He_k(x)$ denote the k th-order probabilist’s Hermite polynomial. It is known that the set of Hermite polynomials $\{He_k(x)\}_{k=0}^\infty$ consists an basis of the square integrable functions under the Gaussian measure. Therefore, we can expand the sigmoid function as

$$\pi(x) = \sum_{k=0}^{\infty} \frac{c_k}{k!} He_k(x); c_k = \mathbb{E}_{x \sim \mathcal{N}(0,1)} \left[\pi^{(k)}(x) \right]$$

Here c_k ’s are the Hermite coefficients of $\pi(x)$. Since

$\sigma(x) = He_3(x)$, $f(\boldsymbol{\theta}, \mathbf{x})$ and $f^*(\mathbf{x})$ has the form

$$f(\boldsymbol{\theta}, \mathbf{x}) = \sum_{k=0}^{\infty} \frac{c_k}{k!} He_k(\bar{\mathbf{v}}_i^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_i^\top \mathbf{x})$$

$$f^*(\mathbf{x}) = \sum_{k=0}^{\infty} \frac{c_k}{k!} He_k(\bar{\mathbf{v}}_i^{*\top} \mathbf{x}) He_3(\bar{\mathbf{w}}_i^{*\top} \mathbf{x})$$

Carrying this idea to the setting of minimizing the MSE in 3, we notice that $\mathcal{L}(\boldsymbol{\theta})$ can be written as

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} [f(\boldsymbol{\theta}, \mathbf{x})^2] - \mathbb{E}_{\mathbf{x}} [f(\boldsymbol{\theta}, \mathbf{x}) f^*(\mathbf{x})] + \frac{1}{2} \mathbb{E}_{\mathbf{x}} [f^*(\mathbf{x})^2]$$

where one could notice that the last term does not depend on $\boldsymbol{\theta}$. However, the first two terms involves second-order terms on $f(\boldsymbol{\theta}, \mathbf{x})$ and $f^*(\mathbf{x})$. As an illustration, we expand the first term as

$$\mathbb{E}_{\mathbf{x}} [f(\boldsymbol{\theta}, \mathbf{x})^2] = \sum_{i,j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} \mathcal{C}_{k,\ell,3,3}^{i,j}$$

$$\mathcal{C}_{k,\ell,3,3}^{i,j} =$$

$$\mathbb{E}_{\mathbf{x}} [He_k(\bar{\mathbf{v}}_i^\top \mathbf{x}) He_\ell(\bar{\mathbf{v}}_j^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_i^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_j^\top \mathbf{x})]$$

Although a large body of prior work has exploited the nice property of Hermite polynomial that $\mathbb{E}_{\mathbf{x}} [He_k(\mathbf{u}_1^\top \mathbf{x}) He_\ell(\mathbf{u}_2^\top \mathbf{x})] = k! (\mathbf{u}_1^\top \mathbf{u}_2)^k \mathbb{I}\{k=\ell\}$ for $\mathbf{u}_1, \mathbf{u}_2$ with unit norm, in our setting we have to deal with the expectation of the product of four Hermite polynomial. Our main tool of handling this difficulty is the lemma below.³

Lemma 1. *Let $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$. For some multi-index $\mathbf{k} \in \mathbb{N}^n$, we define the multi-variate Hermite polynomial as*

$$He_{\mathbf{k}}(\mathbf{x}) = \prod_{i=1}^n He_{\mathbf{k}[i]}(\mathbf{x}[i])$$

Then we have that for $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$,

$$\mathbb{E}_{\mathbf{x}} [He_{\mathbf{k}}(\mathbf{x})] = \left(\prod_{i=1}^n \mathbf{k}[i]! \right) \sum_{\mathbf{M} \in \mathcal{S}} \prod_{i,j=1}^n \frac{\boldsymbol{\Sigma}[i,j]^{\mathbf{M}[i,j]}}{\mathbf{M}[i,j]!}$$

where the set \mathcal{S} is the set of symmetric matrices $\mathbf{M} \in \mathbb{N}^{n \times n}$ satisfying

$$\mathbf{M}[i,i] = 0; \sum_{j=1}^n \mathbf{M}[i,j] = \mathbf{k}[i]; \forall i \in [n]$$

In Lemma 1, each $\mathbf{M} \in \mathcal{S}$ can be considered as the adjacency matrix of a graph with n nodes and integer weights such that the degree of the i th node is

³We are not the first to introduce this result. However, we could not find a formal published source that proves the result.

$k[i]$. Applying Lemma 1 to our case thus only requires to enumerate the graphs of four nodes with degree $(k, \ell, 3, 3)$. With Lemma 1, we are able to derive the form of $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta})$ and $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta})$, whose exact form are omitted from the main text due to its intricacy.

Gradient flow dynamic of target alignment. Recall that our goal is to show that there is some $\bar{\mathbf{v}}_i, \bar{\mathbf{w}}_i$ that converges to $\bar{\mathbf{v}}_j^*, \bar{\mathbf{w}}_j^*$ for $j \in [m^*]$. Thus, it is intuitive to start with tracking the dynamic of $\bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j^*$ and $\bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^*$. With gradient flow, we have that

$$\begin{aligned} \frac{d}{dt} \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j^* &= -\|\mathbf{v}_i(t)\|_2^{-1} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta})^\top \bar{\mathbf{v}}_j^* \\ \frac{d}{dt} \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^* &= -\|\mathbf{w}_i(t)\|_2^{-1} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta})^\top \bar{\mathbf{w}}_j^* \end{aligned} \quad (4)$$

due to the fact that $\|\mathbf{v}_i\|_2$ and $\|\mathbf{w}_i\|_2$ stays constant during the gradient flow process. Utilizing the fact that $\bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j^*, \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^* \ll 1$ when $\bar{\mathbf{v}}_i(t)$ and $\bar{\mathbf{w}}_i(t)$ are near their initialization, we can utilize Lemma 1 to approximate (4) as

$$\begin{aligned} \frac{d}{dt} \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j^* &\propto (\bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^*)^3 \\ \frac{d}{dt} \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^* &\propto (\bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^*)^2 \end{aligned} \quad (5)$$

The approximation above exhibits two interesting behaviors near initialization. First, the improvements in both the router alignment $\bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j^*$ and the expert alignment $\bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^*$ depends on the current magnitude of the expert alignment. This implies that, once the expert alignment reaches a magnitude of $\Omega(1)$, it will take only constant time for the router and expert alignments to grow to a value near one (perfect alignment). This behavior can be observed in Figure 2 where the router alignment values (solid lines except for the green one) follow closely as the expert alignment values (dashed lines) increases.

Second, the quadratic dependency in the expert alignment dynamic induces a sharp phase transition where the alignment value starts off slow for a period of time, and suddenly increases with a fast speed (see dashed lines in Figure 2), as studied in Ren et al. (2025). This sharp phase transition is particularly helpful to prevent multiple experts from the student model to converge to the same expert in the teacher model. As an example, consider dynamics $\bar{\mathbf{w}}_1(t)^\top \bar{\mathbf{w}}_j^*$ and $\bar{\mathbf{w}}_2(t)^\top \bar{\mathbf{w}}_j^*$ with a small different $\Delta \geq 0$ at initialization

$$0 \leq (1 + \Delta) \bar{\mathbf{w}}_1(t)^\top \bar{\mathbf{w}}_j^* \leq \bar{\mathbf{w}}_2(t)^\top \bar{\mathbf{w}}_j^* \leq \tilde{\mathcal{O}} \left(\frac{1}{\sqrt{d}} \right)$$

Solving the ODE in (5) gives that for $i \in \{1, 2\}$

$$\bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^* \approx \left((\bar{\mathbf{w}}_i(0)^\top \bar{\mathbf{w}}_j^*)^{-1} - t \right)^{-1}$$

Thus, the time T required for $\bar{\mathbf{w}}_2(t)^\top \bar{\mathbf{w}}_j^* \geq \frac{1}{2}$ is $T = (\bar{\mathbf{w}}_2(0)^\top \bar{\mathbf{w}}_j^*)^{-1} - 2 \leq (\bar{\mathbf{w}}_2(0)^\top \bar{\mathbf{w}}_j^*)^{-1}$. However, at time T , one can compute that

$$\begin{aligned} \bar{\mathbf{w}}_1(t)^\top \bar{\mathbf{w}}_j^* &\leq \left((\bar{\mathbf{w}}_1(0)^\top \bar{\mathbf{w}}_j^*)^{-1} - (\bar{\mathbf{w}}_2(0)^\top \bar{\mathbf{w}}_j^*)^{-1} \right)^{-1} \\ &\leq \frac{\bar{\mathbf{w}}_2(0)^\top \bar{\mathbf{w}}_j^*}{\Delta} \leq \tilde{\mathcal{O}} \left(\frac{\Delta^{-1}}{\sqrt{d}} \right) \end{aligned}$$

When $\sqrt{d} \gg \Delta^{-1}$, we can conclude that $\bar{\mathbf{w}}_2(T)^\top \bar{\mathbf{w}}_j^* \geq \frac{1}{2}$ while $\bar{\mathbf{w}}_1(T)^\top \bar{\mathbf{w}}_j^* \ll \frac{1}{2}$. This behavior implies that, the expert in the student model that aligns with expert j in the teacher best at initialization will converge to some $\Omega(1)$ quickly while the other experts' alignment remains small. Below, we formalize this dominance determined by the initialization.

Alignment gap at Initialization. We show that there is a set of experts in the student model that aligns with each expert in the teacher model good enough to create a gap compared with other experts in the student model. At a high level, our goal here is to construct the mapping \mathcal{I} in Theorem 1 based on the initialization. Our approach is a greedy forward selection similar to Ren et al. (2025). In particular, we define $\mathcal{R}_\ell = \{i_k^*\}_{k=1}^\ell$ and $\mathcal{C}_\ell = \{j_\ell^*\}_{k=1}^\ell$ recursively as follows

$$i_{\ell+1}^*, j_{\ell+1}^* = \arg \max_{i \in [m] \setminus \mathcal{R}_\ell, j \in [m^*] \setminus \mathcal{C}_\ell} \mathbf{w}_i(0)^\top \bar{\mathbf{w}}_j^* \quad (6)$$

We expect that $\mathcal{I}(j_\ell^*) = i_\ell^*$. Namely, we expect $\bar{\mathbf{v}}_{i_\ell^*}(t)^\top \bar{\mathbf{v}}_{j_\ell^*}^*$ and $\bar{\mathbf{w}}_{i_\ell^*}(t)^\top \bar{\mathbf{w}}_{j_\ell^*}^*$ to converge to 1. The index ℓ denotes the order of the sequential convergence. That is, we expect that $\bar{\mathbf{v}}_{i_1^*}(t)^\top \bar{\mathbf{v}}_{j_1^*}^*$ to grow large first, followed by $\bar{\mathbf{v}}_{i_2^*}(t)^\top \bar{\mathbf{v}}_{j_2^*}^*$, etc. Our theorem below shows that at initialization, the pairs $\bar{\mathbf{w}}_{i_\ell^*}(t)^\top \bar{\mathbf{w}}_{j_\ell^*}^*$ has a gap compared with other alignment values.

Lemma 2. *Let $\mathbf{w}_1, \dots, \mathbf{w}_m \sim \mathcal{N}(0, d^{-1} \mathbf{I}_d)$ be I.I.D. Gaussian random vectors. Define*

$$\begin{aligned} i_\ell^*, j_\ell^* &= \arg \max_{i \in [m] \setminus \mathcal{R}_{\ell-1}, j \in [m^*] \setminus \mathcal{C}_{\ell-1}} \mathbf{w}_i[j] \\ \mathcal{R}_\ell &= \{i_k^*\}_{k=1}^\ell; \mathcal{C}_\ell = \{j_k^*\}_{k=1}^\ell \end{aligned}$$

Let any $\delta_{\mathbb{P}} \in (0, 1/2)$ be given. Then there exists some absolute constant $\beta_2, \beta_4 > 0$ such that if $m \geq \beta_4 m^ \log \frac{m^*}{\delta_{\mathbb{P}}}$, then for $\delta_s = \frac{\beta_2 \delta_{\mathbb{P}}}{m^2}$, with probability at least $1 - 4\delta_{\mathbb{P}}$, it holds that*

- (Row-wise Gap) $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq (1 + 2\delta_s) \mathbf{w}_{i_\ell^*}[j]$ for all $\ell \in [m^*]$ and $j \in [m^*] \setminus \mathcal{C}_\ell$
- (Column-wise Gap) $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq (1 + 2\delta_s) \mathbf{w}_i[j_\ell^*]$ for all $\ell \in [m^*]$ and $i \in [m] \setminus \mathcal{R}_\ell$
- (Threshold Gap) $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq (1 + 2\delta_s) \mathbf{w}_{i_{\ell+1}^*}[j_{\ell+1}^*]^2$ for all $\ell \in [m^* - 1]$

- (Magnitude Lower Bound) $\mathbf{w}_{i_\ell^*}[j_\ell^*]^2 \geq \frac{\log m^*}{d}$ for all $\ell \in [m]$

Since the standard Gaussian distribution is rotational invariant, we can the gaps and lower bound shown in Lemma 2 to the initial alignment scores $\bar{\mathbf{w}}_i(0)^\top \bar{\mathbf{w}}_j^*$. Roughly speaking, the *row-wise gap* facilitates that $\bar{\mathbf{w}}_{i_\ell^*}$ will not converge to $\bar{\mathbf{w}}_j$ s with $j \neq j_\ell^*$; the *column-wise gap* induces the fact that no $\bar{\mathbf{w}}_i$ will converge to $\bar{\mathbf{w}}_{j_\ell^*}$ except for $\bar{\mathbf{w}}_{i_\ell^*}$. Moreover, the *threshold gap* leads to the sequential recovery as stated in Theorem 1. Finally, the *magnitude lower bound* guarantees that at the target alignment values at initialization are not too small for the whole convergence process to be too long.

Induction-based Proof. With the goal of tracking the growth of $\bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j^*$ and $\bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^*$ in mind, however, we also have to track the ‘‘mis-alignments’’ including $\bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{w}}_j^*$, $\bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{v}}_j^*$ and ‘‘self-alignments’’ $\bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j(t)$, $\bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j(t)$, $\bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{w}}_j(t)$ due to the complicated form of the gradient, as can be seen from (10) from the Appendix, so that their value does not interrupt with the target dynamics. To this end, our proof is an induction on $\ell \in [m^*]$ that assumes

- $\bar{\mathbf{v}}_{i_{\ell'}^*}(t)^\top \bar{\mathbf{v}}_j^*$ and $\bar{\mathbf{w}}_{i_{\ell'}^*}(t)^\top \bar{\mathbf{w}}_j^*$ are close to one for $\ell' < \ell$, i.e., the top- $\ell - 1$ router and experts are recovered well while the ℓ th router-expert pair still remains not learned.
- The ‘‘mis-alignments’’ and ‘‘self-alignments’’ associated with the recovered router-expert pairs must be small throughout the process.

to show that $\bar{\mathbf{v}}_{i_\ell^*}(t)^\top \bar{\mathbf{v}}_j^*$ and $\bar{\mathbf{w}}_{i_\ell^*}(t)^\top \bar{\mathbf{w}}_j^*$ converge to a close-to-one value. A formal statement of the inductive hypothesis is provided in Appendix A.1, and the complete proof is provided in Appendix A.

5 PRUNING AND FINE-TUNING

Theorem 1 guarantees that in $\mathcal{O}(\sqrt{d})$ time, the student MoE model trained with gradient flow extracts m^* -pairs of near-perfect features from the teacher model. However, recall that the student have an over-parameterization of $m \geq \Omega\left(m^* \log \frac{m^*}{\delta_p}\right)$. Despite being moderate, the $\log \frac{m^*}{\delta_p}$ factor still leads to a large number of excessive parameters. Continuing to train these unused experts together with their corresponding router parameter results in a wast of the computation resource, regardless of whether they can converge to zero output or not. This theoretical insight corresponds with existing empirical works (Lu et al., 2024; Chowdhury et al., 2024; Zhang et al., 2025d) which discovers the existence of redundant experts in pre-trained LLMs.

5.1 Pruning the Redundant Experts

In this section, we adopt a greedy pruning algorithm based on the test loss similar to Lu et al. (2024) to remove the redundant experts, and show that, if we apply the algorithm at $T^* \leq t \leq T^* + \mathcal{O}\left(\frac{\sqrt{d}}{\delta_p m^2}\right)$, then we can provably remove all the unused experts and keep all the correctly learned router-expert pairs as stated in Theorem 1. To state the algorithm, we first define the sub-model MoE induced by $\mathcal{S} \subseteq [m]$ as

$$f_{\mathcal{S}}(\boldsymbol{\theta}, \mathbf{x}) = \sum_{i \in [m] \setminus \mathcal{S}} \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})$$

We consider the following pruning procedure that iteratively constructs the pruned set \mathcal{S} . In the τ th step, we identify an index $r_\tau \in [m] \setminus \mathcal{S}_{\tau-1}$

$$r_\tau = \arg \min_{r \in [m] \setminus \mathcal{S}_{\tau-1}} \mathbb{E}_{\mathbf{x}} \left[(f_{\mathcal{S}_{\tau-1} \cup \{r\}}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}))^2 \right] \quad (7)$$

$$\mathcal{S}_\tau = \mathcal{S}_{\tau-1} \cup \{r_\tau\}$$

The procedure will stop when pruning one more expert does not improve the population loss. In particular, we define the stopping step τ^* be such that

$$\min_{r \in [m]} \mathbb{E}_{\mathbf{x}} \left[(f_{\mathcal{S}_{\tau^*} \cup \{r\}}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}))^2 \right] \geq \mathbb{E}_{\mathbf{x}} \left[(f_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}))^2 \right] \quad (8)$$

For the simplicity of the analysis, we assume that $\mathcal{I}(i) = i$, as reordering the router-expert pairs does not change $f(\boldsymbol{\theta}, \mathbf{x})$. To facilitate the analysis of the pruning procedure, we make the following assumption

Assumption 4. Let $\{\mathbf{v}_i\}_{i=1}^m$ and $\{\mathbf{w}_i\}_{i=1}^m$ be the router and expert weights of the MoE model in (1). Let $\{\bar{\mathbf{v}}_i^*\}_{i=1}^{m^*}$ and $\{\bar{\mathbf{w}}_i^*\}_{i=1}^{m^*}$ be the router and expert weights of the teacher model in (2). There exists $\varepsilon \leq o\left(\frac{1}{\sqrt{m}}\right)$ such that for all $i \in [m^*]$ it holds that

$$\min \{\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_i^*, \bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_i^*\} \geq 1 - \varepsilon,$$

for all $i_1 \in [m] \setminus [m^*]$, $i_2 \in [m]$, $i_1 \neq i_2$ and $j \in [m^*]$ it holds that

$$\begin{aligned} & |\bar{\mathbf{v}}_{i_1}^\top \bar{\mathbf{v}}_{i_2}|, |\bar{\mathbf{w}}_{i_1}^\top \bar{\mathbf{w}}_{i_2}|, |\bar{\mathbf{v}}_{i_1}^\top \bar{\mathbf{w}}_{i_2}|, |\bar{\mathbf{v}}_{i_2}^\top \bar{\mathbf{w}}_{i_1}|, |\bar{\mathbf{v}}_{i_1}^\top \bar{\mathbf{w}}_{i_1}| \leq \varepsilon \\ & |\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j^*|, |\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j^*|, |\bar{\mathbf{v}}_i^\top \bar{\mathbf{w}}_j^*|, |\bar{\mathbf{w}}_i^\top \bar{\mathbf{v}}_j^*| \leq \varepsilon \end{aligned}$$

Theorem 2. Let $f_{\mathcal{S}}(\boldsymbol{\theta}, \mathbf{x})$, \mathcal{S}_τ , and τ^* be defined above. If Assumption 4 holds, then we have that

$$f_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}, \mathbf{x}) = \sum_{i=1}^{m^*} \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})$$

Theorem 2 states that, after τ^* steps of pruning, the resulting model contains the exact m^* router-expert

pairs with learned features from the teacher model. As a condition of Theorem 2, Assumption 4 is satisfied by Theorem 1 with $\varepsilon = \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right) \leq o\left(\frac{1}{\sqrt{m}}\right)$ since $d \gg m$. This implies that, if we perform the pruning at $T^* \leq t \leq T^* + \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ in the gradient flow process, we are guaranteed to remove all unused router-expert pairs and keep all necessary ones.

Notice that in the pruning procedure we evaluate the model on the population loss. To apply the algorithm in practice, one can effectively approximate the population loss with the sample loss. We use the population loss for the succinctness of the theoretical analysis.

Sketch of Proof. From a high level perspective, our proof relies on the observation that for two nonlinear function $h_1, h_2 : \mathbb{R} \rightarrow \mathbb{R}$ and vectors $\mathbf{u}_1, \mathbf{u}_2$ with $\mathbf{u}_1^\top \mathbf{u}_2 \approx 0$, we have that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}}[h_1(\mathbf{u}_1^\top \mathbf{x}) h_2(\mathbf{u}_2^\top \mathbf{x})] \\ \approx \mathbb{E}_{\mathbf{x}}[h_1(\mathbf{u}_1^\top \mathbf{x})] \mathbb{E}_{\mathbf{x}}[h_2(\mathbf{u}_2^\top \mathbf{x})] \end{aligned}$$

Let $q(\bar{\mathbf{v}}_i, \bar{\mathbf{w}}_i, \mathbf{x}) = \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})$. This allows us to approximate the loss as

$$\begin{aligned} \mathcal{L}(\theta) \approx \mathbb{E}_{\mathbf{x}} \left[\left(\sum_{i=1}^{m^*} (q(\bar{\mathbf{v}}_i, \bar{\mathbf{w}}_i, \mathbf{x}) - q(\bar{\mathbf{v}}_i^*, \bar{\mathbf{w}}_i^*, \mathbf{x})) \right)^2 \right] \\ + \sum_{i=m^*+1}^m \mathbb{E}_{\mathbf{x}} \left[\pi(\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \right] \mathbb{E}_{\mathbf{x}} \left[\sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] \end{aligned}$$

The first term is naturally small due to the fact that $\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_i^*$ and $\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_i^*$ are close to one for $i \in [m^*]$. The second term involves a summation of positive terms, which depends on the redundant router-expert pairs. Thus, removing each one of these will decrease $\mathcal{L}(\theta)$. The proof of Theorem 2 is provided in Appendix B.

5.2 Fine-Tuning the Pruned Model

Recall from Theorem 1 that, although the m^* router-expert pairs in the student model extracted near-perfect features from the teacher model, there is still an $\mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$ error for each $\bar{\mathbf{v}}_{\mathcal{I}(i)}^\top \bar{\mathbf{v}}_i^*$ and $\bar{\mathbf{w}}_{\mathcal{I}(i)}^\top \bar{\mathbf{w}}_i^*$. This results in a non-zero loss even after the pruning in Section 5.1. In this section, we study the convergence guarantee of fine-tuning the pruned model with gradient flow on the population MSE. In particular, we assume that $f(\theta, \mathbf{x})$ is the pruned model from Section 5.1 given by

$$f(\theta, \mathbf{x}) = \sum_{i=1}^{m^*} \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})$$

and the fine-tuning starts at $\theta(T_0)$ learned from Theorem 1 at time $T^* \leq T_0 \leq \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right)$. We slightly abuse

the notation by denoting

$$\theta = [\mathbf{v}_1^\top, \dots, \mathbf{v}_{m^*}^\top, \mathbf{w}_1^\top, \dots, \mathbf{w}_{m^*}^\top]^\top \in \mathbb{R}^{2m^*d}$$

Moreover, we denote the normalized version of θ as

$$\bar{\theta} = [\bar{\mathbf{v}}_1^\top, \dots, \bar{\mathbf{v}}_{m^*}^\top, \bar{\mathbf{w}}_1^\top, \dots, \bar{\mathbf{w}}_{m^*}^\top]^\top \in \mathbb{R}^{2m^*d}$$

The following theorem shows the convergence of gradient flow in the fine-tuning phase.

Theorem 3. *Let $\theta(T_0)$ that satisfy $\|\bar{\mathbf{v}}_i - \bar{\mathbf{v}}_i^*\|_2 \leq \varepsilon$ and $\|\bar{\mathbf{w}}_i - \bar{\mathbf{w}}_i^*\|_2 \leq \varepsilon$ for some $\varepsilon \leq o\left(\frac{1}{m^{*2}}\right)$. Let $C_{S,0} = 2\mathbb{E}_{x \sim \mathcal{N}(0,1)} \left[\pi(x)^2 \right]$ and $C_{S,1} = 6\mathbb{E}_{x \sim \mathcal{N}(0,1)} \left[\pi'(x)^2 \right]$. If $C_{S,0} \geq 1.1C_{S,1}$, then there exists some constant $\kappa > 0$ that only depends on the property of the sigmoid function $\pi(\cdot)$ such that*

$$\|\bar{\theta}(t + T_0) - \theta^*\|_2^2 \leq \exp\left(-\frac{\kappa t}{2}\right) \|\bar{\theta}(T_0) - \theta^*\|_2^2$$

Under the condition that $\bar{\mathbf{v}}_i$'s and $\bar{\mathbf{w}}_i$'s are ε -close to $\bar{\mathbf{v}}_i^*$'s and $\bar{\mathbf{w}}_i^*$'s, Theorem 3 shows a linear convergence rate in terms of the difference between the pruned model's normalized parameters $\bar{\theta}$ and the optimal parameters θ^* . In this fine-tuning stage, the convergence rate κ is independent of the dimension d or the number of experts m^* . Instead, it only depends on the property of the router's non-linear function $\pi(\cdot)$. Since $\theta(T_0)$ is given by the learned result in Theorem 1, the condition that $\|\bar{\mathbf{v}}_i - \bar{\mathbf{v}}_i^*\|_2 \leq \varepsilon$ and $\|\bar{\mathbf{w}}_i - \bar{\mathbf{w}}_i^*\|_2 \leq \varepsilon$ for some $\varepsilon \leq \mathcal{O}\left(\frac{m}{\sqrt{\delta_{\mathbb{P}} d^{\frac{1}{4}}}}\right) \leq o\left(\frac{1}{m^{*2}}\right)$ are automatically satisfied under $d \gg m$, since

$$\|\bar{\mathbf{v}}_i - \bar{\mathbf{v}}_i^*\|_2^2 = 2 - 2\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_i^* \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$$

The assumption that $C_{S,0} \geq 1.1C_{S,1}$ only depends on the property of $\pi(\cdot)$ and is checked in Appendix E.

Sketch of Proof. Our proof relies on the idea that, near the global minimum, the Hessian matrix is positive definite. In particular, we show that for any vector $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^{2m^*d}$ such that $\cos \langle \mathbf{u}_1, \mathbf{u}_2 \rangle \approx 1$, it holds that $\mathbf{u}_1^\top \nabla^2 \mathcal{L}(\theta) \mathbf{u}_2 \geq \kappa \|\mathbf{u}_1\|_2 \|\mathbf{u}_2\|$ for some constant $\kappa > 0$ and θ satisfying $\|\bar{\mathbf{v}}_i - \bar{\mathbf{v}}_i^*\|_2 \leq \varepsilon$ and $\|\bar{\mathbf{w}}_i - \bar{\mathbf{w}}_i^*\|_2 \leq \varepsilon$ with $\varepsilon \leq o\left(\frac{1}{m^{*2}}\right)$. Theorem 4 in Appendix C provides a formal statement of the result. Based on the positive-definiteness of the Hessian matrix, we leverage the classic convex optimization technique to show that the trajectory never leaves the neighborhood near the global minima, and that the distance to the global minimum converges linearly. The full proof is deferred to Appendix C.

6 CONCLUSION

Under the teacher-student set-up, we study the learning dynamics of the sigmoid-routed MoE with non-

linear experts defined in (1) when trained with gradient flow on the population MSE over high dimensional Gaussian inputs. In particular, our main result is a characterization of the feature learning stage, where proper features of the router-expert pairs are discovered in sequential order, with the expert’s recovery leading the router’s recovery. At the end of the feature learning stage, we show that a pruning procedure can be conducted to provably remove all the redundant experts and keep all necessary ones. Lastly, we show a linear convergence rate to the global minima for the the post-pruning fine-tuning with gradient flow. To the best of our knowledge, our work is the first to provide theoretical understanding on the joint training guarantee of MoEs with more than one activated experts and a general data assumption. In general, our paper is a further step into understanding the complicated dynamics of MoE training, and leads to the following open problems:

Online SGD and Sample Complexity. Due to the already sophisticated proof, our study is restricted to the setting of gradient flow on the population loss. However, as the main idea of the proof consists of an ODE based dynamic analysis, one could discretize the dynamic and apply martingale-based analysis to extend the theory to online SGD, as in Ren et al. (2025). This extension may lead to a sample complexity bound of learning m^* experts on d -dimensional data.

Experts with different importance. In our work we considered the teacher’s router and expert parameter $\bar{\mathbf{v}}_i^*$ and $\bar{\mathbf{w}}_i^*$ to be an orthonormal list. Due to the rotational invariance, this set-up puts equal importance to each router-expert pairs. Future work can investigate the scenario where the i th expert is scaled with a factor of α_i , and study the explicit ordering of the recovered experts in the student model.

Relax the dependency of d on m . Our current theory relies on the fact that $d \gg m$. While in the practical application of MoE we rarely set the number of experts to be larger than the input dimension, in most cases the scale of the two remains relatively the same. A meaningful future direction is to bridge the gap by studying the setting where d is only moderately larger than m .

7 Acknowledgements

This work was supported in part by NSF CAREER Award no. 2145629 and the Ken Kennedy Institute at Rice University. AK also acknowledges support from a Microsoft Research Award and an Amazon Research Award.

References

- G erard Ben Arous, Murat A. Erdogdu, N. Mert Vural, and Denny Wu. Learning quadratic neural networks in high dimensions: Sgd dynamics and scaling laws, 2025. URL <https://arxiv.org/abs/2508.03688>.
- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning in the presence of low-dimensional structure: A spiked random matrix perspective. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=H1IAoCHDWW>.
- Umang Bhatt, Sanyam Kapoor, Mihir Upadhyay, Ilya Sucholutsky, Francesco Quinzan, Katherine M. Collins, Adrian Weller, Andrew Gordon Wilson, and Muhammad Bilal Zafar. When should we orchestrate multiple agents?, 2025. URL <https://arxiv.org/abs/2503.13577>.
- Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with shallow neural networks, 2022. URL <https://arxiv.org/abs/2210.15651>.
- Alberto Bietti, Joan Bruna, and Loucas Pillaud-Vivien. On learning gaussian multi-index models with gradient flow, 2023. URL <https://arxiv.org/abs/2310.19793>.
- Enric Boix-Adsera and Philippe Rigollet. The power of fine-grained experts: Granularity boosts expressivity in mixture of experts, 2025. URL <https://arxiv.org/abs/2505.06839>.
- Etienne Boursier and Nicolas Flammarion. Early alignment in two-layer networks training is a two-edged sword. *Journal of Machine Learning Research*, 26(183):1–75, 2025.
- Etienne Boursier, Loucas Pillaud-Vivien, and Nicolas Flammarion. Gradient flow dynamics of shallow ReLU networks for square loss and orthogonal inputs. In *Advances in Neural Information Processing Systems*, volume 35, pages 20105–20118, 2022.
- Zixiang Chen, Yihe Deng, Yue Wu, Quanquan Gu, and Yuanzhi Li. Towards understanding mixture of experts in deep learning, 2022. URL <https://arxiv.org/abs/2208.02813>.
- Mohammed Nowaz Rabbani Chowdhury, Shuai Zhang, Meng Wang, Sijia Liu, and Pin-Yu Chen. Patch-level routing in mixture-of-experts is provably sample-efficient for convolutional neural networks, 2023. URL <https://arxiv.org/abs/2306.04073>.
- Mohammed Nowaz Rabbani Chowdhury, Meng Wang, Kaoutar El Maghraoui, Naigang Wang, Pin-Yu Chen, and Christopher Carothers. A provably effective method for pruning experts in fine-tuned sparse

mixture-of-experts, 2024. URL <https://arxiv.org/abs/2405.16646>.

Alex Damian, Jason D. Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent, 2022. URL <https://arxiv.org/abs/2206.15144>.

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanbiao Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu,

Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025. URL <https://arxiv.org/abs/2412.19437>.

Simon S. Du, Jason D. Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks, 2019. URL <https://arxiv.org/abs/1811.03804>.

William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.

Quentin Fruytier, Aryan Mokhtari, and Sujay Sanghavi. Learning mixtures of experts with em: A mirror descent perspective, 2025. URL <https://arxiv.org/abs/2411.06056>.

Rong Ge, Jason D. Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design, 2017. URL <https://arxiv.org/abs/1711.00501>.

Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems, 2025. URL <https://arxiv.org/abs/2408.08435>.

Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks, 2020. URL <https://arxiv.org/abs/1806.07572>.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

Ryotaro Kawata, Kohsei Matsutani, Yuri Kinoshita, Naoki Nishikawa, and Taiji Suzuki. Mixture of experts provably detect and learn the latent cluster structure in gradient-based learning, 2025. URL <https://arxiv.org/abs/2506.01656>.

Yilun Kong, Guozheng Ma, Qi Zhao, Haoyu Wang, Li Shen, Xueqian Wang, and Dacheng Tao. Mastering massive multi-task reinforcement learning via mixture-of-expert decision transformer, 2025. URL <https://arxiv.org/abs/2505.24378>.

Anastasis Kratsios, Takashi Furuya, Jose Antonio Lara Benitez, Matti Lassas, and Maarten de Hoop. Mixture of experts soften the curse of dimensionality in operator learning, 2024. URL <https://arxiv.org/abs/2404.09101>.

Jason D. Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns low-dimensional polynomials with sgd near the information-theoretic limit, 2024. URL <https://arxiv.org/abs/2406.01581>.

- Binghui Li et al. Feature averaging: An implicit bias of gradient descent leading to non-robustness in neural networks. *arXiv preprint arXiv:2410.10322*, 2024.
- Hongbo Li, Sen Lin, Lingjie Duan, Yingbin Liang, and Ness B. Shroff. Theory on mixture-of-experts in continual learning, 2025. URL <https://arxiv.org/abs/2406.16437>.
- Xudong Lu, Qi Liu, Yuhui Xu, Aojun Zhou, Siyuan Huang, Bo Zhang, Junchi Yan, and Hongsheng Li. Not all experts are equal: Efficient expert pruning and skipping for mixture-of-experts large language models, 2024. URL <https://arxiv.org/abs/2402.14800>.
- Tao Luo et al. Phase diagram for two-layer ReLU neural networks at infinite-width limit. *Journal of Machine Learning Research*, 22(71):1–47, 2021.
- Hartmut Maennel, Olivier Bousquet, and Sylvain Gelly. Gradient Descent Quantizes ReLU network features. *arXiv preprint arXiv:1803.08367*, 2018.
- Hancheng Min, Enrique Mallada, and Rene Vidal. Early neuron alignment in two-layer ReLU networks with small initialization. *arXiv preprint arXiv:2307.12851*, 2023.
- Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A. Erdogdu. Neural networks efficiently learn low-dimensional representations with sgd, 2023a. URL <https://arxiv.org/abs/2209.14863>.
- Alireza Mousavi-Hosseini, Denny Wu, Taiji Suzuki, and Murat A. Erdogdu. Gradient-based feature learning under structured data, 2023b. URL <https://arxiv.org/abs/2309.03843>.
- Basil Mustafa, Carlos Riquelme, Joan Puigcerver, Rodolphe Jenatton, and Neil Houlsby. Multimodal contrastive learning with LIMoE: the language-image mixture of experts. *Advances in Neural Information Processing Systems*, 35:9564–9576, 2022.
- Huy Nguyen, Pedram Akbarian, TrungTin Nguyen, and Nhat Ho. A general theory for softmax gating multinomial logistic mixture of experts, 2024a. URL <https://arxiv.org/abs/2310.14188>.
- Huy Nguyen, Nhat Ho, and Alessandro Rinaldo. Sigmoid gating is more sample efficient than softmax gating in mixture of experts, 2024b. URL <https://arxiv.org/abs/2405.13997>.
- Huy Nguyen, Nhat Ho, and Alessandro Rinaldo. On least square estimation in softmax gating mixture of experts, 2024c. URL <https://arxiv.org/abs/2402.02952>.
- Huy Nguyen, Nhat Ho, and Alessandro Rinaldo. Convergence rates for softmax gating mixture of experts, 2025. URL <https://arxiv.org/abs/2503.03213>.
- Yunwei Ren, Eshaan Nichani, Denny Wu, and Jason D. Lee. Emergence and scaling laws in sgd learning of shallow neural networks, 2025. URL <https://arxiv.org/abs/2504.19983>.
- Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017.
- Zhenmei Shi, Junyi Wei, and Yingyu Liang. A theoretical analysis on feature learning in neural networks: Emergence from inputs and advantage over fixed features, 2022. URL <https://arxiv.org/abs/2206.01717>.
- Zhenmei Shi, Junyi Wei, and Yingyu Liang. Provable guarantees for neural networks via gradient feature learning, 2023. URL <https://arxiv.org/abs/2310.12408>.
- Mingze Wang and Weinan E. On the expressive power of mixture-of-experts for structured complex tasks, 2025. URL <https://arxiv.org/abs/2505.24205>.
- Xiang Wang, Chenwei Wu, Jason D. Lee, Tengyu Ma, and Rong Ge. Beyond lazy training for over-parameterized tensor decomposition, 2020. URL <https://arxiv.org/abs/2010.11356>.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Danyang Zhang, Junhao Song, Ziqian Bi, Yingfang Yuan, Tianyang Wang, Joe Yeong, and Junfeng Hao. Mixture of experts in large language models, 2025a. URL <https://arxiv.org/abs/2507.11181>.

Guibin Zhang, Yanwei Yue, Xiangguo Sun, Guancheng Wan, Miao Yu, Junfeng Fang, Kun Wang, Tianlong Chen, and Dawei Cheng. G-designer: Architecting multi-agent communication topologies via graph neural networks, 2025b. URL <https://arxiv.org/abs/2410.11782>.

Ruichen Zhang, Mufan Qiu, Zhen Tan, Mohan Zhang, Vincent Lu, Jie Peng, Kaidi Xu, Leandro Z. Agudelo, Peter Qian, and Tianlong Chen. Symbiotic cooperation for web agents: Harnessing complementary strengths of large and small llms, 2025c. URL <https://arxiv.org/abs/2502.07942>.

Zeliang Zhang, Xiaodong Liu, Hao Cheng, Chenliang Xu, and Jianfeng Gao. Diversifying the expert knowledge for task-agnostic pruning in sparse mixture-of-experts, 2025d. URL <https://arxiv.org/abs/2407.09590>.

Berfin Şimşek, Amire Bendjeddou, and Daniel Hsu. Learning gaussian multi-index models with gradient flow: Time complexity and directional convergence, 2025. URL <https://arxiv.org/abs/2411.08798>.

Checklist

1. For all models and algorithms presented, check if you include:
 - (a) A clear description of the mathematical setting, assumptions, algorithm, and/or model. [Yes/No/Not Applicable] **Yes**
 - (b) An analysis of the properties and complexity (time, space, sample size) of any algorithm. [Yes/No/Not Applicable] **Yes**
 - (c) (Optional) Anonymized source code, with specification of all dependencies, including external libraries. [Yes/No/Not Applicable] **Yes**
2. For any theoretical claim, check if you include:
 - (a) Statements of the full set of assumptions of all theoretical results. [Yes/No/Not Applicable] **Yes**
 - (b) Complete proofs of all theoretical results. [Yes/No/Not Applicable] **Yes**
 - (c) Clear explanations of any assumptions. [Yes/No/Not Applicable] **Yes**
3. For all figures and tables that present empirical results, check if you include:
 - (a) The code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL). [Yes/No/Not Applicable] **Yes**
 - (b) All the training details (e.g., data splits, hyperparameters, how they were chosen). [Yes/No/Not Applicable] **Yes**
 - (c) A clear definition of the specific measure or statistics and error bars (e.g., with respect to the random seed after running experiments multiple times). [Yes/No/Not Applicable] **No**
 - (d) A description of the computing infrastructure used. (e.g., type of GPUs, internal cluster, or cloud provider). [Yes/No/Not Applicable] **No**
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets, check if you include:
 - (a) Citations of the creator If your work uses existing assets. [Yes/No/Not Applicable] **Not Applicable**
 - (b) The license information of the assets, if applicable. [Yes/No/Not Applicable] **Not Applicable**
 - (c) New assets either in the supplemental material or as a URL, if applicable. [Yes/No/Not Applicable] **Not Applicable**
 - (d) Information about consent from data providers/curators. [Yes/No/Not Applicable] **Not Applicable**
 - (e) Discussion of sensible content if applicable, e.g., personally identifiable information or offensive content. [Yes/No/Not Applicable] **Not Applicable**
5. If you used crowdsourcing or conducted research with human subjects, check if you include:
 - (a) The full text of instructions given to participants and screenshots. [Yes/No/Not Applicable] **Not Applicable**
 - (b) Descriptions of potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable. [Yes/No/Not Applicable] **Not Applicable**
 - (c) The estimated hourly wage paid to participants and the total amount spent on participant compensation. [Yes/No/Not Applicable] **Not Applicable**

Instructions for Paper Submissions to AISTATS 2026: Supplementary Materials

A Proof of Theorem 1

A.1 Proof Outline

Initialization Property. At initialization, the following property needs to be satisfied.

Condition 1 (Initialization). *At initialization $\{\mathbf{w}_i(0)\}_{i=1}^m$ and $\{\mathbf{v}_i(0)\}_{i=1}^m$ satisfies*

- $\mathbf{w}_{i_\ell^*}(0)^\top \bar{\mathbf{w}}_{j_\ell^*}^* \geq (1 + 2\delta_s) \mathbf{w}_{i_\ell^*}(0)^\top \bar{\mathbf{w}}_{j_\ell^*}^*$ for all $\ell \in [m^*]$ and $j \in [m^*] \setminus \mathcal{C}_{\ell+1}$.
- $\mathbf{w}_{i_\ell^*}(0)^\top \bar{\mathbf{w}}_{j_\ell^*}^* \geq (1 + 2\delta_s) \mathbf{w}_i(0)^\top \bar{\mathbf{w}}_{j_\ell^*}^*$ for all $\ell \in [m^*]$ and $i \in [m] \setminus \mathcal{R}_{\ell+1}$.
- $\mathbf{w}_{i_\ell^*}(0)^\top \bar{\mathbf{w}}_{j_{\ell+1}^*}^* \geq (1 + 2\delta_s) \mathbf{w}_{i_{\ell+1}^*}(0)^\top \bar{\mathbf{w}}_{j_{\ell+1}^*}^*$ for all $\ell \in [m^* - 1]$.
- $\left(\mathbf{w}_{i_\ell^*}(0)^\top \bar{\mathbf{w}}_{j_\ell^*}^* \right)^2 \geq \frac{\log m^*}{d}$ for all $\ell \in [m^*]$.
- $\|\mathbf{w}_i\|_2, \|\mathbf{v}_i\|_2 \in [1 - \beta_2\delta_s, 1 + \beta_2\delta_s]$ for all $i \in [m]$.
- $\max \left\{ (\mathbf{v}_i(0)^\top \mathbf{v}_j^*)^2, (\mathbf{v}_i(0)^\top \mathbf{w}_j^*)^2, (\mathbf{v}_i(0)^\top \mathbf{w}_j^*)^2, (\mathbf{w}_i(0)^\top \mathbf{v}_j^*)^2 \right\} \leq \frac{\beta_3}{d} \log \frac{m}{\delta_{\mathbb{P}}}$ for all $i \in [m], j \in [m^*]$.
- $\max \left\{ (\mathbf{v}_i(0)^\top \mathbf{v}_j(0))^2, (\mathbf{w}_i(0)^\top \mathbf{w}_j(0))^2, (\mathbf{v}_i(0)^\top \mathbf{w}_j(0))^2 \right\} \leq \frac{\beta_3}{d} \log \frac{m}{\delta_{\mathbb{P}}}$ for all $i, j \in [m]$ and $i \neq j$. Moreover, $\mathbf{v}_i(0)^\top \mathbf{w}_i(0) = 0$ for all $i \in [m]$.

where $\delta_s = \frac{\beta_1\delta_{\mathbb{P}}}{mm^*}$ for some absolute constant $\beta_1, \beta_3 > 0$ and $\beta_2 \leq o(1)$ and any $\delta_{\mathbb{P}} \in (0, 1/\tau)$.

By Lemma 2 and Lemma 3, the above condition holds with probability at least $1 - 7\delta_{\mathbb{P}}$ as long as $d \geq \frac{\beta_5 m^4}{\delta_{\mathbb{P}}^2} \log \frac{m}{\delta_{\mathbb{P}}}$ and $m \geq \beta_4 m^* \log \frac{m^*}{\delta_{\mathbb{P}}}$.

Inductive Hypothesis. Now we are going to show that $\bar{\mathbf{w}}_{i_\ell^*}(t)^\top \bar{\mathbf{w}}_{j_\ell^*}^*$ converges to at least $1 - \frac{c}{\sqrt{d}}$ for all $\ell \in [m^*]$ by induction. To start, we denote the values of interest as follows

$$\begin{aligned} \gamma_{i,j}^{(1)}(t) &= \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j^*; & \gamma_{i,j}^{(2)}(t) &= \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^*; & \zeta_{i,j}^{(1)}(t) &= \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{w}}_j^*; & \zeta_{i,j}^{(2)}(t) &= \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{v}}_j^* \\ I_{i,j}^{(1)}(t) &= \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j(t); & I_{i,j}^{(2)}(t) &= \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j(t); & I_{i,j}^{(3)}(t) &= \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{w}}_j(t) \end{aligned} \quad (9)$$

To state the inductive hypothesis, we need the following error bounds.

Definition 1 (Error Bounds). *For each $\ell \in [m^*]$, we define the following:*

$$\begin{aligned} \varepsilon_{1,\ell}(t) &:= \max_{i \in [m] \setminus \mathcal{R}_\ell, j \in [m^*]} \left| \gamma_{i,j}^{(1)}(t) \right|; \quad \varepsilon_{2,\ell}(t) := \max_{j \in [m^*] \setminus \{j_\ell^*\}} \left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|; \\ \varepsilon_{3,\ell}(t) &:= \max_{i \in [m] \setminus \mathcal{R}_\ell, j \in [m^*]} \left| \gamma_{i,j}^{(2)}(t) \right|; \quad \varepsilon_{4,\ell}(t) := \max_{j \in [m^*] \setminus \{j_\ell^*\}} \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right|; \quad \varepsilon_{5,\ell}(t) := \left| I_{i_\ell^*,i_\ell^*}^{(3)}(t) \right| \end{aligned}$$

Moreover, we also define the forward error, the backward error, and the aggregated error as

$$\begin{aligned} \varepsilon_{\mathcal{F},\ell}(t) &= \max \{ \varepsilon_{\mathcal{F},\ell,1}(t), \varepsilon_{\mathcal{F},\ell,2}(t) \}; \quad \varepsilon_{\mathcal{B},\ell}^{(1)}(t) = \max \{ \varepsilon_{\mathcal{B},\ell,1}(t), \varepsilon_{\mathcal{B},\ell,2}(t), \varepsilon_{2,\ell}(t) \} \\ \varepsilon_{\mathcal{B},\ell}^{(2)}(t) &= \max \left\{ \varepsilon_{\mathcal{B},\ell}^{(1)}(t), \varepsilon_{1,\ell}(t) \right\}; \quad \varepsilon_{\mathcal{B},\ell}^{(3)}(t) = \max_{i \in [m] \setminus \mathcal{R}_\ell} \left| I_{i,i}^{(3)}(t) \right| \\ \hat{\varepsilon}_{\mathcal{A},\ell}^{(1)}(t) &= \max \left\{ \varepsilon_{4,\ell}(t), \varepsilon_{\mathcal{B},\ell}^{(1)}(t), \varepsilon_{\mathcal{F},\ell}(t) \right\}; \quad \hat{\varepsilon}_{\mathcal{A},\ell}^{(2)}(t) = \max \left\{ \hat{\varepsilon}_{\mathcal{A},\ell}^{(1)}(t), \varepsilon_{1,\ell}(t), \varepsilon_{3,\ell}(t) \right\} \end{aligned}$$

where $\varepsilon_{\mathcal{F},\ell,1}(t), \varepsilon_{\mathcal{F},\ell,2}(t), \varepsilon_{\mathcal{F},\ell,3}(t)$ and $\varepsilon_{\mathcal{B},\ell,1}(t), \varepsilon_{\mathcal{B},\ell,2}(t)$ are defined as

$$\begin{aligned}\varepsilon_{\mathcal{F},\ell,1}(t) &:= \max_{\ell' \leq \ell, j \in [m^*] \setminus \{j_{\ell'}^*\}} \max \left\{ \left| \gamma_{i_{\ell'}, j}^{(1)}(t) \right|, \left| \gamma_{i_{\ell'}, j}^{(2)}(t) \right| \right\}; \\ \varepsilon_{\mathcal{F},\ell,2}(t) &:= \max_{i \in \mathcal{R}_{\ell-1}, j \in [m^*]} \max \left\{ \left| \zeta_{i,j}^{(1)}(t) \right|, \left| \zeta_{i,j}^{(2)}(t) \right| \right\} \\ \varepsilon_{\mathcal{B},\ell,1}(t) &:= \max_{i \in [m] \setminus \mathcal{R}_{\ell-1}, j \in [m^*]} \max \left\{ \left| \zeta_{i,j}^{(1)}(t) \right|, \left| \zeta_{i,j}^{(2)}(t) \right| \right\} \\ \varepsilon_{\mathcal{B},\ell,2}(t) &:= \max_{i,j \in [m], i \neq j} \max \left\{ \left| I_{i,j}^{(1)}(t) \right|, \left| I_{i,j}^{(2)}(t) \right|, \left| I_{i,j}^{(3)}(t) \right| \right\}\end{aligned}$$

Lastly, we are going to define the monotonic upper bound of $\hat{\varepsilon}_{\mathcal{A},\ell}^{(1)}(t)$ and $\hat{\varepsilon}_{\mathcal{A},\ell}^{(2)}(t)$

$$\varepsilon_{\mathcal{A},\ell}^{(1)}(t) = \sup_{t' \in [0,t]} \hat{\varepsilon}_{\mathcal{A},\ell}^{(1)}(t'); \quad \varepsilon_{\mathcal{A},\ell}^{(2)}(t) = \sup_{t' \in [0,t]} \hat{\varepsilon}_{\mathcal{A},\ell}^{(2)}(t')$$

Definition 2 (Recovery Time). Define the ξ -recovery time of $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)$, denoted as $T_{\ell}(\xi)$, as

$$T_{\ell}(\xi) = \min \left\{ t \geq 0 : \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) \geq \xi \right\}$$

Based on $T_{\ell}(\xi)$, we define the **constant-recovery time** and the **near-perfect-recovery time** as $T_{r,\ell} = T_{\ell}(0.9)$ and $T_{p,\ell} = T_{\ell}\left(1 - \frac{\beta_9 m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$, respectively.

Condition 2 (Inductive Hypothesis). Let $\ell \in [m^*]$. Then we have that

- (Sequential recovery) For all $t \geq T_{p,\ell-1}$ such that $\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ we have $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) \geq 1 - \frac{\beta_9 m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}$, and $\gamma_{i_{\ell'}^*, j_{\ell'}^*}^{(1)}(t) \geq 1 - \frac{\beta_9 m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}$ for all $\ell' < \ell$.
- (Error bound of $\gamma_{i,j}^{(1)}, \gamma_{i,j}^{(2)}$) For $t \leq T_{p,\ell-1}$, we have

$$\max_{i \in [m] \setminus \mathcal{R}_{\ell-1}, j \in [m^*]} \max \left\{ \left| \gamma_{i,j}^{(1)}(t) \right|, \left| \gamma_{i,j}^{(2)}(t) \right| \leq \frac{\beta_6 m^2}{\sqrt{d}} \right\}$$

- (Error bound of remaining items) $\varepsilon_{\mathcal{F},\ell}(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ for all t such that $\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$.

Here $\beta_6, \beta_9 > 0$ are some absolute constant.

The proof proceeds by establishing the inductive hypothesis.

A.2 Initialization Property

Lemma 3. Let $\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_m$ and $\mathbf{w}_1, \dots, \mathbf{w}_m$ be I.I.D. random vectors from $\mathcal{N}(\mathbf{0}, d^{-1}\mathbf{I}_d)$. Define $\mathbf{v}_i = \left(\mathbf{I} - \frac{1}{\|\mathbf{w}_i\|_2^2} \mathbf{w}_i \mathbf{w}_i^\top\right) \hat{\mathbf{v}}_i$. Then there exists some absolute constant $\beta_1, \beta_3, \beta_5 > 0, \beta_2 \leq o(1)$, and $\delta_{\mathbb{P}} \in (0, 1/3)$ such that if $d \geq \frac{\beta_5 m^4}{\delta_{\mathbb{P}}^2} \log \frac{m}{\delta_{\mathbb{P}}}$ and $\delta_s = \frac{\beta_1 \delta_{\mathbb{P}}}{m^2}$, with probability at least $1 - 3\delta_{\mathbb{P}}$ we have that

- $\|\mathbf{v}_i\|_2^2, \|\mathbf{w}_i\|_2^2 \in [1 - \beta_2 \delta_s, 1 + \beta_2 \delta_s]$ for all $i \in [m]$;
- $\max \{ \mathbf{v}_i[j]^2, \mathbf{w}_i[j]^2 \} \leq \frac{\beta_3}{d} \log \frac{m}{\delta_{\mathbb{P}}}$ for all $i \in [m], j \in [m^*]$;
- $\max \left\{ (\mathbf{v}_i^\top \mathbf{v}_j)^2, (\mathbf{w}_i^\top \mathbf{w}_j)^2, (\mathbf{v}_i^\top \mathbf{w}_j)^2 \right\} \leq \frac{\beta_3}{d} \log \frac{m}{\delta_{\mathbb{P}}}$ for all $i, j \in [m]$ with $i \neq j$.

Proof. Our proof starts with showing the concentration for $\|\hat{\mathbf{v}}_i\|_2^2$, $\|\mathbf{w}_i\|_2^2$ and $\hat{\mathbf{v}}_i^\top \mathbf{w}_i$, and then moves to the proof of the desired statement.

Concentration of $\|\hat{\mathbf{v}}_i\|_2^2$, $\|\mathbf{w}_i\|_2^2$. Let $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, d^{-1}\mathbf{I}_d)$. Then $\mathbf{v}[i] \sim \mathcal{N}(0, d^{-1})$. Thus, $d\mathbf{v}[i]^2 \in \text{subE}(2, 2)$. Since $\mathbf{v}[i]$'s are I.I.D., we have that

$$d\|\mathbf{v}\|_2^2 = d \sum_{i=1}^d \mathbf{v}[i]^2 \in \text{subE}(2d, 2)$$

Therefore, by the tail bound of sub-exponential random variables, we have that

$$\Pr\left(\left|\|\mathbf{v}\|_2^2 - 1\right| \geq \frac{t}{d}\right) = \Pr\left(\left|d\|\mathbf{v}\|_2^2 - d\mathbb{E}\left[\|\mathbf{v}\|_2^2\right]\right| \geq t\right) \leq 2\exp\left(-\frac{1}{4} \min\left\{\frac{t^2}{d}, t\right\}\right)$$

We are going to focus on the case where $t \leq d$. In particular, we set $t = \frac{1}{2}\beta_2\delta_s d$. Then we have that

$$\Pr\left(\left|\|\mathbf{v}\|_2^2 - 1\right| \geq \frac{1}{2}\beta_2\delta_s\right) \leq 2\exp\left(-\frac{1}{16}\beta_2^2\delta_s^2 d\right)$$

Take a union bound over all $i \in [m]$ for $\hat{\mathbf{v}}_i$ and \mathbf{w}_i gives that, with probability at least $1 - 4m\exp\left(-\frac{1}{16}\beta_2^2\delta_s^2 d\right)$, it holds that

$$\|\hat{\mathbf{v}}_i\|_2^2, \|\mathbf{w}_i\|_2^2 \in \left[1 - \frac{1}{2}\beta_2\delta_s, 1 + \frac{1}{2}\beta_2\delta_s\right]$$

Setting $d \geq \frac{\beta_5 m^4}{\delta_{\mathbb{P}}^2} \log \frac{m}{\delta_{\mathbb{P}}} \geq \frac{16}{\beta_2^2 \delta_s^2} \log \frac{4m}{\delta_{\mathbb{P}}}$ guarantees that the failing probability is within $\delta_{\mathbb{P}}$.

Concentration of $\hat{\mathbf{v}}_i^\top \mathbf{w}_j$, $\hat{\mathbf{v}}_i^\top \hat{\mathbf{v}}_j$, and $\mathbf{w}_i^\top \mathbf{w}_j$. Due to the independence between $\hat{\mathbf{v}}_i$'s and \mathbf{w}_j 's, we have that

$$\mathbb{E}[\hat{\mathbf{v}}_i^\top \mathbf{w}_j] = 0; \quad d\hat{\mathbf{v}}_i[\ell]\mathbf{w}_j[\ell] \in \text{subE}(1, 1).$$

Thus, we have that $\hat{\mathbf{v}}_i^\top \mathbf{w}_j \in \text{subE}(d, 1)$. Applying the tail bound of sub-exponential random variable gives that

$$\Pr\left(\left|\hat{\mathbf{v}}_i^\top \mathbf{w}_j\right| \geq \frac{t}{d}\right) = \Pr\left(d\left|\hat{\mathbf{v}}_i^\top \mathbf{w}_j\right| \geq t\right) \leq 2\exp\left(-\frac{1}{2} \min\left\{\frac{t^2}{d}, t\right\}\right)$$

The same concentration also holds for $\hat{\mathbf{v}}_i^\top \hat{\mathbf{v}}_j$ and $\mathbf{w}_i^\top \mathbf{w}_j$. Again we are going to focus on the case $t \leq d$. Take a union bound over all $i, j \in [m]$ and $\hat{\mathbf{v}}_i, \mathbf{w}_i$ gives that

$$\Pr\left(\max\left\{\left|\hat{\mathbf{v}}_i^\top \mathbf{w}_j\right|, \left|\hat{\mathbf{v}}_i^\top \hat{\mathbf{v}}_j\right|, \left|\mathbf{w}_i^\top \mathbf{w}_j\right|\right\} \geq \frac{t}{d}; \forall i, j \in [m]\right) \leq 8m^2 \exp\left(-\frac{t^2}{2d}\right)$$

Setting the failing probability to $\delta_{\mathbb{P}}$ gives that with probability at least $1 - \delta_{\mathbb{P}}$, we have that

$$\hat{\mathbf{v}}_i^\top \mathbf{w}_j, \hat{\mathbf{v}}_i^\top \hat{\mathbf{v}}_j, \mathbf{w}_i^\top \mathbf{w}_j \in \left[-\frac{1}{\sqrt{d}} \left(\log \frac{8m^2}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}, \frac{1}{\sqrt{d}} \left(\log \frac{8m^2}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}\right]$$

Proof of the first statement. Notice that the bound for $\|\mathbf{w}_i\|_2^2$ is already implied by its concentration property. To prove the bound for $\|\mathbf{v}_i\|_2^2$, we write

$$\|\mathbf{v}_i\|_2^2 = \left\| \left(\mathbf{I} - \frac{1}{\|\mathbf{w}_i\|_2^2} \mathbf{w}_i \mathbf{w}_i^\top \right) \hat{\mathbf{v}}_i \right\|_2^2 = \|\hat{\mathbf{v}}_i\|_2^2 - \frac{1}{\|\mathbf{w}_i\|_2^2} (\mathbf{w}_i^\top \hat{\mathbf{v}}_i)^2$$

By the concentration property of $\|\mathbf{w}_i\|_2^2$ and $\mathbf{w}_i^\top \hat{\mathbf{v}}_j$, we have that

$$\frac{1}{\|\mathbf{w}_i\|_2^2} (\mathbf{w}_i^\top \hat{\mathbf{v}}_i)^2 \leq \frac{1}{d(1 - \frac{1}{2}\beta_s\delta_s)} \log \frac{8m^2}{\delta_{\mathbb{P}}}$$

For $d \geq \frac{\beta_3 m^2}{\delta_{\mathbb{P}}} \log \frac{m}{\delta_{\mathbb{P}}} \geq \frac{2}{\beta_2 \delta_s (1 - \frac{1}{2} \beta_2 \delta_s)} \log \frac{8m^2}{\delta_{\mathbb{P}}}$, we have that $\frac{1}{\|\mathbf{w}_i\|_2^2} (\mathbf{w}_i^\top \hat{\mathbf{v}}_i)^2 \leq \frac{1}{2} \beta_2 \delta_s$. Combined with the concentration property of $\hat{\mathbf{v}}_i$, we have that $\|\mathbf{v}_i\|_2^2 \in [1 - \beta_2 \delta_s, 1 + \beta_2 \delta_s]$.

Proof of second statement. By the tail bound of Gaussian random variable, we have that for all $z \sim \mathcal{N}(0, 1)$, it holds that

$$\Pr(z^2 \geq t) = 2 \Pr\left(z \geq \sqrt{t}\right) \leq \exp\left(-\frac{t}{2}\right)$$

Apply the above to $z = \sqrt{d} \cdot \mathbf{w}_i[j]$ and $\sqrt{d} \cdot \hat{\mathbf{v}}_i[j]$ with a union bound over all $i \in [m]$ and $j \in [m^*]$ gives

$$\Pr\left(\max_{i \in [m], j \in [m^*]} \hat{\mathbf{v}}_i[j]^2 \geq \frac{t}{d}; \max_{i \in [m], j \in [m^*]} \mathbf{w}_i[j]^2 \geq \frac{t}{d}\right) \leq 2mm^* \exp\left(-\frac{t}{2}\right)$$

Set $t = \frac{\beta_3}{4} \log \frac{m}{\delta_{\mathbb{P}}} \geq 2 \log \frac{2mm^*}{\delta_{\mathbb{P}}}$ gives the desired result for $\mathbf{w}_i[j]^2$ s. To bound $\mathbf{v}_i[j]^2$, we notice that

$$|\mathbf{v}_i[j]| \leq |\hat{\mathbf{v}}_i[j]| + \frac{1}{\|\mathbf{w}_i\|_2} |\mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_i[j]|$$

By the previous bounds, we have that

$$|\hat{\mathbf{v}}_i[j]| \leq \frac{\sqrt{\beta_3}}{2\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right); \frac{1}{\|\mathbf{w}_i\|_2} |\mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_i[j]| \leq \frac{1}{\sqrt{d} (1 - \frac{1}{2} \beta_2 \delta_s)} \left(\log \frac{8m^2}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}} \cdot \frac{\sqrt{\beta_3}}{2\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)$$

With $d \geq \frac{\beta_5 m^4}{\delta_{\mathbb{P}}} \log \frac{m}{\delta_{\mathbb{P}}}$ we can guarantee that the latter is also upper bounded by $\frac{\sqrt{\beta_3}}{2\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)$. Combining the two bounds and square both sides gives the desired result for $\mathbf{v}_i[j]^2$.

Proof of the third statement. The bound of $\mathbf{w}_i^\top \mathbf{w}_j$ is again implied by the concentration we showed above. To show the bound of $\mathbf{v}_i^\top \mathbf{v}_j$, we write

$$\begin{aligned} \mathbf{v}_i^\top \mathbf{v}_j &= \hat{\mathbf{v}}_i^\top \left(\mathbf{I} - \frac{1}{\|\mathbf{w}_i\|_2^2} \mathbf{w}_i \mathbf{w}_i^\top \right) \left(\mathbf{I} - \frac{1}{\|\mathbf{w}_j\|_2^2} \mathbf{w}_j \mathbf{w}_j^\top \right) \hat{\mathbf{v}}_j \\ &= \hat{\mathbf{v}}_i^\top \hat{\mathbf{v}}_j - \frac{1}{\|\mathbf{w}_i\|_2^2} \mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_i^\top \hat{\mathbf{v}}_j - \frac{1}{\|\mathbf{w}_j\|_2^2} \mathbf{w}_j^\top \hat{\mathbf{v}}_j \cdot \mathbf{w}_j^\top \hat{\mathbf{v}}_i + \frac{1}{\|\mathbf{w}_i\|_2^2 \|\mathbf{w}_j\|_2^2} \mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_j^\top \hat{\mathbf{v}}_j \cdot \mathbf{w}_i^\top \mathbf{w}_j \end{aligned}$$

By the concentration of the norms and inner-products, we have that

$$\begin{aligned} \frac{1}{\|\mathbf{w}_i\|_2^2} |\mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_i^\top \hat{\mathbf{v}}_j|, \frac{1}{\|\mathbf{w}_j\|_2^2} |\mathbf{w}_j^\top \hat{\mathbf{v}}_j \cdot \mathbf{w}_j^\top \hat{\mathbf{v}}_i| &\leq \frac{1}{d (1 - \frac{1}{2} \beta_s \delta_s)} \log \frac{8m^2}{\delta_{\mathbb{P}}} \\ \frac{1}{\|\mathbf{w}_i\|_2^2 \|\mathbf{w}_j\|_2^2} |\mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_j^\top \hat{\mathbf{v}}_j \cdot \mathbf{w}_i^\top \mathbf{w}_j| &\leq \frac{1}{d^{\frac{3}{2}} (1 - \frac{1}{2} \beta_s \delta_s)^2} \left(\log \frac{8m^2}{\delta_{\mathbb{P}}}\right)^{\frac{3}{2}} \end{aligned}$$

With the condition $d \geq \frac{\beta_5 m^2}{\delta_{\mathbb{P}}} \log \frac{m}{\delta_{\mathbb{P}}}$ we have that

$$\frac{1}{d (1 - \frac{1}{2} \beta_s \delta_s)^2} \log \frac{8m^2}{\delta_{\mathbb{P}}} \leq \frac{1}{6\sqrt{d}} \left(\log \frac{8m^2}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}} \leq \frac{\sqrt{\beta_3}}{6\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}} \leq 1$$

Therefore, we can conclude that

$$\frac{1}{\|\mathbf{w}_i\|_2^2} |\mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_i^\top \hat{\mathbf{v}}_j|, \frac{1}{\|\mathbf{w}_j\|_2^2} |\mathbf{w}_j^\top \hat{\mathbf{v}}_j \cdot \mathbf{w}_j^\top \hat{\mathbf{v}}_i|, \frac{1}{\|\mathbf{w}_i\|_2^2 \|\mathbf{w}_j\|_2^2} |\mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_j^\top \hat{\mathbf{v}}_j \cdot \mathbf{w}_i^\top \mathbf{w}_j| \leq \frac{\sqrt{\beta_3}}{6\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}$$

Combined with the bound on $\hat{\mathbf{v}}_i^\top \hat{\mathbf{v}}_j$ gives that

$$\mathbf{v}_i^\top \mathbf{v}_j \in \left[-\frac{\sqrt{\beta_3}}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}, \frac{\sqrt{\beta_3}}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}} \right]$$

Squaring both sides gives the desired result. Lastly, to bound $\mathbf{v}_i^\top \mathbf{w}_j$, we write

$$\mathbf{v}_i^\top \mathbf{w}_j = \mathbf{w}_j^\top \left(\mathbf{I} - \frac{1}{\|\mathbf{w}_i\|_2^2} \mathbf{w}_i \mathbf{w}_i^\top \right) \hat{\mathbf{v}}_i = \hat{\mathbf{v}}_i^\top \mathbf{w}_j - \frac{1}{\|\mathbf{w}_i\|_2^2} \mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_i^\top \mathbf{w}_j$$

Similar to the above, we have that

$$\frac{1}{\|\mathbf{w}_i\|_2^2} \mathbf{w}_i^\top \hat{\mathbf{v}}_i \cdot \mathbf{w}_i^\top \mathbf{w}_j \leq \frac{1}{d \left(1 - \frac{1}{2} \beta_s \delta_s\right)} \log \frac{8m^2}{\delta_{\mathbb{P}}} \leq \frac{\sqrt{\beta_3}}{6\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}}$$

Therefore, we can conclude that

$$\mathbf{v}_i^\top \mathbf{w}_j \in \left[-\frac{\sqrt{\beta_3}}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}}, \frac{\sqrt{\beta_3}}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right]$$

Squaring both sides gives the desired result. \square

Lemma 4. [Restatement of Lemma 2] Let $\mathbf{w}_1, \dots, \mathbf{w}_m \sim \mathcal{N}(0, d^{-1} \mathbf{I}_d)$ be I.I.D. Gaussian random vectors. Define

$$i_\ell^*, j_\ell^* = \arg \max_{i \in [m] \setminus \mathcal{R}_{\ell-1}, j \in [m^*] \setminus \mathcal{C}_{\ell-1}} \mathbf{w}_i[j]; \mathcal{R}_\ell = \{i_k^*\}_{k=1}^\ell; \mathcal{C}_\ell = \{j_k^*\}_{k=1}^\ell$$

Let any $\delta_{\mathbb{P}} \in (0, 1/2)$ be given. Then there exists some absolute constant $\beta_2, \beta_4 > 0$ such that if $m \geq \beta_4 m^* \log \frac{m^*}{\delta_{\mathbb{P}}}$, then for $\delta_s = \frac{\beta_2 \delta_{\mathbb{P}}}{m^2}$, with probability at least $1 - 4\delta_{\mathbb{P}}$, it holds that

- (Row-wise Gap) $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq (1 + 2\delta_s) \mathbf{w}_{i_\ell^*}[j]$ for all $\ell \in [m^*]$ and $j \in [m^*] \setminus \mathcal{C}_\ell$
- (Column-wise Gap) $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq (1 + 2\delta_s) \mathbf{w}_i[j_\ell^*]$ for all $\ell \in [m^*]$ and $i \in [m] \setminus \mathcal{R}_\ell$
- (Threshold Gap) $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq (1 + 2\delta_s) \mathbf{w}_{i_{\ell+1}^*}[j_{\ell+1}^*]^2$ for all $\ell \in [m^* - 1]$
- (Magnitude Lower Bound) $\mathbf{v}_{i_\ell^*}[j_\ell^*]^2 \geq \frac{\log m^*}{d}$ for all $\ell \in [m]$

Proof. We start by proving an auxiliary result that, with high probability, there are at least $\frac{m}{3}$ out of the m $\mathbf{w}_i[j]$'s that are positive for each $j \in [m]^*$. Define

$$s_{i,j} = \mathbb{I}\{\mathbf{w}_i[j] > 0\}; S_j = \sum_{i=1}^m s_{i,j}$$

Due to the symmetry of Gaussian, we have that $s_{i,j} \sim \text{Bern}(0.5)$ independently. Therefore, $\mathbb{E}[S_j] = \frac{m}{2}$. By Hoeffding's inequality, we have that

$$\Pr\left(S_j - \frac{m}{2} \leq -t\right) \leq \exp\left(-\frac{2t^2}{m}\right)$$

Setting $t = \frac{m}{6}$ and take a union bound over all $j \in [m^*]$ gives that

$$\Pr\left(S_j \geq \frac{m}{3}; \forall j \in [m^*]\right) \leq m^* \exp\left(-\frac{m}{18}\right)$$

Since $m \geq \beta_4 m^* \log \frac{m^*}{\delta_{\mathbb{P}}} \geq 18 \log \frac{m^*}{\delta_{\mathbb{P}}}$, with probability at least $1 - \delta_{\mathbb{P}}$, we have that at least $\frac{m}{3}$ out of $\{\mathbf{w}_i[j]\}_{i=1}^m$ are positive for all $j \in [m]$.

Proof of the first statement. Let any $i_1, i_2 \in [m]$ and $j_1, j_2 \in [m^*]$ such that (i_1, j_1) and (i_2, j_2) differ in at least one coordinate. Then we have that $\mathbf{w}_{i_1}[j_1]$ and $\mathbf{w}_{i_2}[j_2]$ are I.I.D. Gaussian random variables in $\mathcal{N}(0, d^{-1})$. Therefore, $\frac{\mathbf{w}_{i_1}[j_1]}{\mathbf{w}_{i_2}[j_2]}$ is a standard Cauchy random variable. Given the condition that $\mathbf{w}_{i_1}[j_1], \mathbf{w}_{i_2}[j_2] \geq 0$, we

have that $\mathbf{w}_{i_1}[j_1], \mathbf{w}_{i_2}[j_2]$ are half-Gaussian, and thus $\frac{\mathbf{w}_{i_1}[j_1]}{\mathbf{w}_{i_2}[j_2]}$ is half-Cauchy. Using the CDF of Cauchy random variables, we have that for any $\delta \in (0, 1/2)$

$$\begin{aligned}
 & \Pr(\mathbf{w}_{i_1}[j_1] \in (\mathbf{w}_{i_2}[j_2], (1 + \delta)\mathbf{w}_{i_2}[j_2]) \mid \mathbf{w}_{i_1}[j_1], \mathbf{w}_{i_2}[j_2] \geq 0) \\
 &= \Pr\left(\left|\frac{\mathbf{w}_{i_1}[j_1]}{\mathbf{w}_{i_2}[j_2]}\right| \in (1, 1 + \delta)\right) \\
 &= 2 \Pr\left(\frac{\mathbf{w}_{i_1}[j_1]}{\mathbf{w}_{i_2}[j_2]} \in (1, 1 + \delta)\right) \\
 &= \frac{2}{\pi} (\arctan(1 + \delta) - \arctan 1) \\
 &= \frac{2}{\pi} \arctan \frac{\delta}{2 + \delta} \\
 &\leq \frac{\delta}{\pi}
 \end{aligned}$$

where the last inequality follows from $\arctan x \leq x$ with $x \geq 0$ and $\delta > 0$. To prove the first property, we let $j_1 = j_\ell^*, j_2 = j$ and $i_1 = i_2 = i_\ell^*$. Fix any $\ell \in [m^*]$, we have that

$$\Pr(\mathbf{w}_{i_\ell^*}[j_\ell^*] \in (\mathbf{w}_{i_\ell^*}[j], (1 + \delta)\mathbf{w}_{i_\ell^*}[j]); \forall j \in [m^*] \setminus \mathcal{C}_\ell \text{ s.t. } \mathbf{w}_{i_\ell^*}[j] \geq 0) \leq m^* \delta$$

Recall that $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq \mathbf{w}_{i_\ell^*}[j]$ for all $j \in [m^*] \setminus \mathcal{C}_\ell$. If $\mathbf{w}_{i_\ell^*}[j] < 0$, since by definition $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq 0$, it must hold that $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq (1 + 2\delta_s)\mathbf{w}_{i_\ell^*}[j]$. Take a union bound over $\ell \in [m^*]$, and set $\delta = 2\delta_s$ with $\delta_s \leq \frac{\delta_p}{2m^*}$ gives the first property.

Proof of the second statement. To prove the second property, we set $i_1 = i_\ell^*, i_2 = i$ and $j_1 = j_2 = j_\ell^*$. Fix any $\ell \in [m^*]$, we have that

$$\Pr(\mathbf{w}_{i_\ell^*}[j_\ell^*] \in (\mathbf{w}_i[j_\ell^*], (1 + \delta)\mathbf{w}_i[j_\ell^*]); \forall i \in [m] \setminus \mathcal{R}_\ell \text{ s.t. } \mathbf{w}_i[j_\ell^*] \geq 0) \leq m\delta$$

Recall that $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq \mathbf{w}_i[j_\ell^*]$ for all $i \in [m] \setminus \mathcal{R}_\ell$. If $\mathbf{w}_i[j_\ell^*] < 0$, since by definition $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq 0$, it must hold that $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq (1 + 2\delta_s)\mathbf{w}_i[j_\ell^*]$. Take a union bound over $\ell \in [m^*]$, and set $\delta = 2\delta_s$ with $\delta_s \leq \frac{\delta_p}{2m^*}$ gives the second property.

Proof of the third statement. To prove the third property, we set $i_1 = i_\ell^*, i_2 = i_{\ell+1}^*$ and $j_1 = j_\ell^*, j_2 = j_{\ell+1}^*$. Fix any $\ell \in [m^* - 1]$, we have that

$$\Pr(\mathbf{w}_{i_\ell^*}[j_\ell^*] \in (\mathbf{w}_{i_{\ell+1}^*}[j_{\ell+1}^*], (1 + \delta)\mathbf{w}_{i_{\ell+1}^*}[j_{\ell+1}^*])) \leq \delta$$

Recall that $\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq \mathbf{w}_{i_{\ell+1}^*}[j_{\ell+1}^*] \geq 0$ for all $\ell \in [m^* - 1]$. Take a union bound over $\ell \in [m^* - 1]$, and set $\delta = 2\delta_s$ with $\delta_s \leq \frac{\delta_p}{2m^*}$ gives the third property.

Proof of the fourth statement. To show the last result, we notice that by definition of i_ℓ^* , it must holds that

$$\mathbf{w}_{i_\ell^*}[j_\ell^*] \geq \mathbf{w}_i[j_\ell^*]; \quad \forall i \in [m] \setminus \mathcal{R}_\ell$$

This gives that for any $\gamma > 0$

$$\begin{aligned}
 \Pr\left(\mathbf{w}_{i_\ell^*}[j_\ell^*] < \frac{\gamma}{\sqrt{d}}\right) &\leq \prod_{i \in [m] \setminus \mathcal{R}_\ell} \Pr\left(\mathbf{w}_i[j_\ell^*] < \frac{\gamma}{\sqrt{d}}\right) \\
 &= \prod_{i \in [m] \setminus \mathcal{R}_\ell} \Pr\left(\mathbf{w}_i[j_\ell^*] < \frac{\gamma}{\sqrt{d}}\right) \\
 &= \Pr_{z \sim \mathcal{N}(0,1)}(z < \gamma)^{m-m^*} \\
 &\leq \left(1 - \Pr_{z \sim \mathcal{N}(0,1)}(z \geq \gamma)\right)^{m-m^*}
 \end{aligned}$$

By the tail bound of Gaussian random variable, we have that

$$\Pr_{z \sim \mathcal{N}(0,1)} (z \geq \gamma) \geq \frac{1}{\sqrt{2\pi} \cdot \gamma} \exp\left(-\frac{\gamma^2}{2}\right)$$

Therefore

$$\Pr\left(\mathbf{w}_{i_\ell^*}[j_\ell^*] < \frac{\gamma}{\sqrt{d}}\right) \leq \left(1 - \frac{1}{\sqrt{2\pi} \cdot \gamma} \exp\left(-\frac{\gamma^2}{2}\right)\right)^{m-m^*}$$

Take a union bound over all $\ell \in [m^*]$ gives that

$$\Pr\left(\mathbf{w}_{i_\ell^*}[j_\ell^*]^2 \geq \frac{\gamma^2}{d}\right) \geq 1 - m^* \left(1 - \frac{1}{\sqrt{2\pi} \cdot \gamma} \exp\left(-\frac{\gamma^2}{2}\right)\right)^{m-m^*}$$

For the failing probability to be upper bounded by $\delta_{\mathbb{P}}$, we simply need

$$m \geq m^* + \frac{\log \frac{m^*}{\delta_{\mathbb{P}}}}{\log\left(1 - \frac{1}{\sqrt{2\pi} \cdot \gamma} \exp\left(-\frac{\gamma^2}{2}\right)\right)^{-1}}$$

Using $\log \frac{1}{x} \geq 1 - x$ for all $x > 0$, it suffice to guarantee that

$$m \geq m^* + \sqrt{2\pi} \cdot \gamma \exp\left(\frac{\gamma^2}{2}\right) \log \frac{m^*}{\delta_{\mathbb{P}}}$$

Setting $\gamma = (\log m^*)^{\frac{1}{2}}$ gives the desired result. \square

A.3 Hermite Expansion of the Gradient and the Gradient Flow Dynamics

Notice that the population MSE has the following form

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[f(\boldsymbol{\theta}, \mathbf{x})^2 - 2f(\boldsymbol{\theta}, \mathbf{x})f(\boldsymbol{\theta}^*, \mathbf{x}) + f(\boldsymbol{\theta}^*, \mathbf{x})^2 \right]$$

where the last term is independent of $\boldsymbol{\theta}^*$, and is thus omitted from our analysis. Consider the Hermite expansion of $\pi(\cdot)$ and $\sigma(\cdot)$, respectively

$$\pi(x) = \sum_{k=0}^{\infty} c_k H e_k(x); \quad \sigma(x) = H e_3(x)$$

where $c_k = \mathbb{E}_{x \sim \mathcal{N}(0,1)} [\pi^{(k)}(x)]$. Then we have that

$$\begin{aligned} f(\boldsymbol{\theta}, \mathbf{x})^2 &= \sum_{i,j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j) \\ f(\boldsymbol{\theta}, \mathbf{x}) f(\boldsymbol{\theta}^*, \mathbf{x}) &= \sum_{i=1}^m \sum_{j=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j^*) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j^*) \end{aligned}$$

This gives that

$$\frac{\partial}{\partial \mathbf{v}_i} f(\boldsymbol{\theta}, \mathbf{x})^2 = 2 \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \cdot H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j) (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \frac{\mathbf{x}}{\|\mathbf{v}_i\|_2}$$

Taking the expectation gives

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x}} \left[\frac{\partial}{\partial \mathbf{v}_i} f(\boldsymbol{\theta}, \mathbf{x})^2 \right] \\
 &= \frac{2(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j) \mathbf{x}] \\
 &= \frac{2(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_{\ell+1}}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j)] \bar{\mathbf{v}}_j \\
 &\quad + \frac{6(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j)] \bar{\mathbf{w}}_i \\
 &\quad + \frac{6(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_j)] \bar{\mathbf{w}}_j
 \end{aligned}$$

Similar, we can obtain that

$$\begin{aligned}
 & 2\mathbb{E}_{\mathbf{x}} \left[\frac{\partial}{\partial \mathbf{v}_i} f(\boldsymbol{\theta}, \mathbf{x}) f(\boldsymbol{\theta}^*, \mathbf{x}) \right] \\
 &= \frac{2(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \sum_{j=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_{\ell+1}}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j^*) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j^*)] \bar{\mathbf{v}}_j^* \\
 &\quad + \frac{6(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \sum_{j=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j^*) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j^*)] \bar{\mathbf{w}}_i \\
 &\quad + \frac{6(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \sum_{j=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j^*) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_j^*)] \bar{\mathbf{w}}_j^*
 \end{aligned}$$

For the gradient of \mathbf{w}_i , we can compute that

$$\frac{\partial}{\partial \mathbf{w}_i} f(\boldsymbol{\theta}, \mathbf{x})^2 = 6 \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j) (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \frac{\mathbf{x}}{\|\bar{\mathbf{w}}_i\|_2}$$

Taking the expectation gives

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x}} \left[\frac{\partial}{\partial \mathbf{w}_i} f(\boldsymbol{\theta}, \mathbf{x})^2 \right] \\
 &= \frac{6(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{w}_i\|_2} \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j)] \bar{\mathbf{v}}_i \\
 &\quad + \frac{6(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{w}_i\|_2} \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell+1}}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j)] \bar{\mathbf{v}}_j \\
 &\quad + \frac{18(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{w}_i\|_2} \sum_{j=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_j)] \bar{\mathbf{w}}_j
 \end{aligned}$$

Similarly, we have

$$\begin{aligned}
 & 2\mathbb{E}_{\mathbf{x}} \left[\frac{\partial}{\partial \mathbf{w}_i} f(\boldsymbol{\theta}, \mathbf{x}) f(\boldsymbol{\theta}^*, \mathbf{x}) \right] \\
 &= \frac{6(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{w}_i\|_2} \sum_{j=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j^*) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j^*)] \bar{\mathbf{v}}_i \\
 &+ \frac{6(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{w}_i\|_2} \sum_{j=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell+1}}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j^*) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_3(\mathbf{x}^\top \bar{\mathbf{w}}_j^*)] \bar{\mathbf{v}}_j^* \\
 &+ \frac{18(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{w}_i\|_2} \sum_{j=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j^*) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_2(\mathbf{x}^\top \bar{\mathbf{w}}_j^*)] \bar{\mathbf{w}}_j^*
 \end{aligned}$$

This gives that

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x}} \left[\frac{\partial}{\partial \mathbf{v}_i} f(\boldsymbol{\theta}, \mathbf{x}) f(\boldsymbol{\theta}^*, \mathbf{x}) \right]^\top \mathbf{v}_i = \mathbb{E}_{\mathbf{x}} \left[\frac{\partial}{\partial \mathbf{v}_i} f(\boldsymbol{\theta}, \mathbf{x})^2 \right]^\top \mathbf{v}_i = 0 \\
 & \mathbb{E}_{\mathbf{x}} \left[\frac{\partial}{\partial \mathbf{w}_i} f(\boldsymbol{\theta}, \mathbf{x}) f(\boldsymbol{\theta}^*, \mathbf{x}) \right]^\top \mathbf{w}_i = \mathbb{E}_{\mathbf{x}} \left[\frac{\partial}{\partial \mathbf{w}_i} f(\boldsymbol{\theta}, \mathbf{x})^2 \right]^\top \mathbf{w}_i = 0
 \end{aligned}$$

According to the gradient flow dynamics, we can conclude that

$$\frac{d}{dt} \|\mathbf{v}_i(t)\|_2^2 = \frac{d}{dt} \|\mathbf{w}_i(t)\|_2^2 = 0$$

which implies that the norm of each \mathbf{v}_i and \mathbf{w}_i are fixed at initialization. For the convenience of the analysis, we shall denote

$$\begin{aligned}
 \mathcal{C}_{k,\ell,a,b}^{i,j} &= \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j) H e_a(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_b(\mathbf{x}^\top \bar{\mathbf{w}}_j)] \\
 \hat{\mathcal{C}}_{k,\ell,a,b}^{i,j} &= \mathbb{E}_{\mathbf{x}} [H e_k(\mathbf{x}^\top \bar{\mathbf{v}}_i) H e_\ell(\mathbf{x}^\top \bar{\mathbf{v}}_j^*) H e_a(\mathbf{x}^\top \bar{\mathbf{w}}_i) H e_b(\mathbf{x}^\top \bar{\mathbf{w}}_j^*)]
 \end{aligned}$$

and that $\|\bar{\mathbf{v}}_i(t)\|_2 = a_i$, $\|\bar{\mathbf{w}}_i(t)\|_2 = b_i$. Then we have that

$$\begin{aligned}
 \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}) &= \frac{\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top}{\|\mathbf{v}_i\|_2} \left(\sum_{r=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_{\ell+1}}{k! \ell!} \mathcal{C}_{k,\ell,3,3}^{i,r} \bar{\mathbf{v}}_r - \sum_{r=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_{\ell+1}}{k! \ell!} \hat{\mathcal{C}}_{k,\ell,3,3}^{i,r} \bar{\mathbf{v}}_r^* \right) \\
 &+ \frac{3(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \left(\sum_{r=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathcal{C}_{k,\ell,2,3}^{i,r} - \sum_{r=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \hat{\mathcal{C}}_{k,\ell,2,3}^{i,r} \right) \bar{\mathbf{w}}_i \\
 &+ \frac{3(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)}{\|\mathbf{v}_i\|_2} \left(\sum_{r=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathcal{C}_{k,\ell,3,2}^{i,r} \bar{\mathbf{w}}_r - \sum_{r=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \hat{\mathcal{C}}_{k,\ell,3,2}^{i,r} \bar{\mathbf{w}}_r^* \right) \\
 \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}) &= \frac{3(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{w}_i\|_2} \left(\sum_{r=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathcal{C}_{k,\ell,2,3}^{i,r} - \sum_{r=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \hat{\mathcal{C}}_{k,\ell,2,3}^{i,r} \right) \bar{\mathbf{v}}_i \\
 &+ \frac{3(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{v}_i\|_2} \left(\sum_{r=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell+1}}{k! \ell!} \mathcal{C}_{k,\ell,2,3}^{i,r} \bar{\mathbf{v}}_r - \sum_{r=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell+1}}{k! \ell!} \hat{\mathcal{C}}_{k,\ell,2,3}^{i,r} \bar{\mathbf{v}}_r^* \right) \\
 &+ \frac{9(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)}{\|\mathbf{w}_i\|_2} \left(\sum_{r=1}^m \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} \mathcal{C}_{k,\ell,2,2}^{i,r} \bar{\mathbf{w}}_r - \sum_{r=1}^{m^*} \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} \hat{\mathcal{C}}_{k,\ell,2,2}^{i,r} \bar{\mathbf{w}}_r^* \right)
 \end{aligned} \tag{10}$$

In particular, we notice that there are several quantities that appears in the form of the gradient. We make the following definition for the convenience of the analysis

$$\begin{aligned}
 \lambda_{i,j,1} &= \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell+1}}{k!\ell!} \mathcal{C}_{k,\ell,3,3}^{i,j}; & \hat{\lambda}_{i,j,1} &= \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell+1}}{k!\ell!} \hat{\mathcal{C}}_{k,\ell,3,3}^{i,j}; \\
 \lambda_{i,j,2} &= \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell}}{k!\ell!} \mathcal{C}_{k,\ell,2,3}^{i,j}; & \hat{\lambda}_{i,j,2} &= \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell}}{k!\ell!} \hat{\mathcal{C}}_{k,\ell,2,3}^{i,j}; \\
 \lambda_{i,j,3} &= \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell+1}}{k!\ell!} \mathcal{C}_{k,\ell,2,3}^{i,j}; & \hat{\lambda}_{i,j,3} &= \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell+1}}{k!\ell!} \hat{\mathcal{C}}_{k,\ell,2,3}^{i,j}; \\
 \lambda_{i,j,4} &= \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell}}{k!\ell!} \mathcal{C}_{k,\ell,3,2}^{i,j}; & \hat{\lambda}_{i,j,4} &= \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell}}{k!\ell!} \hat{\mathcal{C}}_{k,\ell,3,2}^{i,j}; \\
 \lambda_{i,j,5} &= \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell}}{k!\ell!} \mathcal{C}_{k,\ell,2,2}^{i,j}; & \hat{\lambda}_{i,j,5} &= \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell}}{k!\ell!} \hat{\mathcal{C}}_{k,\ell,2,2}^{i,j};
 \end{aligned}$$

Recall that our goal is to study the dynamics of the following alignment scores

$$\begin{aligned}
 \gamma_{i,j}^{(1)}(t) &= \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j^*; & \gamma_{i,j}^{(2)}(t) &= \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j^*; & \zeta_{i,j}^{(1)}(t) &= \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{w}}_j^*; & \zeta_{i,j}^{(2)}(t) &= \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{v}}_j^* \\
 I_{i,j}^{(1)}(t) &= \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j(t); & I_{i,j}^{(2)}(t) &= \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j(t); & I_{i,j}^{(3)}(t) &= \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{w}}_j(t)
 \end{aligned} \tag{11}$$

This allows us to rewrite the gradient as

$$\begin{aligned}
 \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{a_i} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \left(\sum_{r=1}^m (\lambda_{i,r,1} \bar{\mathbf{v}}_r + 3\lambda_{i,r,2} \bar{\mathbf{w}}_i + 3\lambda_{i,r,4} \bar{\mathbf{w}}_r) - \sum_{r=1}^{m^*} (\hat{\lambda}_{i,r,1} \bar{\mathbf{v}}_r^* + 3\hat{\lambda}_{i,r,2} \bar{\mathbf{w}}_i + 3\hat{\lambda}_{i,r,4} \bar{\mathbf{w}}_r^*) \right) \\
 &= \frac{1}{a_i} \sum_{r=1}^m (\lambda_{i,r,1} \bar{\mathbf{v}}_r + 3\lambda_{i,r,2} \bar{\mathbf{w}}_i + 3\lambda_{i,r,4} \bar{\mathbf{w}}_r) - \frac{1}{a_i} \sum_{r=1}^{m^*} (\hat{\lambda}_{i,r,1} \bar{\mathbf{v}}_r^* + 3\hat{\lambda}_{i,r,2} \bar{\mathbf{w}}_i + 3\hat{\lambda}_{i,r,4} \bar{\mathbf{w}}_r^*) \\
 &\quad - \frac{1}{a_i} \sum_{r=1}^m (\lambda_{i,r,1} I_{i,r}^{(1)} + 3\lambda_{i,r,2} I_{i,i}^{(3)} + 3\lambda_{i,r,4} I_{i,r}^{(3)}) \bar{\mathbf{v}}_i \\
 &\quad + \frac{1}{a_i} \sum_{r=1}^{m^*} (\hat{\lambda}_{i,r,1} \gamma_{i,r}^{(1)} + 3\hat{\lambda}_{i,r,2} I_{i,i}^{(3)} + 3\hat{\lambda}_{i,r,4} \zeta_{i,r}^{(1)}) \bar{\mathbf{v}}_i \\
 \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}) &= \frac{3}{b_i} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \left(\sum_{r=1}^m (\lambda_{i,r,2} \bar{\mathbf{v}}_i + \lambda_{i,r,3} \bar{\mathbf{v}}_r + 3\lambda_{i,r,5} \bar{\mathbf{w}}_r) - \sum_{r=1}^{m^*} (\hat{\lambda}_{i,r,2} \bar{\mathbf{v}}_i + \hat{\lambda}_{i,r,3} \bar{\mathbf{v}}_r^* + 3\hat{\lambda}_{i,r,5} \bar{\mathbf{w}}_r^*) \right) \\
 &= \frac{3}{b_i} \sum_{r=1}^m (\lambda_{i,r,2} \bar{\mathbf{v}}_i + \lambda_{i,r,3} \bar{\mathbf{v}}_r + 3\lambda_{i,r,5} \bar{\mathbf{w}}_r) - \frac{3}{b_i} \sum_{r=1}^{m^*} (\hat{\lambda}_{i,r,2} \bar{\mathbf{v}}_i + \hat{\lambda}_{i,r,3} \bar{\mathbf{v}}_r^* + 3\hat{\lambda}_{i,r,5} \bar{\mathbf{w}}_r^*) \\
 &\quad - \frac{3}{b_i} \sum_{r=1}^m (\lambda_{i,r,2} I_{i,i}^{(3)} + \lambda_{i,r,3} I_{r,i}^{(3)} + 3\lambda_{i,r,5} I_{r,i}^{(2)}) \bar{\mathbf{w}}_i \\
 &\quad + \frac{3}{b_i} \sum_{r=1}^{m^*} (\hat{\lambda}_{i,r,2} I_{i,i}^{(3)} + \hat{\lambda}_{i,r,3} \zeta_{i,r}^{(2)} + 3\hat{\lambda}_{i,r,5} \gamma_{i,r}^{(2)}) \bar{\mathbf{w}}_i
 \end{aligned}$$

Recall that our goal is to study the dynamics in (9). By the gradient flow dynamic, we have that

$$\begin{aligned}
 \frac{d}{dt} \gamma_{i,j}^{(1)}(t) &= -\frac{1}{a_i} \cdot \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j^*; & \frac{d}{dt} \gamma_{i,j}^{(2)}(t) &= -\frac{1}{b_i} \cdot \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j^* \\
 \frac{d}{dt} \zeta_{i,j}^{(1)}(t) &= -\frac{1}{a_i} \cdot \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j^*; & \frac{d}{dt} \zeta_{i,j}^{(2)}(t) &= -\frac{1}{b_i} \cdot \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j^*
 \end{aligned}$$

Moreover, we also have that

$$\begin{aligned}\frac{d}{dt}I_{i,j}^{(1)}(t) &= -\frac{1}{a_i}\nabla_{\mathbf{v}_i}\mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j - \frac{1}{a_j}\nabla_{\mathbf{v}_j}\mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_i \\ \frac{d}{dt}I_{i,j}^{(2)}(t) &= -\frac{1}{b_i}\nabla_{\mathbf{w}_i}\mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j - \frac{1}{b_j}\nabla_{\mathbf{w}_j}\mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_i \\ \frac{d}{dt}I_{i,j}^{(3)}(t) &= -\frac{1}{a_i}\nabla_{\mathbf{v}_i}\mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j - \frac{1}{b_j}\nabla_{\mathbf{w}_j}\mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_i\end{aligned}$$

Therefore, we should consider inner product between the gradient and the vectors $\bar{\mathbf{v}}_i, \bar{\mathbf{w}}_i, \bar{\mathbf{v}}_j^*, \bar{\mathbf{w}}_j^*$ above, which will give us exactly eight terms to analyze. The inner product between $\nabla_{\mathbf{v}_i}\mathcal{L}(\boldsymbol{\theta}(t))$ and $\bar{\mathbf{v}}_j^*, \bar{\mathbf{w}}_j^*$ can be written as

$$\begin{aligned}-\frac{1}{a_i}\nabla_{\mathbf{v}_i}\mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j^* &= \frac{1}{a_i^2}\left(\hat{\lambda}_{i,j,1}(t) - \sum_{r=1}^m \lambda_{i,r,1}(t)\gamma_{r,j}^{(1)}(t)\right) \\ &\quad - \frac{3}{a_i^2}\left(\sum_{r=1}^m \left(\lambda_{i,r,2}(t)\zeta_{i,j}^{(2)}(t) + \lambda_{i,r,4}\zeta_{r,j}^{(2)}(t)\right) - \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t)\zeta_{i,j}^{(2)}(t)\right) \\ &\quad + \frac{1}{a_i^2}\sum_{r=1}^m \left(\lambda_{i,r,1}(t)I_{i,r}^{(1)}(t) + 3\lambda_{i,r,2}(t)I_{i,i}^{(3)}(t) + 3\lambda_{i,r,4}(t)I_{i,r}^{(3)}(t)\right)\gamma_{i,j}^{(1)}(t) \\ &\quad - \frac{1}{a_i^2}\sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,1}(t)\gamma_{i,r}^{(1)}(t) + 3\hat{\lambda}_{i,r,2}(t)I_{i,i}^{(3)}(t) + 3\hat{\lambda}_{i,r,4}(t)\zeta_{i,r}^{(1)}(t)\right)\gamma_{i,j}^{(1)}(t) \\ -\frac{1}{a_i}\nabla_{\mathbf{v}_i}\mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j^* &= \frac{3}{a_i^2}\left(\hat{\lambda}_{i,j,4}(t) - \sum_{r=1}^m \lambda_{i,r,4}(t)\gamma_{r,j}^{(2)}(t)\right) \\ &\quad - \frac{1}{a_i^2}\left(\sum_{r=1}^m \left(\lambda_{i,r,1}(t)\zeta_{r,j}^{(1)}(t) + 3\lambda_{i,r,2}(t)\gamma_{i,j}^{(2)}(t)\right) - 3\sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t)\gamma_{i,j}^{(2)}(t)\right) \\ &\quad + \frac{1}{a_i^2}\sum_{r=1}^m \left(\lambda_{i,r,1}(t)I_{i,r}^{(1)}(t) + 3\lambda_{i,r,2}(t)I_{i,i}^{(3)}(t) + 3\lambda_{i,r,4}(t)I_{i,r}^{(3)}(t)\right)\zeta_{i,j}^{(1)}(t) \\ &\quad - \frac{1}{a_i^2}\sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,1}(t)\gamma_{i,r}^{(1)}(t) + 3\hat{\lambda}_{i,r,2}(t)I_{i,i}^{(3)}(t) + 3\hat{\lambda}_{i,r,4}\zeta_{i,r}^{(1)}(t)\right)\zeta_{i,j}^{(1)}(t)\end{aligned}$$

The inner product between $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))$ and $\bar{\mathbf{v}}_j^*, \bar{\mathbf{w}}_j^*$ can be written as

$$\begin{aligned}
 -\frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j^* &= \frac{9}{b_i^2} \left(\hat{\lambda}_{i,j,5}(t) - \sum_{r=1}^m \lambda_{i,r,5}(t) \gamma_{r,j}^{(2)}(t) \right) \\
 &\quad - \frac{3}{b_i^2} \left(\sum_{r=1}^m \left(\lambda_{i,r,2}(t) \zeta_{i,j}^{(1)}(t) + \lambda_{i,r,3}(t) \zeta_{r,j}^{(1)}(t) \right) - \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) \zeta_{i,j}^{(1)}(t) \right) \\
 &\quad + \frac{3}{b_i^2} \sum_{r=1}^m \left(\lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) + \lambda_{i,r,3}(t) I_{r,i}^{(3)}(t) + 3\lambda_{i,r,5}(t) I_{r,i}^{(2)}(t) \right) \gamma_{i,j}^{(2)}(t) \\
 &\quad - \frac{3}{b_i^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) + \hat{\lambda}_{i,r,3}(t) \zeta_{i,r}^{(2)}(t) + 3\hat{\lambda}_{i,r,5}(t) \gamma_{i,r}^{(2)}(t) \right) \gamma_{i,j}^{(2)}(t) \\
 -\frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j^* &= \frac{3}{b_i^2} \left(\hat{\lambda}_{i,j,3}(t) - \sum_{r=1}^m \lambda_{i,r,3}(t) \gamma_{r,j}^{(1)}(t) \right) \\
 &\quad - \frac{3}{b_i^2} \left(\sum_{r=1}^m \left(\lambda_{i,r,2}(t) \gamma_{i,j}^{(1)}(t) + 3\lambda_{i,r,5}(t) \zeta_{r,j}^{(2)}(t) \right) - \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) \gamma_{i,j}^{(1)}(t) \right) \\
 &\quad + \frac{3}{b_i^2} \sum_{r=1}^m \left(\lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) + \lambda_{i,r,3}(t) I_{r,i}^{(3)}(t) + 3\lambda_{i,r,5}(t) I_{r,i}^{(2)}(t) \right) \zeta_{i,j}^{(2)}(t) \\
 &\quad - \frac{3}{b_i^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) + \hat{\lambda}_{i,r,3}(t) \zeta_{i,r}^{(2)}(t) + 3\hat{\lambda}_{i,r,5}(t) \gamma_{i,r}^{(2)}(t) \right) \zeta_{i,j}^{(2)}(t)
 \end{aligned}$$

The inner product between $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))$ and $\bar{\mathbf{v}}_j, \bar{\mathbf{w}}_j$ can be written as

$$\begin{aligned}
 -\frac{1}{a_i} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j &= \frac{1}{a_i^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,1}(t) \gamma_{j,r}^{(1)}(t) + 3\hat{\lambda}_{i,r,2}(t) I_{j,i}^{(3)}(t) + 3\hat{\lambda}_{i,r,4}(t) \zeta_{j,r}^{(1)}(t) \right) \\
 &\quad - \frac{1}{a_i^2} \sum_{r=1}^m \left(\lambda_{i,r,1}(t) I_{r,j}^{(1)}(t) + 3\lambda_{i,r,2}(t) I_{j,i}^{(3)}(t) + 3\lambda_{i,r,4}(t) I_{j,r}^{(3)}(t) \right) \\
 &\quad + \frac{1}{a_i^2} \sum_{r=1}^m \left(\lambda_{i,r,1}(t) I_{i,r}^{(1)}(t) + 3\lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) + 3\lambda_{i,r,4}(t) I_{i,r}^{(3)}(t) \right) I_{i,j}^{(1)}(t) \\
 &\quad - \frac{1}{a_i^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,1}(t) \gamma_{i,r}^{(1)}(t) + 3\hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) + 3\hat{\lambda}_{i,r,4}(t) \zeta_{i,r}^{(1)}(t) \right) I_{i,j}^{(1)}(t) \\
 -\frac{1}{a_i} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j &= \frac{1}{a_i^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,1}(t) \zeta_{j,r}^{(2)} + 3\hat{\lambda}_{i,r,2}(t) I_{i,j}^{(2)} + 3\hat{\lambda}_{i,r,4}(t) \gamma_{j,r}^{(2)}(t) \right) \\
 &\quad - \frac{1}{a_i^2} \sum_{r=1}^m \left(\lambda_{i,r,1}(t) I_{r,j}^{(3)}(t) + 3\lambda_{i,r,2}(t) I_{i,j}^{(2)}(t) + 3\lambda_{i,r,4}(t) I_{r,j}^{(2)}(t) \right) \\
 &\quad + \frac{1}{a_i^2} \sum_{r=1}^m \left(\lambda_{i,r,1}(t) I_{i,r}^{(1)}(t) + 3\lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) + 3\lambda_{i,r,4}(t) I_{i,r}^{(3)}(t) \right) I_{i,j}^{(3)}(t) \\
 &\quad - \frac{1}{a_i^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,1}(t) \gamma_{i,r}^{(1)}(t) + 3\hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) + 3\hat{\lambda}_{i,r,4}(t) \zeta_{i,r}^{(1)}(t) \right) I_{i,j}^{(3)}(t)
 \end{aligned}$$

Lastly, the inner product between $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))$ and $\bar{\mathbf{v}}_j, \bar{\mathbf{w}}_j$ can be written as

$$\begin{aligned}
 -\frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j &= \frac{3}{b_j^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,2}(t) I_{i,j}^{(1)}(t) + \hat{\lambda}_{i,r,3}(t) \gamma_{j,r}^{(1)}(t) + 3\hat{\lambda}_{i,r,5}(t) \zeta_{j,r}^{(1)}(t) \right) \\
 &\quad - \frac{3}{b_j^2} \sum_{r=1}^m \left(\lambda_{i,r,2}(t) I_{i,j}^{(1)}(t) + \lambda_{i,r,3}(t) I_{j,r}^{(1)}(t) + 3\lambda_{i,r,5}(t) I_{j,r}^{(3)}(t) \right) \\
 &\quad + \frac{3}{b_j^2} \sum_{r=1}^m \left(\lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) + \lambda_{i,r,3}(t) I_{r,i}^{(3)}(t) + 3\lambda_{i,r,5}(t) I_{r,i}^{(2)}(t) \right) I_{i,j}^{(3)}(t) \\
 &\quad - \frac{3}{b_j^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) + \hat{\lambda}_{i,r,3}(t) \zeta_{i,r}^{(2)}(t) + 3\hat{\lambda}_{i,r,5}(t) \gamma_{i,r}^{(2)}(t) \right) I_{i,j}^{(3)}(t) \\
 -\frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j &= \frac{3}{b_i^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,2}(t) I_{i,j}^{(3)}(t) + \hat{\lambda}_{i,r,3}(t) \zeta_{j,r}^{(2)}(t) + 3\hat{\lambda}_{i,r,5}(t) \gamma_{j,r}^{(2)}(t) \right) \\
 &\quad - \frac{3}{b_i^2} \sum_{r=1}^m \left(\lambda_{i,r,2}(t) I_{i,j}^{(3)}(t) + \lambda_{i,r,3}(t) I_{r,j}^{(3)}(t) + 3\lambda_{i,r,5}(t) I_{r,j}^{(2)}(t) \right) \\
 &\quad + \frac{3}{b_i^2} \sum_{r=1}^m \left(\lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) + \lambda_{i,r,3}(t) I_{r,i}^{(3)}(t) + 3\lambda_{i,r,5}(t) I_{r,i}^{(2)}(t) \right) I_{i,j}^{(2)}(t) \\
 &\quad - \frac{3}{b_i^2} \sum_{r=1}^{m^*} \left(\hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) + \hat{\lambda}_{i,r,3}(t) \zeta_{i,r}^{(2)}(t) + 3\hat{\lambda}_{i,r,5}(t) \gamma_{i,r}^{(2)}(t) \right) I_{i,j}^{(2)}(t)
 \end{aligned}$$

A.4 Approximating the Gradient Flow Dynamics

In order to understand the property of the GF induced dynamics given in the previous section, we need to first compute $\lambda_{i,j,\ell s}$ and $\hat{\lambda}_{i,j,\ell s}$. The following lemma provides such property.

Lemma 5. Fix $i \in [m], j \in [m^*]$ and $t \geq 0$. If for any $\left| \zeta_{i,j}^{(1)}(t) \right|, \left| \zeta_{i,j}^{(2)}(t) \right|, \left| I_{i,i}^{(3)}(t) \right| \leq \delta_r$ for some $\delta_r > 0$, then we have that

$$\begin{aligned}
 \hat{\lambda}_{i,j,1}(t) &= 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^3 \pm \mathcal{O}(\delta_r^2) \gamma_{i,j}^{(2)}(t)^2 \pm \mathcal{O}(\delta_r^4) \\
 \hat{\lambda}_{i,j,2}(t) &= 6 \sum_{k=1}^{\infty} \frac{c_{k+2} c_k}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \zeta_{i,j}^{(1)}(t) \pm \mathcal{O}(\delta_r^3) \\
 \hat{\lambda}_{i,j,3}(t) &= 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \zeta_{i,j}^{(1)}(t) \pm \mathcal{O}(\delta_r^3) \\
 \hat{\lambda}_{i,j,4}(t) &= 6 \sum_{k=1}^{\infty} \frac{c_{k+2} c_k}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 I_{i,i}^{(3)}(t) + 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \zeta_{i,j}^{(2)}(t) \pm \mathcal{O}(\delta_r^3) \\
 \hat{\lambda}_{i,j,5}(t) &= 2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \pm \mathcal{O}(\delta_r^2) \gamma_{i,j}^{(2)}(t) \pm \mathcal{O}(\delta_r^4)
 \end{aligned}$$

Proof. We are going to Lemma 21 with $\mathbf{v}_1 = \bar{\mathbf{v}}_i(t), \mathbf{v}_2 = \bar{\mathbf{v}}_j^*$ and $\mathbf{w}_1 = \bar{\mathbf{w}}_i(t), \mathbf{w}_2 = \bar{\mathbf{w}}_j^*$. In this case, we have that $\mathbf{v}_1^\top \mathbf{v}_2 = \gamma_{i,j}^{(1)}(t)$ and $\mathbf{w}_1^\top \mathbf{w}_2 = \gamma_{i,j}^{(2)}(t)$. We start with $\hat{\lambda}_{i,j,1}(t)$. By definition, we have that

$$\hat{\lambda}_{i,j,1}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_{\ell+1}}{k! \ell!} \mathbb{E} \left[H e_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) H e_\ell(\bar{\mathbf{v}}_j^{*\top} \mathbf{x}) H e_3(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) H e_3(\bar{\mathbf{w}}_j^{*\top} \mathbf{x}) \right]$$

Here, invoking Lemma 21 with $h_k = c_{k+1}, h'_\ell = c_{\ell+1}$ gives

$$\hat{\lambda}_{i,j,1}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^3 \pm \mathcal{O}(\delta_r^2) \gamma_{i,j}^{(2)}(t)^2 \pm \mathcal{O}(\delta_r^4)$$

For $\hat{\lambda}_{i,j,5}(t)$, by definition we have

$$\hat{\lambda}_{i,j,5}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} \mathbb{E} [He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_\ell(\bar{\mathbf{v}}_j^*{}^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_j^*{}^\top \mathbf{x})]$$

Invoking Lemma 21 with $h_k = c_k, h'_\ell = c_\ell$ gives

$$\hat{\lambda}_{i,j,5}(t) = 2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \pm \mathcal{O}(\delta_r^2) \gamma_{i,j}^{(2)}(t) \pm \mathcal{O}(\delta_r^4)$$

For $\hat{\lambda}_{i,j,2}(t)$, by definition we have

$$\hat{\lambda}_{i,j,2}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E} [He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_\ell(\bar{\mathbf{v}}_j^*{}^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_j^*{}^\top \mathbf{x})]$$

Here, invoking Lemma 21 with $h_k = c_{k+1}, h'_\ell = c_\ell$, and noticing that $\bar{\mathbf{v}}_j^*{}^\top \bar{\mathbf{w}}_j = 0$, gives

$$\hat{\lambda}_{i,j,2}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+2} c_k}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \zeta_{i,j}^{(1)}(t) \pm \mathcal{O}(\delta_r^3)$$

Noticing that $c_2 = 0$ gives that

$$\hat{\lambda}_{i,j,2}(t) = 6 \sum_{k=1}^{\infty} \frac{c_{k+2} c_k}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \zeta_{i,j}^{(1)}(t) \pm \mathcal{O}(\delta_r^3)$$

For $\hat{\lambda}_{i,j,3}(t)$, by definition, we have that

$$\hat{\lambda}_{i,j,3}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell+1}}{k! \ell!} \mathbb{E} [He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_\ell(\bar{\mathbf{v}}_j^*{}^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_j^*{}^\top \mathbf{x})]$$

Invoking Lemma 21 with $h_k = c_k, h'_\ell = c_{\ell+1}$, and noticing that $\bar{\mathbf{v}}_j^*{}^\top \bar{\mathbf{w}}_j = 0$, gives

$$\hat{\lambda}_{i,j,3}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \zeta_{i,j}^{(1)}(t) \pm \mathcal{O}(\delta_r^3)$$

Lastly, for $\hat{\lambda}_4$, by definition we have

$$\hat{\lambda}_{i,j,4}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_{k+1} c_\ell}{k! \ell!} \mathbb{E} [He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_\ell(\bar{\mathbf{v}}_j^*{}^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_j^*{}^\top \mathbf{x})]$$

Therefore, we need to consider $\mathbf{v}_1 = \bar{\mathbf{v}}_i(t), \mathbf{v}_2 = \bar{\mathbf{v}}_j^*$ and $\mathbf{w}_1 = \bar{\mathbf{w}}_j^*, \mathbf{w}_2 = \bar{\mathbf{w}}_i(t)$. Moreover, we need to set $h_k = c_{k+1}$ and $h'_\ell = c_\ell$. In this case, $\mathbf{v}_1^\top \mathbf{w}_2 = I_{i,i}^{(3)}(t)$ and $\mathbf{v}_2^\top \mathbf{w}_2 = \zeta_{i,j}^{(2)}(t)$. Therefore

$$\hat{\lambda}_{i,j,4}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+2} c_k}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 I_{i,i}^{(3)}(t) + 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \zeta_{i,j}^{(2)}(t) \pm \mathcal{O}(\delta_r^3)$$

Noticing that $c_2 = 0$ gives that

$$\hat{\lambda}_{i,j,4}(t) = 6 \sum_{k=1}^{\infty} \frac{c_{k+2} c_k}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 I_{i,i}^{(3)}(t) + 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 \zeta_{i,j}^{(2)}(t) \pm \mathcal{O}(\delta_r^3)$$

□

Lemma 6. Fix $i, j \in [m]$ and $t \geq 0$. If for any $|I_{i,j}^{(1)}(t)|, |I_{i,j}^{(2)}(t)|, |I_{i,j}^{(3)}(t)| \leq \delta_r$ when $i \neq j$, and $|I_{i,i}^{(3)}(t)| \leq \delta_p \leq \mathcal{O}(\delta_r)$, then we have that

$$\begin{aligned}\lambda_{i,j,1}(t) &= \begin{cases} \pm \mathcal{O}(\delta_r^3) & \text{if } i \neq j \\ 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \pm \mathcal{O}(\delta_p^2) & \text{if } i = j \end{cases} \\ \lambda_{i,j,2}(t) &= \begin{cases} \pm \mathcal{O}(\delta_r^3) & \text{if } i \neq j \\ 6C_{S,2}I_{i,i}^{(3)}(t) \pm \mathcal{O}(\delta_p^3) & \text{if } i = j \end{cases} \\ \lambda_{i,j,3}(t) &= \begin{cases} \pm \mathcal{O}(\delta_r^3) & \text{if } i \neq j \\ 6C_{S,2}I_{i,i}^{(3)}(t) \pm \mathcal{O}(\delta_p^3) & \text{if } i = j \end{cases} \\ \lambda_{i,j,4}(t) &= \begin{cases} \pm \mathcal{O}(\delta_r^3) & \text{if } i \neq j \\ 6C_{S,2}I_{i,i}^{(3)}(t) \pm \mathcal{O}(\delta_p^3) & \text{if } i = j \end{cases} \\ \lambda_{i,j,5}(t) &= \begin{cases} \pm \mathcal{O}(\delta_r^2) & \text{if } i \neq j \\ 2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \pm \mathcal{O}(\delta_p^2) & \text{if } i = j \end{cases}\end{aligned}$$

Here $C_{S,2} = \sum_{k=0}^{\infty} \frac{c_{k+1}^2 + c_k c_{k+2}}{k!}$.

Proof. We are going to use Lemma 21 with $\mathbf{v}_1 = \bar{\mathbf{v}}_i(t)$, $\mathbf{v}_2 = \bar{\mathbf{v}}_j(t)$ and $\mathbf{w}_1 = \bar{\mathbf{w}}_i(t)$, $\mathbf{w}_2 = \bar{\mathbf{w}}_j(t)$. In this case, we have that $\mathbf{v}_1^\top \mathbf{v}_2 = I_{i,j}^{(1)}(t)$, $\mathbf{w}_1^\top \mathbf{w}_2 = I_{i,j}^{(2)}(t)$, and $\mathbf{v}_1^\top \mathbf{w}_2 = I_{i,j}^{(3)}(t)$, $\mathbf{v}_2^\top \mathbf{w}_1 = I_{j,i}^{(3)}(t)$. Now, for $\lambda_{i,j,1}$, by definition, we have that

$$\lambda_{i,j,1}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell+1}}{k!\ell!} \mathbb{E}[He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_\ell(\bar{\mathbf{v}}_j(t)^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_j(t)^\top \mathbf{x})]$$

Invoking Lemma 21 with $h_k = c_{k+1}$, $h'_\ell = c_{\ell+1}$ gives

$$\lambda_{i,j,1}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} I_{i,j}^{(1)}(t)^k I_{i,j}^{(2)}(t)^3 \pm \mathcal{O}(\delta_r) I_{i,j}^{(2)}(t)^2 \pm \mathcal{O}(\delta_r^4)$$

Since $|I_{i,j}^{(2)}| \leq \delta_r$, we have that

$$\lambda_{i,j,1}(t) = \pm \mathcal{O}(\delta_r^3)$$

In the special case where $i = j$, we have that $I_{i,j}^{(1)} = I_{i,j}^{(2)} = 1$, and $\mathbf{v}_1^\top \mathbf{w}_1 = \mathbf{v}_1^\top \mathbf{w}_2 = \mathbf{v}_2^\top \mathbf{w}_1 = \mathbf{v}_2^\top \mathbf{w}_2 = I_{i,i}^{(3)}$. Therefore

$$\lambda_{i,j,1}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \pm \mathcal{O}(\delta_p^2)$$

For $\lambda_{i,j,5}$, by definition, we have that

$$\lambda_{i,j,5}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k!\ell!} \mathbb{E}[He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_\ell(\bar{\mathbf{v}}_j(t)^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_j(t)^\top \mathbf{x})]$$

Invoking Lemma 21 with $h_k = c_k$, $h'_\ell = c_k$ gives

$$\lambda_{i,j,5}(t) = 2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} I_{i,j}^{(1)}(t)^k I_{i,j}^{(2)}(t)^2 \pm \mathcal{O}(\delta_r^2) = \pm \mathcal{O}(\delta_r^2)$$

In the case where $i = j$, we have that

$$\lambda_{i,j,5}(t) = 2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \pm \mathcal{O}(\delta_p^2)$$

For $\lambda_{i,j,2}$, by definition, we have that

$$\lambda_{i,j,2}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell}}{k!\ell!} \mathbb{E} [He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_{\ell}(\bar{\mathbf{v}}_j(t)^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_j(t)^\top \mathbf{x})]$$

Invoking Lemma 21 with $h_k = c_{k+1}, h'_{\ell} = c_{\ell}$ gives

$$\lambda_{i,j,2}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+2}c_k}{k!} I_{i,j}^{(1)}(t)^k I_{i,j}^{(2)}(t)^2 I_{i,j}^{(3)}(t) + 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} I_{i,j}^{(1)}(t)^k I_{i,j}^{(2)}(t)^2 I_{j,j}^{(3)}(t) \pm \mathcal{O}(\delta_r^3) = \pm \mathcal{O}(\delta_r^3)$$

In the case where $i = j$, we have that $I_{i,j}^{(1)}(t) = I_{i,j}^{(2)}(t) = 1$. Therefore

$$\lambda_{i,j,2}(t) = 6I_{i,i}^{(3)}(t) \sum_{k=0}^{\infty} \frac{c_k c_{k+2} + c_{k+1}^2}{k!} \pm \mathcal{O}(\delta_p^3)$$

For $\lambda_{i,j,3}$, by definition, we have that

$$\lambda_{i,j,3}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_k c_{\ell+1}}{k!\ell!} \mathbb{E} [He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_{\ell}(\bar{\mathbf{v}}_j(t)^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_j(t)^\top \mathbf{x})]$$

Invoking Lemma 21 with $h_k = c_k, h'_{\ell} = c_{\ell+1}$ gives

$$\lambda_{i,j,3}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} I_{i,j}^{(1)}(t)^k I_{i,j}^{(2)}(t)^2 I_{i,j}^{(3)}(t) + 6 \sum_{k=0}^{\infty} \frac{c_{k+2}c_k}{k!} I_{i,j}^{(1)}(t)^k I_{i,j}^{(2)}(t)^2 I_{j,j}^{(3)}(t) \pm \mathcal{O}(\delta_r^3) = \pm \mathcal{O}(\delta_r^3)$$

In the case where $i = j$, we have that $I_{i,j}^{(1)}(t) = I_{i,j}^{(2)}(t) = 1$. Therefore

$$\lambda_{i,j,3}(t) = 6I_{i,i}^{(3)}(t) \sum_{k=0}^{\infty} \frac{c_k c_{k+2} + c_{k+1}^2}{k!} \pm \mathcal{O}(\delta_p^3)$$

Lastly, for $\lambda_{i,j,4}$, we have that

$$\lambda_{i,j,4}(t) = \sum_{k,\ell=0}^{\infty} \frac{c_{k+1}c_{\ell}}{k!\ell!} \mathbb{E} [He_k(\bar{\mathbf{v}}_i(t)^\top \mathbf{x}) He_{\ell}(\bar{\mathbf{v}}_j(t)^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_i(t)^\top \mathbf{x}) He_2(\bar{\mathbf{w}}_j(t)^\top \mathbf{x})]$$

Here we need to apply Lemma 21 with $\mathbf{v}_1 = \bar{\mathbf{v}}_i(t), \mathbf{v}_2 = \bar{\mathbf{v}}_j(t), \mathbf{w}_1 = \bar{\mathbf{w}}_j(t), \mathbf{w}_2 = \bar{\mathbf{w}}_i(t)$ and $h_k = c_{k+1}, h'_{\ell} = c_{\ell}$. This gives that

$$\lambda_{i,j,4}(t) = 6 \sum_{k=0}^{\infty} \frac{c_{k+2}c_k}{k!} I_{i,j}^{(1)}(t)^k I_{i,j}^{(2)}(t)^2 I_{i,i}^{(1)}(t) + 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} I_{i,j}^{(1)}(t)^k I_{i,j}^{(2)}(t)^2 I_{j,i}^{(3)}(t) \pm \mathcal{O}(\delta_r^3) = \pm \mathcal{O}(\delta_r^3)$$

In the case where $i = j$, we have that $I_{i,j}^{(1)}(t) = I_{i,j}^{(2)}(t) = 1$. Therefore

$$\lambda_{i,j,4}(t) = 6I_{i,i}^{(3)}(t) \sum_{k=0}^{\infty} \frac{c_k c_{k+2} + c_{k+1}^2}{k!} \pm \mathcal{O}(\delta_p^3)$$

□

With the above lemmas that studies $\lambda_{i,j,\ell}$ s and $\hat{\lambda}_{i,j,\ell}$ s, we are ready to analyze the dynamics of $\gamma_{i,j}^{(1)}, \gamma_{i,j}^{(2)}, \zeta_{i,j}^{(1)}, \zeta_{i,j}^{(2)}$, and $I_{i,j}^{(1)}, I_{i,j}^{(2)}, I_{i,j}^{(3)}$. In particular, we will fix any $\ell \in [m^*]$, and assumes the inductive hypothesis.

Lemma 7. Let $\varepsilon_{\mathcal{A},\ell}^{(2)}(t), \varepsilon_{\mathcal{B},\ell}^{(3)}(t)$, and $\varepsilon_{5,\ell}(t)$ be defined in Definition 1 with $\varepsilon_{5,\ell}(t) \leq \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)\right)$. Then the gradient alignment $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^*$ and $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^*$ satisfies

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^* = -\frac{1}{a_i} \cdot \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ -\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i,j}^{(1)}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^*, j \in [m^*] \setminus \mathcal{C}_\ell \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)\varepsilon_{5,\ell}(t)\right) & \\ \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^2\right) \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^*, j = j_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)\varepsilon_{5,\ell}(t)\right) & \end{cases}$$

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^* = -\frac{9}{b_i} \cdot \begin{cases} \hat{\lambda}_{i,j,5}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) & \\ \hat{\lambda}_{i,j,5}(t) - \lambda_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \gamma_{i,j}^{(2)}(t) & \text{if } i = i_\ell^*, j \in [m^*] \setminus \mathcal{C}_\ell \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) & \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)\varepsilon_{5,\ell}(t)\right) & \\ \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2\right) \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) & \text{if } i = i_\ell^*, j = j_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)\varepsilon_{5,\ell}(t)\right) & \end{cases}$$

and in particular, for the case $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_{j_\ell^*}^*$, we have that

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_{j_\ell^*}^* = \frac{9}{b_i} \lambda_{i, i_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) - \frac{9}{b_i} \cdot \begin{cases} \hat{\lambda}_{i,j,5}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i,j,5}(t) - \lambda_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \gamma_{i,j}^{(2)}(t) & \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)\varepsilon_{5,\ell}(t)\right) & \end{cases}$$

Further more, the mis-alignment terms $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^*$ and $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^*$ satisfies

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^* = -\frac{1}{a_i} \cdot \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i,j}^{(2)}(t) - \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \zeta_{i,j}^{(1)}(t) & \\ -36C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) & \text{if } i = i_\ell^*, j \neq j_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \\ 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) + \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) - \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) & \\ -36C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) & \text{if } i = i_\ell^*, j = j_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \end{cases}$$

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^* = \frac{3}{b_i} \cdot \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i,j}^{(1)}(t) - 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \zeta_{i,j}^{(2)}(t) & \\ -12C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) & \text{if } i = i_\ell^*, j \neq j_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) + \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) - 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) & \\ -12C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) & \text{if } i = i_\ell^*, j = j_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \end{cases}$$

The self-alignments $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j, \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j$, in the case of $i \neq j$, are given by

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j = -\frac{1}{a_i} \cdot \begin{cases} \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{j_\ell^*}^{(1)}(t) - \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) I_{i,j}^{(1)}(t) \\ \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{5,\ell}(t) \right) & \end{cases}$$

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j = -\frac{9}{b_i} \cdot \begin{cases} \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) & \text{if } i, j \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{j_\ell^*}^{(2)}(t) - \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i,j}^{(2)}(t) \\ \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{5,\ell}(t) \right) \\ \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) & \text{if } j = i_\ell^* \\ \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) & \end{cases}$$

The self-alignments $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j, \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j$, in the case of $i \neq j$, are given by

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j = -\frac{1}{a_i} \cdot \begin{cases} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \left(\zeta_{j_\ell^*}^{(2)}(t) - \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) I_{i,j}^{(3)}(t) \right) \\ \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{5,\ell}(t) \right) \\ \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases}$$

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j = -\frac{9}{b_i} \cdot \begin{cases} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \left(\zeta_{j_\ell^*}^{(1)}(t) - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i,j}^{(3)}(t) \right) \\ \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{5,\ell}(t) \right) \\ \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases}$$

Lastly, the self-alignments $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j, \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j$, in the case of $i = j \in [m] \setminus \mathcal{R}_{\ell-1}$, are given by

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j = -\frac{1}{a_i} \left(\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \zeta_{j_\ell^*}^{(2)}(t) + 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) + 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right) \mathbb{I}\{i = i_\ell^*\}$$

$$- \frac{1}{a_i} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \mathbb{I}\{i = i_\ell^*\} - \frac{36}{a_i} C_{S,2} I_{i,i}^{(3)}(t) \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right)$$

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j = -\frac{3}{b_i} \cdot \left(3\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \zeta_{j_\ell^*}^{(1)}(t) + \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) + \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right) \mathbb{I}\{i = i_\ell^*\}$$

$$- \frac{9}{b_i} \cdot \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \mathbb{I}\{i = i_\ell^*\} - 36C_{S,2} I_{i,i}^{(3)}(t) \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right)$$

Proof. In the following of the proof we will assume that $i \in [m] \setminus \mathcal{R}_{\ell-1}$ and $j \in [m^*]$. For $\hat{\lambda}_{i,j,1}, \dots, \hat{\lambda}_{i,j,5}$, we apply Lemma 5 with

$$\delta_r = \varepsilon_{\mathcal{B},\ell}^{(1)}(t); \text{ for } i \in [m] \setminus \mathcal{R}_{\ell-1}, j \in [m^*];$$

For $\lambda_{i,j,1}, \dots, \lambda_{i,j,5}$, we apply

$$\delta_r = \begin{cases} \varepsilon_{\mathcal{B},\ell}^{(1)}(t); & \text{if } i, j \in [m] \setminus \mathcal{R}_{\ell-1}, i \neq j \\ \varepsilon_{\mathcal{F},\ell}(t); & \text{if } i \in \mathcal{R}_{\ell-1} \vee j \in \mathcal{R}_{\ell-1}, i \neq j \end{cases} \leq \varepsilon_{\mathcal{A},\ell}^{(1)}(t); \quad \delta_p = \begin{cases} \varepsilon_{\mathcal{B},\ell}^{(3)}(t) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \varepsilon_{5,\ell}(t) & \text{if } i = i_\ell^* \end{cases}$$

For $\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}, \dots, \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}$ for $\ell' < \ell$, we have that

$$\delta_r = \varepsilon_{\mathcal{F},\ell}(t) \leq \varepsilon_{\mathcal{A},\ell}^{(1)}(t);$$

For $\lambda_{i_{\ell'}^*, i_{\ell'}^*, 1}, \dots, \lambda_{i_{\ell'}^*, i_{\ell'}^*, 5}$, we have that

$$\delta_r = \delta_p = \varepsilon_{\mathcal{F}, \ell}(t) \leq \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)$$

Moreover, we also have that

$$\begin{aligned} \left| \gamma_{i,j}^{(1)}(t) \right|, \left| \gamma_{i,j}^{(2)}(t) \right| &\leq \begin{cases} \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) & \text{if } i \neq i_{\ell}^* \\ \varepsilon_{\mathcal{A}, \ell}^{(1)}(t) & \text{if } i = i_{\ell}^* \end{cases} \forall j \in [m^*], (i, j) \neq (i_{\ell'}^*, j_{\ell'}^*) \forall \ell' \leq \ell \\ \left| \zeta_{i,j}^{(1)}(t) \right|, \left| \zeta_{i,j}^{(2)}(t) \right| &\leq \varepsilon_{\mathcal{A}, \ell}^{(1)}(t); \forall i \in [m], j \in [m^*] \\ \left| I_{i,j}^{(1)}(t) \right|, \left| I_{i,j}^{(2)}(t) \right|, \left| I_{i,j}^{(3)}(t) \right| &\leq \varepsilon_{\mathcal{A}, \ell}^{(1)}(t); \forall i, j \in [m], i \neq j \\ \left| I_{i,i}^{(3)}(t) \right| &\leq \begin{cases} \varepsilon_{\mathcal{F}, \ell}(t) & \text{if } i \in \mathcal{R}_{\ell-1} \\ \varepsilon_{5, \ell}(t) & \text{if } i = i_{\ell}^* \\ \varepsilon_{\mathcal{B}, \ell}^{(3)}(t) & \text{if } i \in [m] \setminus \mathcal{R}_{\ell} \end{cases} \end{aligned}$$

This gives that for all $j \in [m^*]$ and such that $(i, j) \neq (i_{\ell}^*, j_{\ell}^*)$

$$\hat{\lambda}_{i,j,1}(t) \leq \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3\right) & \text{if } i \neq i_{\ell}^* \\ \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) & \text{if } i = i_{\ell}^* \end{cases}; \hat{\lambda}_{i,j,5}(t) \leq \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2\right) & \text{if } i \neq i_{\ell}^* \\ \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2\right) & \text{if } i = i_{\ell}^* \end{cases}$$

We will analyze each dynamic separately. However, we should notice some common terms that appears in the dynamics.

Common terms. To start, let's tackle some common terms in the dynamics we are interest in. In particular, we have that for $i \in [m] \setminus \mathcal{R}_{\ell-1}$,

$$\begin{aligned} \sum_{r=1}^m \lambda_{i,r,1}(t) I_{i,r}^{(1)}(t) &= \lambda_{i,i,1}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) \\ \sum_{r=1}^m \lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) &= \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{B}, \ell}^{(3)}(t)^2 + m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_{\ell} \\ \mathcal{O}\left(\varepsilon_{5, \ell}(t)^2 + m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) & \text{if } i = i_{\ell}^* \end{cases} \\ \sum_{r=1}^m \lambda_{i,r,4}(t) I_{i,r}^{(3)}(t) &= \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{B}, \ell}^{(3)}(t)^2 + m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_{\ell} \\ \mathcal{O}\left(\varepsilon_{5, \ell}(t)^2 + m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) & \text{if } i = i_{\ell}^* \end{cases} \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,1}(t) \gamma_{i,r}^{(t)}(t) &= \begin{cases} \lambda_{i_{\ell}^*, j_{\ell}^*, 1}(t) \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) & \text{if } i = i_{\ell}^* \\ \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_{\ell} \end{cases} \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) &= \begin{cases} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2\right) \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) & \text{if } i = i_{\ell}^* \\ \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_{\ell} \end{cases} \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,4}(t) \zeta_{i,r}^{(1)}(t) &= \begin{cases} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2\right) \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3\right) & \text{if } i = i_{\ell}^* \\ \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_{\ell} \end{cases} \end{aligned}$$

Therefore, we have that

$$\begin{aligned}
 & \sum_{r=1}^m \lambda_{i,r,1}(t) I_{i,r}^{(1)}(t) + 3 \sum_{r=1}^m \lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) + 3 \sum_{r=1}^m \lambda_{i,r,4}(t) I_{i,r}^{(3)}(t) \\
 &= \lambda_{i,i,1} \pm \begin{cases} \mathcal{O}\left(\varepsilon_{5,\ell}(t)^2 + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \\ \mathcal{O}\left(\varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2 + m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases} \\
 & \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,1}(t) \gamma_{i,r}^{(t)}(t) + 3 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) + \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,4}(t) \zeta_{i,r}^{(1)}(t) \\
 &= \begin{cases} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases}
 \end{aligned} \tag{12}$$

Moreover, we can also compute that

$$\begin{aligned}
 & \sum_{r=1}^m \lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) = \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2 + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(\varepsilon_{5,\ell}(t)^2 + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \end{cases} \\
 & \sum_{r=1}^m \lambda_{i,r,3}(t) I_{r,i}^{(3)}(t) = \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2 + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(\varepsilon_{5,\ell}(t)^2 + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \end{cases} \\
 & \sum_{r=1}^m \lambda_{i,r,5}(t) I_{r,i}^{(2)}(t) = \lambda_{i,i,5}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) \\
 & \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) = \begin{cases} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases} \\
 & \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,3}(t) \zeta_{i,r}^{(2)}(t) = \begin{cases} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases} \\
 & \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,5}(t) \gamma_{i,r}^{(2)}(t) = \begin{cases} \lambda_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases}
 \end{aligned}$$

This gives that

$$\begin{aligned}
 & \sum_{r=1}^m \lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) + \sum_{r=1}^m \lambda_{i,r,3}(t) I_{r,i}^{(3)}(t) + 3 \sum_{r=1}^m \lambda_{i,r,5}(t) I_{r,i}^{(2)}(t) \\
 &= \lambda_{i,i,5}(t) \pm \begin{cases} \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{5,\ell}(t)^2\right) & \text{if } i = i_\ell^* \end{cases} \\
 & \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) + \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,3}(t) \zeta_{i,r}^{(2)}(t) + \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,5}(t) \gamma_{i,r}^{(2)}(t) \\
 &= \begin{cases} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \\ \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases}
 \end{aligned} \tag{13}$$

Now we are ready to analyze the dynamics.

Analysis of $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^*$. To analyze $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^*$, we first compute the following quantities for $i \in [m] \setminus$

$\mathcal{R}_{\ell-1}$

$$\begin{aligned} \hat{\lambda}_{i,j,1}(t) &= \begin{cases} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^*, j \neq j_\ell^* \end{cases}; \\ \sum_{r=1}^m \lambda_{i,r,1}(t) \gamma_{r,j}^{(1)}(t) &= \lambda_{i,i,1}(t) \gamma_{i,j}^{(1)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) \\ \sum_{r=1}^m \lambda_{i,r,2}(t) \zeta_{i,j}^{(2)}(t) &= \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\ \sum_{r=1}^m \lambda_{i,r,4}(t) \zeta_{r,j}^{(2)}(t) &= \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) \zeta_{i,j}^{(2)}(t) &= \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases} \end{aligned}$$

Therefore, combining with (12), and noticing that the term $\lambda_{i,i,1}(t) \gamma_{i,j}^{(1)}(t)$ cancels out, we have that

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^* = -\frac{1}{a_i} \cdot \begin{cases} \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ -\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i,j}^{(1)}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^*, j \in [m^*] \setminus \mathcal{C}_\ell \\ \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t) \varepsilon_{5,\ell}(t)\right) & \\ \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^2\right) \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^*, j = j_\ell^* \\ \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t) \varepsilon_{5,\ell}(t)\right) & \end{cases}$$

Analysis of $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^*$. To analyze $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^*$, we first compute the following quantities

$$\begin{aligned} \sum_{r=1}^m \lambda_{i,r,5}(t) \gamma_{r,j}^{(2)}(t) &= \lambda_{i,i,5}(t) \gamma_{i,j}^{(5)}(t) + \lambda_{i,i_\ell^*,5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \mathbb{I}\{j = j_\ell^*, \ell' < \ell\} \\ &\pm \begin{cases} \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \mathbb{I}\{j = j_\ell^*\} & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\ \sum_{r=1}^m \lambda_{i,r,2}(t) \zeta_{i,j}^{(1)}(t) &= \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\ \sum_{r=1}^m \lambda_{i,r,3}(t) \zeta_{r,j}^{(1)}(t) &= \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) \zeta_{i,j}^{(1)}(t) &= \begin{cases} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases} \end{aligned}$$

Therefore, combining with (13), and noticing that the term $\lambda_{i,i,5}(t) \gamma_{i,j}^{(2)}(t)$ cancels out, we have that

$$\begin{aligned} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_{j_\ell'}^* &= \frac{9}{b_i} \lambda_{i,i_\ell^*,5}(t) \gamma_{i_\ell^*, j_\ell'}^{(2)}(t) \\ &- \frac{9}{b_i} \cdot \begin{cases} \hat{\lambda}_{i,j,5}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i,j,5}(t) - \lambda_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \gamma_{i,j}^{(2)}(t) & \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) & \end{cases} \end{aligned}$$

Moreover, we have that

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^* = -\frac{9}{b_i} \cdot \begin{cases} \hat{\lambda}_{i,j,5}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) \\ \quad \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i,j,5}(t) - \lambda_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \gamma_{i,j}^{(2)}(t) \\ \quad \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) & \text{if } i = i_\ell^*, j \in [m^*] \setminus \mathcal{C}_\ell \\ \quad \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)\varepsilon_{5,\ell}(t)\right) \\ \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)\right) \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) & \text{if } i = i_\ell^*, j = j_\ell^* \\ \quad \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)\varepsilon_{5,\ell}(t)\right) \end{cases}$$

Analysis of $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^*$. To analyze $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^*$, we first compute the following quantities for $i \in [m] \setminus \mathcal{R}_{\ell-1}$:

$$\begin{aligned} \hat{\lambda}_{i,j,4}(t) &= \begin{cases} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^*, j \neq j_\ell^* \end{cases} \\ \sum_{r=1}^m \lambda_{i,r,4}(t) \gamma_{r,j}^{(2)}(t) &= 6C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\right) \\ \sum_{r=1}^m \lambda_{i,r,1}(t) \zeta_{r,j}^{(1)}(t) &= \lambda_{i,i,1}(t) \zeta_{i,j}^{(1)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) \\ \sum_{r=1}^m \lambda_{i,r,2}(t) \gamma_{i,j}^{(2)}(t) &= 6C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\right) \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) \gamma_{i,j}^{(2)}(t) &= \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i,j}^{(2)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \end{cases} \end{aligned}$$

Combining with (12) and noticing that the term $\lambda_{i,i,1}(t) \zeta_{i,j}^{(1)}(t)$ cancels out, we have that

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{w}}_j^* = -\frac{1}{a_i} \cdot \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i,j}^{(2)}(t) - \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \zeta_{i,j}^{(1)}(t) \\ \quad - 36C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) & \text{if } i = i_\ell^*, j \neq j_\ell^* \\ \quad \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) \\ 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) + \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) - \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \\ \quad - 36C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) & \text{if } i = i_\ell^*, j = j_\ell^* \\ \quad \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) \end{cases}$$

Analysis of $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^*$. To analyze $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^*$, we first compute the following for $i \in [m] \setminus \mathcal{R}_{\ell-1}$:

$$\begin{aligned} \hat{\lambda}_{i,j,3}(t) &= \begin{cases} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^*, j \neq j_\ell^* \end{cases} \\ \sum_{r=1}^m \lambda_{i,r,3}(t) \gamma_{r,j}^{(1)}(t) &= 6C_{S,2} \gamma_{i,j}^{(1)}(t) I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\right) \\ \sum_{r=1}^m \lambda_{i,r,2}(t) \gamma_{i,j}^{(1)}(t) &= 6C_{S,2} \gamma_{i,j}^{(1)}(t) I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\right) \\ \sum_{r=1}^m \lambda_{i,r,5}(t) \zeta_{r,j}^{(2)}(t) &= \lambda_{i,i,5}(t) \zeta_{i,j}^{(2)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) \gamma_{i,j}^{(1)}(t) &= \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i,j}^{(1)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) & \text{if } i = i_\ell^* \end{cases} \end{aligned}$$

Combining with (13) and noticing that the term $\lambda_{i,i,5} \zeta_{i,j}^{(2)}(t)$ cancels out, we have that

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t)) \bar{\mathbf{v}}_j^* = \frac{3}{b_i} \cdot \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i,j}^{(1)}(t) - 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \zeta_{i,j}^{(2)}(t) \\ \quad - 12C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) & \text{if } i = i_\ell^*, j \neq j_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) + \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) - 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \\ \quad - 12C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)}(t) & \text{if } i = i_\ell^*, j = j_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2\right) \end{cases}$$

Analysis of $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j$. To analyze $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j$, we first compute that for $i, j \in [m] \setminus \mathcal{R}_{\ell-1}, i \neq j$

$$\begin{aligned} \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,1}(t) \gamma_{j,r}^{(1)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{j, j_\ell^*}^{(1)}(t) \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{j,i}^{(3)}(t) &= \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,4}(t) \zeta_{j,r}^{(1)}(t) &= \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right) \\ \sum_{r=1}^m \lambda_{i,r,2}(t) I_{j,i}^{(3)}(t) &= \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \neq i_\ell^* \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\ \sum_{r=1}^m \lambda_{i,r,4}(t) I_{j,r}^{(3)}(t) &= \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \neq i_\ell^* \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \end{aligned}$$

Combining with (12) and noticing that the term $\lambda_{i,i,1}(t) I_{i,j}^{(1)}(t)$ and the term $\lambda_{j,j,1}(t) I_{i,j}^{(1)}(t)$ cancels out, we have

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j = -\frac{1}{a_i} \cdot \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{j, j_\ell^*}^{(1)}(t) - \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) I_{i,j}^{(1)}(t) \\ \quad \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) \end{cases}$$

Analysis of $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j$. To analyze $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j$, we first compute that for $i, j \in [m] \setminus \mathcal{R}_{\ell-1}, i \neq j$

$$\begin{aligned}
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{i,j}^{(3)}(t) &= \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,3}(t) \zeta_{j,r}^{(2)}(t) &= \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,5}(t) \gamma_{j,r}^{(2)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 5} \gamma_{j_\ell^*}^{(2)}(t) \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \mathbb{I}\{j = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\
 \sum_{r=1}^m \lambda_{i,r,2}(t) I_{i,j}^{(3)}(t) &= \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \neq i_\ell^* \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\
 \sum_{r=1}^m \lambda_{i,r,3}(t) I_{r,j}^{(3)}(t) &= \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \neq i_\ell^* \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\
 \sum_{r=1}^m \lambda_{i,r,5}(t) I_{r,j}^{(2)}(t) &= \lambda_{i,i,5}(t) I_{i,j}^{(2)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right)
 \end{aligned}$$

Combining with (13) and noticing that the term $\lambda_{i,i,5}(t) I_{i,j}^{(2)}(t)$ and the term $\lambda_{j,j,5}(t) I_{i,j}^{(2)}(t)$ cancels out, we have

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j = -\frac{9}{b_i} \cdot \begin{cases} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i, j \in [m] \setminus \mathcal{R}_\ell \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{j_\ell^*}^{(2)}(t) - \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i,j}^{(2)}(t) \\ \quad \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{5,\ell}(t)\right) & \\ \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) & \text{if } j = i_\ell^* \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \end{cases}$$

Analysis of $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j$. The analysis of $\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j$ will be separated into two cases. First, regardless of the cases, we have that

$$\begin{aligned}
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,1}(t) \zeta_{j,r}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \zeta_{j_\ell^*}^{(2)}(t) \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\
 \sum_{r=1}^m \lambda_{i,r,1}(t) I_{r,j}^{(3)}(t) &= \lambda_{i,i,1}(t) I_{i,j}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right)
 \end{aligned}$$

We first analyze the case $i = j$, and then we dive into $i \neq j$. In the case where $i = j$, we have that

$$\begin{aligned}
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{i,j}^{(2)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,4}(t) \gamma_{j,r}^{(2)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\
 \sum_{r=1}^m \lambda_{i,r,2}(t) I_{i,j}^{(2)}(t) &= 6C_{S,2} I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\
 \sum_{r=1}^m \lambda_{i,r,4}(t) I_{r,j}^{(2)}(t) &= 6C_{S,2} I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right)
 \end{aligned}$$

Thus, for the case $i = j$, we have that

$$\begin{aligned} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j &= -\frac{1}{a_i} \left(\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \zeta_{j, j_\ell^*}^{(2)}(t) + 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) + 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right) \mathbb{I}\{i = i_\ell^*\} \\ &\quad - \frac{1}{a_i} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \mathbb{I}\{i = i_\ell^*\} - \frac{36}{a_i} C_{S,2} I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \end{aligned}$$

Now, for the case $i \neq j$, we can compute that

$$\begin{aligned} \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{i,j}^{(2)}(t) &= \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,4}(t) \gamma_{j,r}^{(2)}(t) &= \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\ \sum_{r=1}^m \lambda_{i,r,2}(t) I_{i,j}^{(2)}(t) &= \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \\ \sum_{r=1}^m \lambda_{i,r,4}(t) I_{r,j}^{(2)}(t) &= \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \pm \begin{cases} \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) & \text{if } i = i_\ell^* \end{cases} \end{aligned}$$

Therefore, combining with (12) gives

$$\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j = -\frac{1}{a_i} \cdot \begin{cases} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \left(\zeta_{j, j_\ell^*}^{(2)}(t) - \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) I_{i,j}^{(3)}(t) \right) \\ \quad \pm \mathcal{O}\left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \quad \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t)\right) \\ \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t)\right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases}$$

Analysis of $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j$. The analysis of $\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j$ will also be separated into two cases. First, regardless of the cases, we have that

$$\begin{aligned} \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,5}(t) \zeta_{j,r}^{(1)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \zeta_{j, j_\ell^*}^{(1)}(t) \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\ \sum_{r=1}^m \lambda_{i,r,5}(t) I_{j,r}^{(3)}(t) &= \lambda_{i,i,5}(t) I_{j,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \end{aligned}$$

We first analyze the case $i = j$, and then we dive into $i \neq j$. In the case where $i = j$, we have that

$$\begin{aligned} \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{j,i}^{(1)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,3}(t) \gamma_{j,r}^{(1)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\ \sum_{r=1}^m \lambda_{j,r,2}(t) I_{i,j}^{(1)}(t) &= 6C_{S,2} I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \\ \sum_{r=1}^m \lambda_{j,r,3}(t) I_{i,r}^{(1)}(t) &= 6C_{S,2} I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3\right) \end{aligned}$$

Thus, for the case $i = j$, we have that

$$\begin{aligned} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j &= -\frac{3}{b_i} \cdot \left(3\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \zeta_{j, j_\ell^*}^{(1)}(t) + \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) + \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right) \mathbb{I}\{i = i_\ell^*\} \\ &\quad - \frac{9}{b_i} \cdot \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \mathbb{I}\{i = i_\ell^*\} - 36C_{S,2} I_{i,i}^{(3)}(t) \pm \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2\right) \end{aligned}$$

Now, for the case $i \neq j$, we can compute that

$$\begin{aligned}
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{j,i}^{(1)}(t) &= \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,3}(t) \gamma_{j,r}^{(1)}(t) &= \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \mathbb{I}\{i = i_\ell^*\} \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\
 \sum_{r=1}^m \lambda_{i,r,2}(t) I_{j,i}^{(1)}(t) &= \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \pm \begin{cases} \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t) \right) & \text{if } i = i_\ell^* \end{cases} \\
 \sum_{r=1}^m \lambda_{i,r,3}(t) I_{j,r}^{(1)}(t) &= \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \pm \begin{cases} \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \\ \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t) \right) & \text{if } i = i_\ell^* \end{cases}
 \end{aligned}$$

Therefore, combining with (12) and (13) gives

$$\nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j = -\frac{9}{b_i} \cdot \begin{cases} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \left(\zeta_{j, j_\ell^*}^{(1)}(t) - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i, j}^{(3)}(t) \right) \\ \quad \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 & \text{if } i = i_\ell^* \\ \quad \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t) \right) \\ \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) & \text{if } i \in [m] \setminus \mathcal{R}_\ell \end{cases}$$

□

A.5 Establishing the Inductive Hypothesis: Phase 1

Starting from this section, we assume that for a fixed $\ell \in [m^*]$, the inductive hypothesis holds. That is, the condition of Lemma 7 holds. Then we shall analyze the convergence for that ℓ to prove the inductive hypothesis and establish convergence. Notice that, by the statement of the inductive hypothesis, the case $\ell = 1$ naturally satisfies it, thus requiring no additional proof. Therefore, we focus on the case of a general fixed ℓ , and the proof will be divided into two phase. Phase 1 (this section) show that there exist some T^* such that $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(T^*)$ and $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(T^*)$ are some constant close to 1, while $\varepsilon_{\mathcal{A},\ell}^{(2)}(t)$, $\varepsilon_{\mathcal{B},\ell}^{(3)}(t)$, and $\varepsilon_{5,\ell}(t)$ are $\mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$.

We denote the following time

$$\begin{aligned}
 T_c &:= \min \{t \geq 0 : A_1(t) \wedge A_2(t) \wedge A_3(t)\} \\
 A_1(t) &= \left\{ \min \left\{ \varepsilon_{\mathcal{B},\ell}^{(2)}(t), \varepsilon_{3,\ell}(t), \varepsilon_{4,\ell}(t) \right\} > \frac{\beta_6 m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right\}; \quad A_2(t) = \left\{ \varepsilon_{\mathcal{B},\ell}^{(3)}(t) \geq \beta_6 m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right\} \\
 A_3(t) &= \left\{ \varepsilon_{5,\ell}(t) \geq \beta_6 \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \right) \right\}
 \end{aligned} \tag{14}$$

By the initialization property, we have that

$$\varepsilon_{\mathcal{B},\ell}^{(1)}(t), \varepsilon_{3,\ell}(t), \varepsilon_{4,\ell}(t) \leq \frac{\beta_3}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}}; \quad \varepsilon_{\mathcal{B},\ell}^{(3)}(0) = \varepsilon_{5,\ell}(0) = 0$$

Therefore, $T_c > 0$. Moreover, for all $t \leq T_c$, by the inductive hypothesis, we have that

$$\varepsilon_{\mathcal{B},\ell}^{(1)}(t), \varepsilon_{3,\ell}(t), \varepsilon_{4,\ell}(t) \leq \frac{\beta_6 m^2}{\delta_{\mathbb{P}} \sqrt{d}} \Rightarrow \varepsilon_{\mathcal{F},\ell}(t) \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$$

which implies that $\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$. Also, by definition we have that $\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \leq \varepsilon_{\mathcal{A},\ell}^{(2)}(t)$.

Growth of $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)$. We will start with analyzing $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)$. For any $\xi \in (0, 1)$, recall the definition of $T_\ell(\xi)$ in Definition 2. We should notice that, although by definition we have $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(0) > 0$, it is not always the case that

$\gamma_{i_\ell^*, j_\ell^*}^{(1)}(0)$ will also be positive. Therefore, we also need to control the "negativeness" of $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(0)$ when analyzing the growth of $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)$. To do this, we present the following lemma.

Lemma 8. *Suppose that the inductive hypothesis in Condition 2 and the initialization condition in Condition 1 holds. Let T_c be defined in (14) and $T(\xi)$ in Definition 2. If $T_c \geq T_\ell(\frac{1}{2})$, then we have that*

$$\begin{aligned} \frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) &\geq 0; \lambda_{i_\ell^*, j_\ell^*, 5}(t) \geq 0; \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \geq -\frac{\beta_6}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}}; \forall t \leq \min\{T_c, T(\xi)\} \\ \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t + T_0) &\geq \left(\gamma_{i_\ell^*, j_\ell^*}^{(2)}(T_0)^{-1} - \frac{18}{b_{i_\ell^*}^2} \left(c_0^2 (1 - \xi^2) - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right) t \right)^{-1} \\ T_\ell(\xi) &\leq \frac{b_{i_\ell^*}^2}{18} \left(c_0^2 (1 - \xi) - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right)^{-1} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^{-1} \end{aligned}$$

for all ξ such that $(1 - \xi^2)^{-1} \leq \mathcal{O}(1)$.

Proof. To start, we lower bound the time-derivative of $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)$. Let T'_c be defined as

$$T'_c = \min \left\{ t \geq 0 : \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) < -\frac{\beta_6}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right\}$$

Since $(\mathbf{v}_{i_\ell^*}^\top \bar{\mathbf{v}}_j^*) \leq \frac{\beta_3}{d} \log \frac{m}{\delta_{\mathbb{P}}}$ and $\|\mathbf{v}_{i_\ell^*}\|_2 \geq 1 - \beta_2 \delta_s$ for some $\beta_3 > 0, \beta_2 \leq o(1)$ and $\delta_s \leq \mathcal{O}(\frac{1}{m^2})$, we have that $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(0)^2 \leq \frac{2\beta_3}{d} \log \frac{m}{\delta_{\mathbb{P}}}$. For $\beta_6 \geq 4\sqrt{\beta_3}$, we must have that $T'_c > 0$. By Lemma 30, we have that for all $t \leq T'_c$, it holds that

$$\sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \in c_0^2 \pm \mathcal{O} \left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right); \quad \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \in c_1^2 \pm \mathcal{O} \left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right)$$

By Lemma 7, we have that

$$\begin{aligned} \frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) &= -\frac{1}{b_{i_\ell^*}^2} \nabla_{\mathbf{w}_{i_\ell^*}} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_{j_\ell^*}^* \\ &= \frac{9}{b_{i_\ell^*}^2} \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \varepsilon_{5, \ell}(t) \right) \end{aligned}$$

Therefore, with the form of $\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}$ from Lemma 5, we have that

$$\begin{aligned} \frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) &= \frac{18}{b_{i_\ell^*}^2} \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \varepsilon_{5, \ell}(t) \right) \\ &\quad \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \\ &\geq \frac{18}{b_{i_\ell^*}^2} \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 - \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 \sqrt{d}} \right) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \left(c_0^2 - \mathcal{O} \left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right) \right) \\ &\quad - \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \varepsilon_{5, \ell}(t) \right) \\ &\geq \frac{18}{b_{i_\ell^*}^2} \left(1 - \xi^2 - \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 \sqrt{d}} \right) \right) \left(c_0^2 - \mathcal{O} \left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \\ &\geq \frac{18}{b_{i_\ell^*}^2} \left(c_0^2 (1 - \xi^2) - \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 \sqrt{d}} \right) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \end{aligned} \tag{15}$$

Let ξ be given such that $(1 - \xi^2)^{-1} \leq \mathcal{O}(1)$. By the lower bound that $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^2 \geq \frac{\log m^*}{d}$, we have that

$$\frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0) \geq \beta_{\text{temp}} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^2 - \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \geq 0$$

for $d \geq \beta_5 \delta_{\mathbb{P}}^6 m^{16}$ for some $\beta_{\text{temp}} > 0$. This shows that $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \geq \frac{\log m^*}{d}$ for all $t \leq \min\{T_c, T'_c, T_\ell(\xi)\}$ for any $(1 - \xi^2)^{-1} \leq \mathcal{O}(1)$, and thus $\frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0$ for the same choice of t . This shows the first property. Moreover, by Lemma 5,

$$\begin{aligned} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) &= 2\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \pm \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^4\right) \\ &\geq 2\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \left(c_0^2 - \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 \sqrt{d}}\right)\right) - \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \\ &\geq 0 \end{aligned}$$

as $d \geq \beta_5 \delta_{\mathbb{P}}^6 m^{16}$. This shows the second property. To lower bound $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)$, we first write that, by Lemma 7

$$\begin{aligned} \frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) &= -\frac{1}{a_i} \nabla_{\mathbf{v}_{i_\ell^*}} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_{j_\ell^*}^* \\ &= \frac{1}{a_i^2} \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^2\right) \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \pm \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \varepsilon_{5, \ell}(t)\right) \end{aligned}$$

Therefore, we could notice that, since

$$\sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \geq c_1^2 - \mathcal{O}\left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}\right) \geq 0$$

for $d \geq \beta_5 m^{14} \geq \beta_5 \log \frac{m}{\delta_{\mathbb{P}}}^4$, we must have that for all $t \leq \min\{T_c, T'_c, T_\ell(\xi)\}$ for any $(1 - \xi^2)^{-1} \leq \mathcal{O}(1)$, either $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \geq \frac{1}{2}$, or

$$\begin{aligned} \frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) &\geq 6 \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k - \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \varepsilon_{5, \ell}(t)\right) \\ &\quad - \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \\ &\geq c_1^2 \left(\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) - \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2\right)\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 - \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \\ &\geq -\mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \end{aligned} \tag{16}$$

where the last inequality follows from the fact that

$$\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2 \leq \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \leq \left(\frac{\log m^*}{d}\right)^{\frac{1}{2}} \leq \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0) \leq \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)$$

Therefore, for all $t \leq \min\{T_c, T'_c, T_\ell(\xi)\}$ for any $(1 - \xi^2)^{-1} \leq \mathcal{O}(1)$ we must have that

$$\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \geq \gamma_{i_\ell^*, j_\ell^*}^{(1)}(0) - \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \cdot t \geq -\frac{\sqrt{2\beta_4}}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}$$

⁴since we require $m \geq \frac{\beta_4 \log m^*}{\delta_{\mathbb{P}}}$

Choosing $\beta_6 \geq 4\sqrt{\beta_4}$ gives that

$$T_{c'} \geq \min \left\{ T_c, T_\ell(\xi), \Omega \left(\frac{d}{m^7} \right) \right\}$$

Now we are going to lower bound $T_\ell(\xi)$. When $d \geq \beta_5 m^{15}$,

$$\frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq \frac{18}{b_{i_\ell^*}^2} \left(c_0^2 (1 - \xi^2) - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2$$

Solving the differential equation gives that

$$\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t + T_0) \geq \left(\gamma_{i_\ell^*, j_\ell^*}^{(2)}(T_0)^{-1} - \frac{18}{b_{i_\ell^*}^2} \left(c_0^2 (1 - \xi^2) - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right) t \right)^{-1}$$

for any $T_0 \geq 0$ and $t + T_0 \leq T_c$. This gives that, if $T(\xi) \leq T_c$, then

$$\begin{aligned} T_\ell(\xi) &\leq \frac{b_{i_\ell^*}^2}{18} \left(c_0^2 (1 - \xi) - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right)^{-1} \left(\gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^{-1} - \zeta^{-1} \right) \\ &\leq \frac{b_{i_\ell^*}^2}{18} \left(c_0^2 (1 - \xi) - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right)^{-1} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^{-1} \end{aligned}$$

This shows the last two properties. Moreover, we have that $T_\ell(\frac{1}{2}) \leq \mathcal{O} \left(\left(\frac{d}{\log m^*} \right)^{\frac{1}{2}} \right)$. However, at $T_\ell(\frac{1}{2})$, by (16), we have that

$$\frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \geq c_1^2 \left(\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) - \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2 \right) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \geq 0$$

which implies that $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \geq -\frac{\beta_6}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}}$ for all $t \leq \min \{T_c, T_\ell(\xi)\}$ for all ξ such that $(1 - \xi^2)^{-1} \leq \mathcal{O}(1)$, as long as $T(\frac{1}{2}) \leq T_c$. \square

Upper bounding $\gamma_{i,j}^{(2)}(t)$ for $i \in [m] \setminus \mathcal{R}_{\ell-1}$ and $j \in [m^*] \setminus \mathcal{C}_{\ell-1}$ with $(i, j) \neq (i_\ell^*, j_\ell^*)$. Here we are going to show that the growth of $\gamma_{i,j}^{(2)}(t)$ with $i \in [m] \setminus \mathcal{R}_{\ell-1}, j \in [m^*] \setminus \mathcal{C}_{\ell-1}$ and $(i, j) \neq (i_\ell^*, j_\ell^*)$ is slow in terms of when $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)$ reaches ξ , $\gamma_{i,j}^{(2)}(t)$ is no bigger than $\mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$.

Lemma 9. *Suppose that the inductive hypothesis in Condition 2, and the initialization condition in Condition 1. Let T_c be defined in (14) and $T_\ell(\xi)$ in Definition 2. Then there exists some constant $\beta_7 > 0$ such that for all $t \leq \min \left\{ T_c, T \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} \right) + \frac{\beta_7 \delta_{\mathbb{P}} \sqrt{d}}{m^2} \right\}$, we have that $|\gamma_{i,j}^{(2)}(t)| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$ for all $(i, j) \neq (i_\ell^*, j_\ell^*)$.*

Proof. For $t \leq T_c$, by Lemma 7, we write out the dynamic of $\gamma_{i,j}^{(2)}(t)$ as

$$\begin{aligned} \frac{d}{dt} \gamma_{i,j}^{(2)}(t) &= -\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \mathbf{v}_j^* \\ &= \frac{18}{b_i^2} \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i,j}^{(1)}(t)^k \gamma_{i,j}^{(2)}(t)^2 - \frac{9}{b_i^2} \cdot \hat{\lambda}_{i_\ell^*, j_\ell^*} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \gamma_{i,j}^{(2)}(t) \mathbb{I}\{i = i_\ell^*\} \\ &\quad \pm \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \pm \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d} \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \mathbb{I}\{i = i_\ell^*\} \end{aligned}$$

As shown in Lemma 8, $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0, \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \geq 0$ for all $t \leq T_c$. Therefore, we have that

$$\begin{aligned} \frac{d}{dt} \left| \gamma_{i,j}^{(2)}(t) \right| &\leq \frac{18}{b_i^2} \left| \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i,j}^{(1)}(t)^k \right| \cdot \left| \gamma_{i,j}^{(2)}(t) \right|^2 - \frac{9}{b_i^2} \cdot \hat{\lambda}_{i_\ell^*, j_\ell^*, 5} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \left| \gamma_{i,j}^{(2)}(t) \right| \mathbb{I}\{i = i_\ell^*\} \\ &\quad + \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) + \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \\ &\leq \frac{18}{b_i^2} \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \left| \gamma_{i,j}^{(1)}(t) \right|^k \left| \gamma_{i,j}^{(2)}(t) \right|^2 + \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) + \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \\ &\leq \frac{18}{b_i^2} \left(c_0^2 + \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right) \right) \left| \gamma_{i,j}^{(2)}(t) \right|^2 + \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) + \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \end{aligned}$$

where the last inequality is because $\left| \gamma_{i,j}^{(1)}(t) \right| \leq \varepsilon_1(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right)$ for $t \leq T_c$. For any $t \leq T\left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}}\right)$, we must have that

$$\mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \leq \mathcal{O}\left(\frac{m^4}{d^{\frac{7}{5}}}\right)$$

Therefore, for any $t \leq T\left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}}\right)$, we have

$$\frac{d}{dt} \left| \gamma_{i,j}^{(2)}(t) \right| \leq \frac{18}{b_i^2} \left(c_0^2 + \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right) \right) \left| \gamma_{i,j}^{(2)}(t) \right|^2 + \mathcal{O}\left(\frac{m^4}{d^{\frac{7}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$$

Notice that $\left| \gamma_{i,j}^{(2)}(t) \right|$ must be upper bounded by $\hat{\gamma}_{i,j}^{(2)}(t)$ given by

$$\frac{d}{dt} \hat{\gamma}_{i,j}^{(2)}(t) = \frac{18}{b_i^2} \left(c_0^2 + \mathcal{O}\left(\frac{m^2}{\sqrt{d}}\right) \right) \hat{\gamma}_{i,j}^{(2)}(t)^2 + \mathcal{O}\left(\frac{m^4}{d^{\frac{7}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right); \quad \hat{\gamma}_{i,j}^{(2)}(0) = \left| \gamma_{i,j}^{(2)}(0) \right|$$

Observe that for any $t \leq T_c$, $\hat{\gamma}_{i,j}^{(2)}(t)$ grows monotonically as $\hat{\gamma}_{i,j}^{(2)}(0)$ grows. By the initialization property, we have that $\hat{\gamma}_{i,j}^{(2)}(0) = \left| \gamma_{i,j}^{(2)}(0) \right| \leq \frac{\sqrt{\beta_3}}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}}$. We can observe that $\hat{\gamma}_{i,j}^{(2)}(t)$ increases as $\hat{\gamma}_{i,j}^{(2)}(0)$ becomes larger. Thus, it suffice to consider $\hat{\gamma}_{i,j}^{(2)}(0) = \frac{\sqrt{\beta_3}}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}}$. In this case,

$$\frac{d}{dt} \hat{\gamma}_{i,j}^{(2)}(t) \leq \frac{18}{b_i^2} \left(c_0^2 + \mathcal{O}\left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}}\right) \right) \hat{\gamma}_{i,j}^{(2)}(t)^2$$

Solving the differential equation gives that

$$\left| \gamma_{i,j}^{(2)}(t) \right| \leq \hat{\gamma}_{i,j}^{(2)}(t) \leq \left(\left| \gamma_{i,j}^{(2)}(0) \right|^{-1} - \frac{18}{b_i^2} \left(c_0^2 + \mathcal{O}\left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}}\right) \right) t \right)^{-1}$$

Recall that $(1 + \delta_s)^2 b_i^2 \gamma_{i,j}^{(2)}(0)^2 \leq b_{i_\ell^*} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^2$ by the initialization property. Thus, $t \leq T \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} \right)$, it holds that

$$\begin{aligned}
 |\gamma_{i,j}^{(2)}(t)| &\leq \left(|\gamma_{i,j}^{(2)}(0)|^{-1} - \frac{18}{b_i^2} \left(c_0^2 + \mathcal{O} \left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right) T \left(\frac{m^{\frac{3}{2}}}{d^{\frac{1}{4}}} \right) \right)^{-1} \\
 &\leq \left(|\gamma_{i,j}^{(2)}(0)|^{-1} - \frac{b_{i_\ell^*}^2}{b_i^2} \left(c_0^2 + \mathcal{O} \left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right) \left(c_0^2 \left(1 - \frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} \right) - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}} \sqrt{d}} \right) \right)^{-1} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^{-1} \right)^{-1} \\
 &\leq \left(|\gamma_{i,j}^{(2)}(0)|^{-1} - \frac{b_{i_\ell^*}^2}{b_i^2} \cdot \frac{c_0^2 + \mathcal{O} \left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)}{c_0^2 - \mathcal{O} \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^{-1} \right)^{-1} \\
 &\leq \left(\frac{b_i}{b_{i_\ell^*}} (1 + \delta_s) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^{-1} - \frac{b_{i_\ell^*}^2}{b_i^2} \cdot \frac{c_0^2 + \mathcal{O} \left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)}{c_0^2 - \mathcal{O} \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0)^{-1} \right)^{-1} \\
 &\leq \frac{b_{i_\ell^*}}{b_i} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0) \left(1 + \delta_s - \frac{b_{i_\ell^*}^3}{b_i^3} \cdot \frac{c_0^2 + \mathcal{O} \left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)}{c_0^2 - \mathcal{O} \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)} \right)^{-1}
 \end{aligned}$$

By the initialization property, we have that $b_i, b_{i_\ell^*} \in [1 - \beta_2 \delta_s, 1 + \beta_2 \delta_s]$ for $\beta_2 \leq o(1)$. Therefore, we have that

$$\begin{aligned}
 \frac{b_{i_\ell^*}^3}{b_i^3} \cdot \frac{c_0^2 + \mathcal{O} \left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)}{c_0^2 - \mathcal{O} \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)} &\leq \left(\frac{1 + \beta_2 \delta_s}{1 - \beta_2 \delta_s} \right)^3 \frac{c_0^2 + \mathcal{O} \left(\frac{m^4}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)}{c_0^2 - \mathcal{O} \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} + \frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right)} \\
 &\leq \left(\frac{1 + \beta_2 \delta_s}{1 - \beta_2 \delta_s} \right)^4 \\
 &\leq 1 + \frac{1}{2} \delta_s
 \end{aligned}$$

where the second inequality is due to $d \geq \beta_5 m^{16}$ and the last inequality due to $\beta_2 \leq o(1)$. Therefore

$$|\gamma_{i,j}^{(2)}(t)| \leq \frac{b_{i_\ell^*}}{b_i} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(0) \cdot \frac{2}{\delta_s} \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$$

Let $T_1 = T \left(\frac{m^{\frac{3}{2}}}{d^{\frac{1}{4}}} \right)$. Then for $t \geq T_1$, the dynamic of $\gamma_{i,j}^{(2)}(t)$ is upper bounded by

$$\frac{d}{dt} \hat{\gamma}_{i,j}^{(2)}(t) \leq \frac{18}{b_i^2} \left(c_0^2 + \mathcal{O} \left(\frac{m^2}{\sqrt{d}} \right) \right) |\gamma_{i,j}^{(2)}(t)|^2 + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d} \right)$$

Let T_2 be the smallest $t \geq 0$ such that $\hat{\gamma}_{i,j}^{(2)}(t) \geq 2\hat{\gamma}_{i,j}^{(2)}(T_1)$. Then for all $t \leq T_2$, we have $\frac{d}{dt} \hat{\gamma}_{i,j}^{(2)}(t) \leq \mathcal{O} \left(\frac{m^4}{d} \right)$. Therefore

$$\hat{\gamma}_{i,j}^{(2)}(t + T_1) \leq \hat{\gamma}_{i,j}^{(2)}(T_1) + t \cdot \max_{t \leq T_2} \frac{d}{dt} \hat{\gamma}_{i,j}^{(2)}(t) \leq \hat{\gamma}_{i,j}^{(2)}(T_1) + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d} \right) t$$

Thus, we must have that $T_2 \geq \mathcal{O} \left(\frac{\delta_{\mathbb{P}} \sqrt{d}}{m^2} \right) + T_1$. Therefore, we can conclude that there exists some constant $\beta_7 > 0$ such that for all $t \leq T \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} \right) + \frac{\beta_7 \delta_{\mathbb{P}} \sqrt{d}}{m^2}$ and $t \leq T_c$, we have that

$$|\gamma_{i,j}^{(2)}(t)| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$$

for all $(i, j) \neq (i_\ell^*, j_\ell^*)$. □

Upper bounding $\gamma_{i,j}^{(2)}(t)$ for $i \in [m] \setminus \mathcal{R}_{\ell-1}, j \in \mathcal{C}_{\ell-1}$. In this section, we show that the alignment of $\bar{\mathbf{w}}_i$ with previously recovered components $\bar{\mathbf{w}}_j^*$ must be small.

Lemma 10. *Suppose that the inductive hypothesis in Condition 2 and the initialization condition in Condition 1 holds. Let T_c be defined in (14) and $T_\ell(\xi)$ in Definition 2. Then we have that for all $t \leq \min \left\{ T_c, T_\ell(\xi) + \mathcal{O} \left(\frac{\delta_{\mathbb{P}} \sqrt{d}}{m^2} \right), \mathcal{O} \left(\frac{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}}{m^5} \right) \right\}$ for any ξ such that $(1 - \xi)^{-1} \leq \mathcal{O}(1)$, it holds that*

$$\left| \gamma_{i,j_{\ell'}}^{(2)}(t) \right| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right); \forall i \in [m] \setminus \mathcal{R}_{\ell-1}, \ell' < \ell$$

Proof. By the inductive hypothesis, we have that $\gamma_{i_{\ell'}, j_{\ell'}}^{(2)}(t) \geq 1 - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right)$. Fix any $i \in [m] \setminus \mathcal{R}_{\ell-1}$. By Lemma 7, we have that

$$\begin{aligned} \frac{d}{dt} \gamma_{i,j_{\ell'}}^{(2)}(t) &= -\frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_{j_{\ell'}}^* \\ &= \frac{9}{b_i^2} \left(\hat{\lambda}_{i,j_{\ell'},5}(t) - \lambda_{i,i_{\ell'},5}(t) \gamma_{i_{\ell'},j_{\ell'}}^{(2)}(t) \right) - \frac{9}{b_i^2} \hat{\lambda}_{i_{\ell'},j_{\ell'},5}(t) \gamma_{i_{\ell'},j_{\ell'}}^{(2)}(t) \gamma_{i,j_{\ell'}}^{(2)}(t) \mathbb{I}\{i = i_{\ell'}^*\} \\ &\quad \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_{\ell'},j_{\ell'}}^{(2)}(t)^2 \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t) + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) \end{aligned}$$

Due to the same reasoning as in the previous lemma, we have that $\hat{\lambda}_{i_{\ell'},j_{\ell'},5}(t) \gamma_{i_{\ell'},j_{\ell'}}^{(2)}(t) \geq 0$. Therefore, for all $t \leq T_c$, we have that

$$\begin{aligned} \frac{d}{dt} \left| \gamma_{i,j_{\ell'}}^{(2)}(t) \right| &\leq \frac{9}{b_i^2} \left| \hat{\lambda}_{i,j_{\ell'},5}(t) - \lambda_{i,i_{\ell'},5}(t) \gamma_{i_{\ell'},j_{\ell'}}^{(2)}(t) \right| + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d} \right) \gamma_{i,j_{\ell'}}^{(2)}(t)^2 + \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \\ &\leq \frac{9c_0^2}{b_i^2} \left| \gamma_{i,j_{\ell'}}^{(2)}(t)^2 - I_{i,i_{\ell'}}^{(2)}(t)^2 \gamma_{i_{\ell'},j_{\ell'}}^{(2)}(t) \right| + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d} \right) \gamma_{i,j_{\ell'}}^{(2)}(t)^2 + \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \\ &\leq \frac{9c_0^2}{b_i^2} \left(1 - \gamma_{i_{\ell'},j_{\ell'}}^{(2)}(t) \right) \gamma_{i,j_{\ell'}}^{(2)}(t)^2 + \frac{9c_0^2}{b_i^2} \gamma_{i_{\ell'},j_{\ell'}}^{(2)}(t) \left| \gamma_{i,j_{\ell'}}^{(2)}(t)^2 - I_{i,i_{\ell'}}^{(2)}(t)^2 \right| \\ &\quad + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d} \right) \gamma_{i,j_{\ell'}}^{(2)}(t)^2 + \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \end{aligned}$$

Recall that $\gamma_{i,j_{\ell'}}^{(2)}(t) = \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_{j_{\ell'}}^*$ and $I_{i,i_{\ell'}} = \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_{i_{\ell'}}^*$. Therefore,

$$\begin{aligned} \left| \gamma_{i,j_{\ell'}}^{(2)}(t)^2 - I_{i,i_{\ell'}}^{(2)}(t)^2 \right| &= \left| \gamma_{i,j_{\ell'}}^{(2)}(t) + I_{i,i_{\ell'}}^{(2)}(t) \right| \cdot \left| \gamma_{i,j_{\ell'}}^{(2)}(t) - I_{i,i_{\ell'}}^{(2)}(t) \right| \\ &\leq 2 \left| \gamma_{i,j_{\ell'}}^{(2)}(t) \right| \cdot \left| \gamma_{i,j_{\ell'}}^{(2)}(t) - I_{i,i_{\ell'}}^{(2)}(t) \right| + \left(\gamma_{i,j_{\ell'}}^{(2)}(t) - I_{i,i_{\ell'}}^{(2)}(t) \right)^2 \\ &\leq 2 \left| \gamma_{i,j_{\ell'}}^{(2)}(t) \right| \cdot \left| \bar{\mathbf{v}}_i^\top (\bar{\mathbf{v}}_{j_{\ell'}}^* - \bar{\mathbf{v}}_{i_{\ell'}}^*) \right| + \left(\bar{\mathbf{v}}_i^\top (\bar{\mathbf{v}}_{j_{\ell'}}^* - \bar{\mathbf{v}}_{i_{\ell'}}^*) \right)^2 \\ &\leq 2 \left| \gamma_{i,j_{\ell'}}^{(2)}(t) \right| \cdot \left\| \bar{\mathbf{v}}_{j_{\ell'}}^* - \bar{\mathbf{v}}_{i_{\ell'}}^* \right\|_2 + \left\| \bar{\mathbf{v}}_{j_{\ell'}}^* - \bar{\mathbf{v}}_{i_{\ell'}}^* \right\|_2^2 \\ &\leq \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}} \right) \left| \gamma_{i,j_{\ell'}}^{(2)}(t) \right| + \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \end{aligned}$$

Moreover, since $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 1 - \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$, for $T_{p, \ell-1} \leq t \leq T_c$, we must have that

$$\begin{aligned} \frac{d}{dt} \left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| &\leq \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}}\right) \left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| + \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \\ &\leq \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}}\right) \left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| + \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \\ &\leq \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) + \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{4}}}\right) \end{aligned}$$

where the last inequality follows from Lemma 8. This gives that

$$\begin{aligned} \left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t + T_{p, \ell-1}) \right| &\leq \left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(T_{p, \ell-1}) \right| + \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) + \mathcal{O}\left(\frac{m^7 t}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \\ &\leq \left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(T_{p, \ell-1}) \right| + \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d} + \frac{m^7 t}{\delta_{\mathbb{P}}^3 d^{\frac{3}{4}}}\right) \\ &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} + \frac{m^7 t}{\delta_{\mathbb{P}}^3 d^{\frac{5}{4}}}\right) \end{aligned}$$

where the last step follows from the inductive hypothesis. Further requiring that $t \leq \mathcal{O}\left(\frac{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}}{m^5}\right)$ keeps $\left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t + T_{p, \ell-1}) \right| \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right)$. When $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq \xi$, we can still obtain that

$$\frac{d}{dt} \left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| \leq \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right)$$

Thus, for all $t \leq \min\left\{T_\ell(\xi) + \mathcal{O}\left(\frac{\delta_{\mathbb{P}} \sqrt{d}}{m^2}\right), \mathcal{O}\left(\frac{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}}{m^5}\right)\right\}$ we can guarantee that $\left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t + T_{p, \ell-1}) \right| \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right)$. \square

Bounding $\varepsilon_{1, \ell}(t)$, $\varepsilon_{2, \ell}(t)$, and $\varepsilon_{\mathcal{B}, \ell}^{(1)}(t)$. Here we are going to upper bound $\varepsilon_{1, \ell}(t)$, $\varepsilon_{2, \ell}(t)$ and $\varepsilon_{\mathcal{B}, \ell}^{(1)}(t)$. In particular, we are going to analyze $\gamma_{i, j}^{(1)}(t)$ for $i \in [m] \setminus \mathcal{R}_{\ell-1}$ and $j \in [m^*]$ where $(i, j) \neq i_\ell^*, j_\ell^*$, $\zeta_{i, j}^{(1)}, \zeta_{i, j}^{(2)}$ for $i \in [m] \setminus \mathcal{R}_{\ell-1}$ and $j \in [m^*]$, and also $I_{i, j}^{(1)}, I_{i, j}^{(2)}, I_{i, j}^{(3)}$ for $i, j \in [m]$ and $i \neq j$.

Lemma 11. *Suppose that the inductive hypothesis in Condition 2 and the initialization condition in Condition 1 holds. Let T_c be defined in (14) and $T(\xi)$ in Definition 2. Then there exists some constant $\beta_7 > 0$ such that for all $\beta_8 \leq \mathcal{O}(1)$, for all $t \leq \min\left\{T_c, T\left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}}\right) + \frac{\beta_7 \delta_{\mathbb{P}} \sqrt{d}}{m}, T(\xi)\right\} + \beta_8$, we have that $\varepsilon_{\mathcal{A}, \ell}^{(2)} \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right)$. Moreover, for all $t \leq \min\left\{\mathcal{O}\left(\frac{\delta_{\mathbb{P}}^2 d}{m^5}\right), T_c\right\}$ we shall have that*

$$\begin{aligned} \max\left\{\left|\gamma_{i, j}^{(1)}(t)\right|, \left|\zeta_{i, j}^{(1)}(t)\right|, \left|\zeta_{i, j}^{(2)}(t)\right|\right\} &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right); \forall i \in [m] \setminus \mathcal{R}_\ell, j \in [m^*] \\ \max\left\{\left|I_{i, j}^{(1)}(t)\right|, \left|I_{i, j}^{(2)}(t)\right|, \left|I_{i, j}^{(3)}(t)\right|\right\} &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right); \forall i, j \in [m] \setminus \mathcal{R}_\ell, i \neq j \end{aligned}$$

Proof. Throughout the proof, we will relax the upper bound in terms of $\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)$ into the upper bound in terms of $\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)$.

Bounding $\zeta_{i, j}^{(1)}(t), \zeta_{i, j}^{(2)}(t)$. First, we are going to derive some rough estimation for $\zeta_{i, j}^{(1)}(t), \zeta_{i, j}^{(2)}(t)$ for $i \in [m] \setminus$

$\mathcal{R}_{\ell-1}, j \in [m^*]$. By Lemma 7, we have that

$$\begin{aligned} \frac{d}{dt} \left| \zeta_{i,j}^{(1)}(t) \right| &= -\frac{1}{a_i} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j^* \cdot \text{sign} \left(\zeta_{i,j}^{(1)}(t) \right) \\ &\leq \frac{3}{a_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 4} \right| + \frac{1}{a_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| - \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| \mathbb{I}\{i = i_\ell^*\} \\ &\quad + \frac{36}{a_i^2} \left| C_{S,2} \gamma_{i,j}^{(2)}(t) I_{i,i}^{(3)} \right| + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2 + \varepsilon_{5,\ell}(t)^2 \right) \end{aligned}$$

Diving into the details of the first three terms, we have that

$$\begin{aligned} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) \right| &\leq 6 \left| I_{i_\ell^*, i_\ell^*}^{(3)} \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_{k+2} c_k}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^k \\ &\quad + 6 \left| \zeta_{i,j}^{(2)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^k + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{B},\ell}^{(1)}(t) + \varepsilon_{5,\ell}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| &\leq 6 \left| \zeta_{i,j}^{(1)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^k + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{B},\ell}^{(1)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ -\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \left| \zeta_{i,j}^{(1)}(t) \right| &\leq 6 \left| \zeta_{i,j}^{(1)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^k \\ &\quad + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 + \varepsilon_{\mathcal{B},\ell}^{(1)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

Also, we notice that $\left| \gamma_{i,j}^{(2)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{5,\ell}(t) + \varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right)$. Applying the fact that $a_i^{-1} \leq \mathcal{O}(1)$, we have that

$$\begin{aligned} \frac{d}{dt} \left| \zeta_{i,j}^{(1)}(t) \right| &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 + \varepsilon_{\mathcal{B},\ell}^{(1)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(3)}(t) + \varepsilon_{5,\ell}(t) \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 + \varepsilon_{\mathcal{B},\ell}^{(1)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

Similarly, for $\zeta_{i,j}^{(2)}(t)$, by Lemma 7, we have that

$$\begin{aligned} \frac{d}{dt} \left| \zeta_{i,j}^{(2)}(t) \right| &= -\frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j^* \cdot \text{sign} \left(\zeta_{i,j}^{(2)}(t) \right) \\ &\leq \frac{3}{b_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \right| + \frac{3}{b_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| - 3 \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \left| \zeta_{i,j}^{(2)}(t) \right| \mathbb{I}\{i = i_\ell^*\} \\ &\quad + \frac{36}{b_i^2} \left| C_{S,2} \gamma_{i,j}^{(1)}(t) I_{i,i}^{(3)}(t) \right| + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{B},\ell}^{(3)}(t)^2 + \varepsilon_{5,\ell}(t)^2 \right) \end{aligned}$$

The third term is apparently negative due to the fact that $\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \geq 0$, $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0$. For the first two, we have that

$$\begin{aligned} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \right| &\leq 6 \left| \zeta_{i,j}^{(1)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^k + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{B},\ell}^{(1)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| &\leq 6 \left| \zeta_{i,j}^{(1)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_{k+2} c_k}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^{k+1} + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{B},\ell}^{(1)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

As before, we have that $\left| I_{i,i}^{(3)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{B},\ell}^{(3)}(t) + \varepsilon_{5,\ell}(t) \right)$. Therefore, we have that for $t \leq T_c$,

$$\frac{d}{dt} \left| \zeta_{i,j}^{(2)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{B},\ell}^{(1)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right)$$

Bounding $I_{i,j}^{(1)}(t)$, $I_{i,j}^{(2)}(t)$, and $I_{i,j}^{(3)}(t)$ for $i, j \in [m] \setminus \mathcal{R}_{\ell-1}$, and $i \neq j$. Next, we will look at $I_{i,j}^{(1)}$. By Lemma 7, we have that

$$\begin{aligned} \frac{d}{dt} \left| I_{i,j}^{(1)}(t) \right| &= -\text{sign} \left(I_{i,j}^{(1)}(t) \right) \left(\frac{1}{a_i} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j + \frac{1}{a_j} \nabla_{\mathbf{v}_j} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_i \right) \\ &\leq \frac{1}{a_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \right| \left(\left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| + \left| \gamma_{j_\ell^*, i_\ell^*}^{(1)}(t) \right| \right) - \frac{1}{a_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \left| I_{i,j}^{(1)}(t) \right| \mathbb{I} \{ i = i_\ell^* \vee j = i_\ell^* \} \\ &\quad + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{5,\ell}(t) + \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) \end{aligned}$$

Since $i \neq j$, for the first term we have that

$$\begin{aligned} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \right| \left(\left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| + \left| \gamma_{j_\ell^*, i_\ell^*}^{(1)}(t) \right| \right) &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \right| \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

For the second term, we have that

$$\begin{aligned} -\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \left| I_{i,j}^{(1)}(t) \right| &\leq 6 \left| I_{i,j}^{(1)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \sum_{k=0}^{\infty} \frac{C_{k+1}^2}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^k \\ &\quad + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

Therefore, we have that

$$\frac{d}{dt} \left| I_{i,j}^{(1)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right)$$

For $I_{i,j}^{(2)}(t)$, we have that

$$\begin{aligned} \frac{d}{dt} \left| I_{i,j}^{(2)}(t) \right| &= -\text{sign} \left(I_{i,j}^{(2)}(t) \right) \left(\frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j + \frac{1}{b_j} \nabla_{\mathbf{w}_j} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_i \right) \\ &\leq \frac{9}{b_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \right| \left(\left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| + \left| \gamma_{j_\ell^*, i_\ell^*}^{(2)}(t) \right| \right) - \frac{9}{b_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \left| I_{i,j}^{(2)}(t) \right| \mathbb{I} \{ i = i_\ell^* \vee j = i_\ell^* \} \\ &\quad + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \frac{9}{b_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \right| \left(\left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| + \left| \gamma_{j_\ell^*, i_\ell^*}^{(2)}(t) \right| \right) + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

since $\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0$. For the first term, we have that

$$\left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \right| \left(\left| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| + \left| \gamma_{j_\ell^*, i_\ell^*}^{(2)}(t) \right| \right) \leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right)$$

Therefore, we have that

$$\frac{d}{dt} \left| I_{i,j}^{(2)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right)$$

For $I_{i,j}^{(3)}$, by Lemma 7, we also have that

$$\begin{aligned} \frac{d}{dt} \left| I_{i,j}^{(3)}(t) \right| &= -\text{sign} \left(\frac{1}{a_i} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j + \frac{1}{b_j} \nabla_{\mathbf{w}_j} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_i \right) \\ &\leq \frac{1}{a_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \zeta_{j, j_\ell^*}^{(2)}(t) \right| - \frac{1}{a_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \left| I_{i,j}^{(3)}(t) \right| \mathbb{I}\{i = i_\ell^*\} \\ &\quad + \frac{9}{b_j^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \zeta_{i, i_\ell^*}^{(1)}(t) \right| - \frac{9}{b_j^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \left| I_{i,j}^{(3)}(t) \right| \mathbb{I}\{j = i_\ell^*\} \\ &\quad + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right) \end{aligned}$$

where for the first, second, and third term, we have

$$\begin{aligned} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \zeta_{j, j_\ell^*}^{(2)}(t) \right| &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right) \\ \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \zeta_{j, j_\ell^*}^{(1)}(t) \right| &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right) \\ -\hat{\lambda}_{i_\ell^*, j_\ell^*, 1} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \left| I_{i,j}^{(3)}(t) \right| &\leq 6 \left| I_{i,j}^{(3)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^k \\ &\quad + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right) \end{aligned}$$

Noticing that $\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t), \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0$ for all $t \leq T_c$, we have that

$$\frac{d}{dt} \left| I_{i,j}^{(3)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right)$$

Bounding $\gamma_{i,j}^{(1)}(t)$. Finally, for $\gamma_{i,j}^{(1)}(t)$, in the case of $(i, j) \neq (i_\ell^*, j_\ell^*)$, we have that

$$\begin{aligned} \frac{d}{dt} \left| \gamma_{i,j}^{(1)}(t) \right| &= -\frac{1}{a_i} \text{sign} \left(\nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j^* \right) \\ &\leq -\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \left| \gamma_{i,j}^{(1)}(t) \right| \mathbb{I}\{i = i_\ell^*\} + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right) \end{aligned}$$

where since $(i, j) \neq (i_\ell^*, j_\ell^*)$, we must have that

$$-\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \left| \gamma_{i,j}^{(1)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right)$$

Therefore

$$\frac{d}{dt} \left| \gamma_{i,j}^{(1)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right)$$

Bounding $I_{i,j}^{(1)}(t), I_{i,j}^{(2)}(t)$, and $I_{i,j}^{(3)}(t)$ for $i, j \in [m] \setminus \mathcal{R}_{\ell-1}$ for $i \in \mathcal{R}_{\ell-1}$ or $j \in \mathcal{R}_{\ell-1}$. Recall the definition of $I_{i,j}^{(1)}(t), I_{i,j}^{(2)}(t)$, and $I_{i,j}^{(3)}(t)$

$$I_{i,j}^{(1)}(t) = \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{v}}_j(t); I_{i,j}^{(2)}(t) = \bar{\mathbf{w}}_i(t)^\top \bar{\mathbf{w}}_j(t); I_{i,j}^{(3)}(t) = \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{w}}_j(t)$$

For the sake of convenience, we let $i \in \mathcal{R}_{\ell-1}$ and study $I_{i,j}^{(1)}(t), I_{i,j}^{(2)}(t)$, and $I_{i,j}^{(3)}(t), I_{j,i}^{(3)}(t)$. Since $i \in \mathcal{R}_{\ell-1}$, there exists $\ell' < \ell$ such that $i = i_{\ell'}^*$. Therefore, for $I_{i,j}^{(1)}(t)$, when $t \geq T_{p, \ell-1}$, we have that

$$\begin{aligned} \left| I_{i,j}^{(1)}(t) \right| &= \left| \bar{\mathbf{v}}_j(t)^\top \bar{\mathbf{v}}_{j_{\ell'}}^* + \bar{\mathbf{v}}_j(t)^\top \left(\bar{\mathbf{v}}_{i_{\ell'}^*}(t) - \bar{\mathbf{v}}_{j_{\ell'}^*}^* \right) \right| \\ &\leq \left| \gamma_{j, j_{\ell'}^*}^{(1)}(t) \right| + \left\| \bar{\mathbf{v}}_{i_{\ell'}^*}(t) - \bar{\mathbf{v}}_{j_{\ell'}^*}^* \right\|_2 \\ &\leq \left| \gamma_{j, j_{\ell'}^*}^{(1)}(t) \right| + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}} \right) \end{aligned} \tag{17}$$

for all $t \geq T_{p,\ell-1}$. Similarly, for $I_{i,j}^{(2)}(t)$, when $t \geq T_{p,\ell-1}$, we have that

$$\begin{aligned} \left| I_{i,j}^{(2)}(t) \right| &= \left| \bar{\mathbf{w}}_j(t)^\top \bar{\mathbf{w}}_{j_{\ell'}^*}^* + \bar{\mathbf{w}}_j(t)^\top \left(\bar{\mathbf{w}}_{i_{\ell'}^*}(t) - \bar{\mathbf{w}}_{j_{\ell'}^*}^* \right) \right| \\ &\leq \left| \gamma_{j,j_{\ell'}^*}^{(2)}(t) \right| + \left\| \bar{\mathbf{w}}_{i_{\ell'}^*}(t) - \bar{\mathbf{w}}_{j_{\ell'}^*}^* \right\|_2 \\ &\leq \left| \gamma_{j,j_{\ell'}^*}^{(2)}(t) \right| + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}} \right) \end{aligned} \quad (18)$$

For $I_{i,j}^{(3)}(t)$ and $I_{j,i}^{(3)}(t)$, when $t \geq T_{p,\ell-1}$, we have that

$$\begin{aligned} \left| I_{i,j}^{(3)}(t) \right| &= \left| \bar{\mathbf{w}}_j(t)^\top \bar{\mathbf{v}}_{j_{\ell'}^*}^* + \bar{\mathbf{w}}_j(t)^\top \left(\bar{\mathbf{v}}_{i_{\ell'}^*}(t) - \bar{\mathbf{v}}_{j_{\ell'}^*}^* \right) \right| \\ &\leq \left| \zeta_{j,j_{\ell'}^*}^{(2)}(t) \right| + \left\| \bar{\mathbf{v}}_{i_{\ell'}^*}(t) - \bar{\mathbf{v}}_{j_{\ell'}^*}^* \right\|_2 \\ &\leq \left| \zeta_{j,j_{\ell'}^*}^{(2)}(t) \right| + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}} \right) \\ \left| I_{j,i}^{(3)}(t) \right| &= \left| \bar{\mathbf{v}}_i(t)^\top \bar{\mathbf{w}}_{j_{\ell'}^*}^* + \bar{\mathbf{v}}_i(t)^\top \left(\bar{\mathbf{w}}_{i_{\ell'}^*}(t) - \bar{\mathbf{w}}_{j_{\ell'}^*}^* \right) \right| \\ &\leq \left| \zeta_{i,j_{\ell'}^*}^{(1)}(t) \right| + \left\| \bar{\mathbf{w}}_{i_{\ell'}^*}(t) - \bar{\mathbf{w}}_{j_{\ell'}^*}^* \right\|_2 \\ &\leq \left| \zeta_{i,j_{\ell'}^*}^{(1)}(t) \right| + \mathcal{O} \left(\frac{m^4}{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}} \right) \end{aligned} \quad (19)$$

Moreover, by the inductive hypothesis, we have that for $t \leq T_{p,\ell-1}$

$$\left| I_{i,j}^{(1)}(t) \right|, \left| I_{i,j}^{(2)}(t) \right|, \left| I_{i,j}^{(3)}(t) \right|, \left| I_{j,i}^{(3)}(t) \right| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$$

Gathering the results. Defining $\hat{\varepsilon}_\ell(t)$ to be the maximum of $\left| \zeta_{i,j}^{(1)}(t) \right|$, $\left| \zeta_{i,j}^{(2)}(t) \right|$, $\left| \gamma_{i,j}^{(1)}(t) \right|$ for i, j bounded above excluding $(i, j) = (i_{\ell'}^*, j_{\ell'}^*)$ for $\gamma_{i,j}^{(1)}(t)$, and $\left| I_{i,j}^{(1)}(t) \right|$, $\left| I_{i,j}^{(2)}(t) \right|$, $\left| I_{i,j}^{(3)}(t) \right|$ for i, j bounded above excluding $i \in \mathcal{R}_{\ell-1}$ or $j \in \mathcal{R}_{\ell-1}$. Gathering the results, we have that

$$\frac{d}{dt} \hat{\varepsilon}_\ell(t) \leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \gamma_{i_{\ell'}^*, j_{\ell'}^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right)$$

for $t \leq T_c$. In the meantime, by Lemma 9 we have that

$$\left| \gamma_{i,j}^{(2)}(t) \right| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right); \forall t \leq \min \left\{ T_c, T \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} \right) + \frac{\beta_7 \delta_{\mathbb{P}} \sqrt{d}}{m^2} \right\}$$

where $i \in [m] \setminus \mathcal{R}_{\ell-1}$, $j \in [m^*] \setminus \mathcal{C}_{\ell-1}$ and $(i, j) \neq (i_{\ell'}^*, j_{\ell'}^*)$. Therefore, we have that

$$\begin{aligned} \frac{d}{dt} \hat{\varepsilon}_\ell(t) &\leq \mathcal{O} \left(\hat{\varepsilon}_\ell(t) + \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right) \right) \gamma_{i_{\ell'}^*, j_{\ell'}^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\hat{\varepsilon}_\ell(t) + \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right) \right) \frac{d}{dt} \gamma_{i_{\ell'}^*, j_{\ell'}^*}^{(2)}(t) + \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\ &\leq \mathcal{O} \left(\hat{\varepsilon}_\ell(t) + \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right) \right) \frac{d}{dt} \gamma_{i_{\ell'}^*, j_{\ell'}^*}^{(2)}(t) + \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \end{aligned}$$

where the last inequality follows from Lemma 8. Solving the differential equation gives

$$\hat{\varepsilon}_\ell(t) \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right) \hat{\varepsilon}_\ell(0) \exp \left(\mathcal{O} \left(\gamma_{i_{\ell'}^*, j_{\ell'}^*}^{(2)}(t) \right) \right) + \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) \exp \left(\mathcal{O} \left(\gamma_{i_{\ell'}^*, j_{\ell'}^*}^{(2)}(t) \right) \right) \cdot t$$

Imposing $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \leq 1$ gives the desired result.

Fine-Grained Result for $i \in [m] \setminus \mathcal{R}_\ell$. By Lemma 7, we have that for all $i \in [m] \setminus \mathcal{R}_\ell$ and $j \in [m^*]$, it holds that

$$\max \left\{ \left| \gamma_{i,j}^{(1)}(t) \right|, \left| \zeta_{i,j}^{(1)}(t) \right|, \left| \zeta_{i,j}^{(2)}(t) \right| \right\} \leq \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) \leq \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^2 d^{\frac{3}{2}}} \right)$$

For $i, j \in [m] \setminus \mathcal{R}_\ell$ with $i \neq j$, we have that

$$\max \left\{ \left| I_{i,j}^{(1)}(t) \right|, \left| I_{i,j}^{(2)}(t) \right|, \left| I_{i,j}^{(3)}(t) \right| \right\} \leq \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(2)}(t)\varepsilon_{\mathcal{B},\ell}^{(3)}(t) \right) \leq \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right)$$

Thus, for all $t \leq \mathcal{O} \left(\frac{\delta_{\mathbb{P}}^2 d}{m^5} \right)$ we shall have that

$$\begin{aligned} \max \left\{ \left| \gamma_{i,j}^{(1)}(t) \right|, \left| \zeta_{i,j}^{(1)}(t) \right|, \left| \zeta_{i,j}^{(2)}(t) \right| \right\} &\leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right); \forall i \in [m] \setminus \mathcal{R}_\ell, j \in [m^*] \\ \max \left\{ \left| I_{i,j}^{(1)}(t) \right|, \left| I_{i,j}^{(2)}(t) \right|, \left| I_{i,j}^{(3)}(t) \right| \right\} &\leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right); \forall i, j \in [m] \setminus \mathcal{R}_\ell, i \neq j \end{aligned}$$

□

Bounding $I_{i,i}^{(3)}(t)$. Here we are going to upper bound $I_{i,i}^{(3)}(t)$ for $i \in [m] \setminus \mathcal{R}_{\ell-1}$

Lemma 12. *Suppose that the inductive hypothesis in Condition 2, and the initialization condition in Condition 1. Let T_c be defined in (14) and $T(\xi)$ in Definition 2. Then for all $t \leq T_c$, we have that*

$$\left| I_{i,i}^{(3)}(t) \right| \leq \begin{cases} \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) & \text{if } i \neq i_\ell^* \\ \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(2)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \right) & \text{if } i = i_\ell^* \end{cases}$$

Proof. By Lemma 7, we have that, in the case where $i \neq i_\ell^*$

$$\begin{aligned} \frac{d}{dt} \left| I_{i,i}^{(3)}(t) \right| &= -\text{sign} \left(I_{i,i}^{(3)}(t) \right) \left(\frac{1}{a_i} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j + \frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j \right) \\ &\leq -36C_{S,2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) \left| I_{i,i}^{(3)}(t) \right| + \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

Since $\varepsilon_{\mathcal{A},\ell}^{(2)}(t)$ is monotonic non-decreasing, we must have that $\left| I_{i,i}^{(3)}(t) \right| \leq \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right)$ for all $t \geq 0$. Similarly, in the case $i = i_\ell^*$, we have that

$$\begin{aligned} \frac{d}{dt} \left| I_{i,i}^{(3)}(t) \right| &= -\text{sign} \left(I_{i,i}^{(3)}(t) \right) \left(\frac{1}{a_i} \nabla_{\mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{w}}_j + \frac{1}{b_i} \nabla_{\mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}(t))^\top \bar{\mathbf{v}}_j \right) \\ &= \frac{1}{a_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| + \frac{9}{b_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| + 3 \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \right| \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) \\ &\quad + \frac{3}{a_i^2} \text{sign} \left(I_{i,i}^{(3)}(t) \right) \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) + \frac{3}{b_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| \right| \\ &\quad - 36C_{S,2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) \left| I_{i,i}^{(3)}(t) \right| - \frac{1}{a_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| \left| I_{i,i}^{(3)}(t) \right| \\ &\quad - \frac{9}{b_i^2} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right| \left| I_{i,i}^{(3)}(t) \right| + \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

As in the previous proof, we have that

$$\begin{aligned} \left| \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| &\leq 6 \left| \zeta_{i,j}^{(1)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{C_k C_{k+2}}{k!} \left| \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|^{k+1} + \mathcal{O} \left(\varepsilon_1(t)^3 \right) \\ &\leq \mathcal{O} \left(\varepsilon_{\mathcal{B},\ell}^{(1)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \end{aligned}$$

The second and third term requires a more careful analysis. First, by the definition of $\hat{\lambda}_{i_\ell^*, j_\ell^*, 4}$, we have that

$$\begin{aligned}
 & \hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \cdot \text{sign} \left(I_{i,i}^{(3)}(t) \right) \\
 & \leq 18 \left| I_{i,i}^{(3)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \sum_{k=0}^{\infty} \frac{c_{k+2} c_k}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k + 18 \left| \zeta_{i,j}^{(2)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \right|^3 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \\
 & \leq 18 \left| I_{i,i}^{(3)}(t) \right| \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \sum_{k=0}^{\infty} \frac{c_{k+2} c_k}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \\
 & \leq \mathcal{O} \left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right) \left| I_{i,i}^{(3)}(t) \right| + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3
 \end{aligned}$$

Moreover, by the definition of $\hat{\lambda}_{i_\ell^*, j_\ell^*, 1}$, we have that

$$\begin{aligned}
 -\lambda_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) & = -6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^{k+1} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \\
 & = -6 \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^3 \left(c_1^2 \pm \mathcal{O} \left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right) \right) \\
 & \leq \mathcal{O} \left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right)
 \end{aligned}$$

Recalling that $\lambda_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)} \geq 0$ for all t , we have that

$$\begin{aligned}
 \frac{d}{dt} \left| I_{i,i}^{(3)}(t) \right| & \leq - \left(36 C_{S,2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) - \mathcal{O} \left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}} \right)^{\frac{1}{2}} \right) \right) \left| I_{i,i}^{(3)}(t) \right| \\
 & \quad + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right) \\
 & \leq -35 C_{S,2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) \left| I_{i,i}^{(3)}(t) \right| + \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right)
 \end{aligned}$$

Since $\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)$ is monotonically non-decreasing, we must have that $I_{i,i}^{(3)} \leq \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right)$ for all $t \geq 0$. \square

With the above preparation work, we are ready to derive the result for phase 1 convergence.

Lemma 13. *Suppose that the inductive hypothesis in Condition 2 and the initialization condition in Condition 1 holds. Let T_c be defined in (14) and $T(\xi)$ be defined in Definition 2. Then there exists constant $\beta_7 \geq 0$ such that for all ξ satisfying $(1 - \xi^2)^{-1} \leq \mathcal{O}(1)$ and all constant $\beta_8 \geq 0$, it holds that $T_c \geq \min \left\{ T \left(\frac{\delta_{\mathbb{P}}^2}{d^{\frac{2}{5}}} \right) + \frac{\beta_7 \delta_{\mathbb{P}} \sqrt{d}}{m^2}, T(\xi) \right\} + \beta_8$, and there exists $T_1 \leq \min \left\{ T_c, \mathcal{O}(\sqrt{d}) \right\}$ such that for all $T_1 \leq t \leq T_c$ we have that $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0.9$.*

Proof. By Lemma 10, Lemma 11 and Lemma 12, we have that $T_c \geq \min \left\{ T \left(\frac{\delta_{\mathbb{P}}}{d^{\frac{2}{5}}} \right) + \frac{\beta_7 \delta_{\mathbb{P}} \sqrt{d}}{m}, T(\xi) \right\} + \beta_8$. Thus, it remains to show that there exists some $T_1 \leq T_c$ such that $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0.9$. To start, we choose $\xi = \frac{1}{2}$. By Lemma 8 we have that

$$\begin{aligned}
 \frac{1}{2} & = \gamma_{i_\ell^*, j_\ell^*}^{(2)} \left(T \left(\frac{1}{2} \right) \right) \\
 & \geq \left(\gamma_{i_\ell^*, j_\ell^*}^{(2)} \left(T \left(\frac{\delta_{\mathbb{P}}}{d^{\frac{2}{5}}} \right) \right) \right)^{-1} - \frac{18}{b_{i_\ell^*}} \left(c_0^2 \left(1 - \frac{1}{4} \right) - \mathcal{O} \left(\frac{m^7}{\delta_{\mathbb{P}}^3 \sqrt{d}} \right) \right) \left(T \left(\frac{1}{2} \right) - T \left(\frac{\delta_{\mathbb{P}}}{d^{\frac{2}{5}}} \right) \right)^{-1}
 \end{aligned}$$

This gives that

$$T\left(\frac{1}{2}\right) \leq \mathcal{O}\left(\frac{d^{\frac{2}{5}}}{\delta_{\mathbb{P}}^2}\right) + T\left(\frac{\delta_{\mathbb{P}}}{d^{\frac{2}{5}}}\right) \leq T\left(\frac{\delta_{\mathbb{P}}}{d^{\frac{2}{5}}}\right) + \frac{\beta_7 \delta_{\mathbb{P}} \sqrt{d}}{m}$$

when $d \geq \delta_{\mathbb{P}}^{10} m^5$. Thus, it suffice to show that a $T_1 \leq T\left(\frac{1}{2}\right) + \beta_8$ achieves $\gamma_{i_{\ell}^*, j_{\ell}^*}(1), \gamma_{i_{\ell}^*, j_{\ell}^*}(2) \geq 0.9$ for some constant β_8 . Notice that this choice of T_1 satisfies $T_1 \leq \mathcal{O}(\sqrt{d})$. That is, if we define

$$T_{1,1} = \min\left\{t \geq 0 : \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t) \geq 0.9\right\}; T_{1,2} = \min\left\{t \geq 0 : \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) \geq 0.9\right\}$$

Then it suffice to show that $T_{1,1}, T_{1,2} \leq T\left(\frac{1}{2}\right) + \beta_8$. By Lemma 7, we have that for $T\left(\frac{1}{2}\right) \leq t \leq \min\{T_{1,2}, T\left(\frac{1}{2}\right) + \beta_8\}$, it holds that

$$\begin{aligned} \frac{d}{dt} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) &\geq \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2 \left(1 - \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2\right) \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t)^k - \mathcal{O}\left(\varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) \\ &\geq \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2 \left(1 - \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2\right) \left(c_0^2 - \mathcal{O}\left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}\right)\right) - \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right) \\ &\geq \frac{c_0^2}{2} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2 \left(1 - \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2\right) - \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right) \end{aligned}$$

This shows that $\frac{d}{dt} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) \geq 0$ as long as $1 - \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2 \geq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right)$. Thus, we can conclude that $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) \geq \frac{1}{2}$ for all $T\left(\frac{1}{2}\right) \leq t \leq \min\{T_{1,2}, T\left(\frac{1}{2}\right) + \beta_8\}$. For the dynamic of $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)$, this implies that for $T\left(\frac{1}{2}\right) \leq t \leq \min\{T_{1,2}, T\left(\frac{1}{2}\right) + \beta_8\}$

$$\frac{d}{dt} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) \geq \frac{c_0^2}{42} - \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right) \geq \frac{c_0^2}{50}$$

Thus $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}\left(t + T\left(\frac{1}{2}\right)\right) \geq \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}\left(T\left(\frac{1}{2}\right)\right) + \frac{c_0^2}{50} \cdot t$, which implies that $T_{1,2} \leq T\left(\frac{1}{2}\right) + \frac{20}{c_0^2} \leq T\left(\frac{1}{2}\right) + \beta_8$ for some $\beta_8 \geq 0$. Next, for the dynamic of $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)$, we have that

$$\begin{aligned} \frac{d}{dt} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t) &\geq 6 \left(1 - \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t)^2\right) \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^3 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t)^k - \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right) \\ &\geq \frac{3}{4} \left(1 - \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t)^2\right) \left(c_1^2 - \mathcal{O}\left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}\right)\right) - \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right) \\ &\geq \frac{c_1^2}{8} \end{aligned}$$

Thus $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}\left(t + T\left(\frac{1}{2}\right)\right) \geq \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}\left(T\left(\frac{1}{2}\right)\right) + \frac{c_1^2}{8} \cdot t$, which implies that $T_{1,1} \leq T\left(\frac{1}{2}\right) + \frac{16}{c_1^2}$ since $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}\left(T\left(\frac{1}{2}\right)\right) \geq -\mathcal{O}\left(\frac{1}{\sqrt{d}} \left(\log \frac{m}{\delta_{\mathbb{P}}}\right)^{\frac{1}{2}}\right)$ by Lemma 8. Therefore, we can conclude that $T_{1,1}, T_{1,2} \leq T\left(\frac{1}{2}\right) + \beta_8$ for some constant β_8 .

Now, we consider the dynamic of $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t)$ and $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)$ for $t \geq \max\{T_{1,1}, T_{1,2}\}$. As before, we have that

$$\begin{aligned} \frac{d}{dt} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) &\geq \frac{c_1^2}{2} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^3 \left(1 - \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t)^2\right) - \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right) \\ \frac{d}{dt} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) &\geq \frac{c_0^2}{2} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2 \left(1 - \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t)^2\right) - \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right) \end{aligned}$$

We can observe that $\frac{d}{dt} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t), \frac{d}{dt} \gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t)$ if

$$\gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t), \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) \leq 1 - \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}}\right)$$

Thus, for all $T_1 \leq t \leq T_c$ we have that $\gamma_{i_{\ell}^*, j_{\ell}^*}^{(1)}(t), \gamma_{i_{\ell}^*, j_{\ell}^*}^{(2)}(t) \geq 0.9$ □

A.6 Establishing the Inductive Hypothesis: Phase 2

A.7 Phase 2: Growth to Near-Perfect Alignment.

In this section, our goal is to show complete the inductive hypothesis.

The decay of $\left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|$, $\left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|$ **and** $\left| I_{i_\ell^*, j_\ell^*}^{(3)}(t) \right|$.

Lemma 14. *Suppose that the inductive hypothesis in Condition 2, the initialization condition in Condition 1, and the condition on the sigmoid function in Assumption 3 holds. Then for $t \geq T_1$, where T_1 is defined in Lemma 13, we have that*

$$\max \left\{ \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|, \left| I_{i_\ell^*, j_\ell^*}^{(3)}(t) \right| \right\} \leq \begin{cases} \mathcal{O} \left(m\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 + \frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right) & \text{for } T_1 \leq t \leq T_c \\ \mathcal{O} \left(m\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right) & \text{for } T_1 + \mathcal{O}(\log d) \leq t \leq T_c \end{cases}$$

Proof. To ease the analysis in this section, we are going to define

$$\begin{aligned} Q_0(t) &= \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k; \\ Q_1(t) &= \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k; \\ Q_2(t) &= \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \end{aligned}$$

By Lemma 5, we have that

$$\begin{aligned} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) &= 6\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)Q_1(t) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3 \right) \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) &= 6\zeta_{i_\ell^*, j_\ell^*}^{(1)}(t)Q_2(t) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3 \right) \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) &= 6\zeta_{i_\ell^*, j_\ell^*}^{(1)}(t)Q_1(t) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3 \right) \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) &= 6I_{i_\ell^*, j_\ell^*}^{(3)}(t)Q_2(t) + 6\zeta_{i_\ell^*, j_\ell^*}^{(2)}(t)Q_1(t) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3 \right) \\ \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) &= 2Q_0(t) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3 \right) \end{aligned}$$

We start with a detailed analysis of $I_{i_\ell^*, j_\ell^*}^{(3)}(t)$. Recall that

$$\begin{aligned} \frac{d}{dt} I_{i_\ell^*, j_\ell^*}^{(3)}(t) &= \frac{1}{a_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) + \frac{9}{b_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) + 3\hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) \\ &\quad + \frac{3}{a_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 4}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) + \frac{3}{b_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \\ &\quad - 36C_{S,2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) I_{i_\ell^*, i_\ell^*}^{(3)}(t) - \frac{1}{a_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) I_{i_\ell^*, i_\ell^*}^{(3)}(t) \\ &\quad - \frac{9}{b_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i_\ell^*, i_\ell^*}^{(3)}(t) \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \\ &= \frac{6}{a_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) + \frac{18}{b_i^2} Q_0(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) + 18Q_2(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) \\ &\quad + \frac{18}{a_i^2} Q_2(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i_\ell^*, j_\ell^*}^{(3)}(t) + \frac{18}{a_i^2} Q_1(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) + \frac{18}{b_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_1(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \\ &\quad - 36C_{S,2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) I_{i_\ell^*, i_\ell^*}^{(3)}(t) - \frac{6}{a_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) I_{i_\ell^*, i_\ell^*}^{(3)}(t) \\ &\quad - \frac{18}{b_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_0(t) I_{i_\ell^*, i_\ell^*}^{(3)}(t) \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \end{aligned}$$

Next, we will look into $\zeta_{i_\ell^*, j_\ell^*}^{(1)}(t), \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t)$. In particular, for $\zeta_{i_\ell^*, j_\ell^*}^{(1)}(t)$, we have that

$$\begin{aligned} \frac{d}{dt} \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) &= \frac{3}{a_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 4} + \frac{1}{a_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) - \frac{1}{a_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 1}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \\ &\quad - \frac{36}{a_i^2} C_{S,2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i_\ell^*, j_\ell^*}^{(3)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) \\ &= \frac{18}{a_i^2} Q_2(t) I_{i_\ell^*, j_\ell^*}^{(3)}(t) + \frac{18}{a_i^2} Q_1(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) + \frac{6}{a_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_2(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \\ &\quad - \frac{6}{a_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) - \frac{36}{a_i^2} C_{S,2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i_\ell^*, j_\ell^*}^{(3)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) \end{aligned}$$

Similarly, for $\zeta_{i_\ell^*, j_\ell^*}^{(2)}(t)$, we have that

$$\begin{aligned} \frac{d}{dt} \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) &= \frac{3}{b_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) + \frac{3}{b_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) - \frac{9}{b_i^2} \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \\ &\quad - \frac{36}{b_i^2} C_{S,2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i_\ell^*, j_\ell^*}^{(3)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) \\ &= \frac{18}{b_i^2} Q_1(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) + \frac{18}{b_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_2(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) - \frac{18}{b_i^2} Q_0(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \\ &\quad - \frac{36}{b_i^2} C_{S,2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) I_{i_\ell^*, j_\ell^*}^{(3)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) \end{aligned}$$

Then we can write $\zeta_{i_\ell^*, j_\ell^*}^{(1)}(t), \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t)$, and $I_{i_\ell^*, j_\ell^*}^{(3)}(t)$ into a system given by

$$\begin{aligned} \frac{d}{dt} I_{i_\ell^*, j_\ell^*}^{(3)}(t) &= -\alpha_1 I_{i_\ell^*, j_\ell^*}^{(3)}(t) + \iota_1 \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) + \rho_1 \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) \\ \frac{d}{dt} \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) &= -\alpha_2 I_{i_\ell^*, j_\ell^*}^{(3)}(t) - \iota_2 \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) + \rho_2 \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) \\ \frac{d}{dt} \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) &= -\alpha_3 I_{i_\ell^*, j_\ell^*}^{(3)}(t) + \iota_3 \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) + \rho_3 \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O}\left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)\right) \end{aligned}$$

where

$$\begin{aligned} \alpha_1 &= 36 C_{S,2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) + \frac{6}{a_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) + \frac{18}{b_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_1(t) - \frac{18}{a_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_2(t) \\ \alpha_2 &= \frac{36}{a_i^2} C_{S,2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) - \frac{18}{a_i^2} Q_2(t); \quad \alpha_3 = \frac{36}{b_i^2} C_{S,2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \\ \iota_1 &= \frac{18}{b_i^2} \left(Q_0(t) + \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_1(t) \right) + 18 \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) Q_2(t); \\ \iota_2 &= \frac{6}{a_i^2} \left(\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_2(t) \right) \\ \iota_3 &= \frac{18}{b_i^2} \left(Q_1(t) + \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_2(t) \right); \quad \rho_1 = \frac{24}{a_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_1(t) \\ \rho_2 &= \frac{18}{a_i^2} Q_1(t); \quad \rho_3 = \frac{18}{b_i^2} Q_0(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \end{aligned}$$

□

By Lemma 28, we first need to check that $\iota_2 \rho_3 \geq \iota_3 \rho_2$

$$\begin{aligned} \iota_2 \rho_3 - \iota_3 \rho_2 &= \frac{108}{a_i^2 b_i^2} \left(\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_2(t) \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_0(t) \\ &\quad - \frac{324}{a_i^2 b_i^2} \left(Q_1(t) + \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_2(t) \right) Q_1(t) \end{aligned}$$

By Lemma 29, we have that $Q_2(t) \leq 0$ for all $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0$. Therefore

$$\iota_2 \rho_3 - \iota_3 \rho_2 \geq \frac{108}{a_i^2 b_i^2} Q_1(t) \left(\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 Q_0(t) - 3Q_1(t) \right)$$

Thus the condition holds if $Q_0(t) \geq 4.12Q_1(t)$ given that $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0.9$, which translates to

$$\sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \geq 4.12 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k$$

where we can use

$$\sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k \geq c_0^2 = 0.25$$

and noticing that $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^{-2} Q_1(t) \leq \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!} = \mathbb{E}_{x \sim \mathcal{N}(0,1)} [\pi'(x)^2] \leq 0.05$. Then, we need to check that

$$\alpha_1^2 \iota_2 + \alpha_1^2 \rho_3 + \alpha_1 \alpha_2 \iota_1 + \alpha_1 \alpha_3 \rho_1 + \alpha_1 \iota_2^2 + \alpha_2 \iota_1 \iota_2 + \alpha_1 \rho_3^2 + \alpha_3 \rho_1 \rho_3 > \alpha_2 \iota_3 \rho_1 + \alpha_3 \iota_1 \rho_2$$

To start, we notice that

$$\alpha_1 \alpha_2 - \alpha_3 \rho_2 = \frac{36}{a_i^2} C_{S,2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \left(\alpha_1 - \frac{a_i^2}{b_i^2} \cdot \rho_2 \right) = \frac{36}{a_i^2} C_{S,2} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \left(\alpha_1 - \frac{18}{b_i^2} Q_1(t) \right) > 0$$

when $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0.9$. Thus, we have that $\alpha_1 \alpha_2 \iota_1 > \alpha_3 \iota_1 \rho_2$. Lastly, we also notice that

$$\begin{aligned} \iota_1 \iota_2 - \iota_3 \rho_1 &\geq \frac{108}{a_i^2 b_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \left(Q_0(t) + \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_1(t) \right) Q_1(t) \\ &\quad + \frac{108}{a_i^2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) Q_2(t) - \frac{432}{a_i^2 b_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_1(t)^2 \\ &\geq \frac{108}{a_i^2 b_i^2} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) Q_1(t) \left(\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_0(t) + \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) - 4Q_1(t) \right) \\ &\quad + \frac{108}{a_i^2} \left(\frac{1}{a_i^2} + \frac{1}{b_i^2} \right) \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) Q_1(t) Q_2(t) \end{aligned}$$

which can be numerically verified as $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0.9$. Therefore, by Lemma 28 we have that

$$\begin{aligned} &\max \left\{ \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(T_1 + t) \right|, \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(T_1 + t) \right|, \left| I_{i_\ell^*, i_\ell^*}^{(3)}(T_1 + t) \right| \right\} \\ &\leq e^{-\Omega(t)} \left(\left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(T_1) \right| + \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(T_1) \right| + \left| I_{i_\ell^*, i_\ell^*}^{(3)}(T_1) \right| \right) + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(2)}(t)^3 \right) \\ &\leq e^{-\Omega(t)} \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right) + \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A}, \ell}^{(2)}(t) \right) \end{aligned}$$

Thus, for $t \geq T_1 + \mathcal{O}(\log d)$ we shall have that

$$\max \left\{ \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|, \left| I_{i_\ell^*, i_\ell^*}^{(3)}(t) \right| \right\} \leq \mathcal{O} \left(m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right)$$

The same holds for $\zeta_{i_\ell^*, j}^{(1)}(t), \zeta_{i_\ell^*, j}^{(1)}(t), I_{i, i}^{(3)}(t)$ as their growth is upper bounded by the above.

Bounding $\zeta_{i_\ell^*, j}^{(1)}(t), \zeta_{i_\ell^*, j}^{(2)}(t)$ and $\gamma_{i_\ell^*, j}^{(1)}(t), \gamma_{i_\ell^*, j}^{(1)}(t)$ for $j \in [m^*] \setminus \{j_\ell^*\}$

Lemma 15. *Suppose that the inductive hypothesis in Condition 2 and the initialization condition in Condition 1 hold. Then for all $t \geq 0$ we have that*

$$\begin{aligned} \left| \zeta_{i_\ell^*, j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*, j}^{(2)}(t) \right| &\leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} + m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \right); \forall j \in [m^*] \\ \left| \gamma_{i_\ell^*, j}^{(1)}(t) \right|, \left| \gamma_{i_\ell^*, j}^{(2)}(t) \right| &\leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} + m \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A}, \ell}^{(1)}(t)^2 \right); \forall j \neq j_\ell^* \end{aligned}$$

Proof. For $j \in [m^*] \setminus \{j_\ell^*\}$, we write out the dynamic of $\zeta_{i_\ell^*,j}^{(1)}(t), \zeta_{i_\ell^*,j}^{(2)}(t)$ from Lemma 7

$$\begin{aligned} \frac{d}{dt} \left| \zeta_{i_\ell^*,j}^{(1)}(t) \right| &= -\frac{1}{a_{i_\ell^*}^2} \hat{\lambda}_{i_\ell^*,j_\ell^*,1}(t) \gamma_{i_\ell^*,j_\ell^*}^{(1)}(t) \left| \zeta_{i_\ell^*,j}^{(1)}(t) \right| + \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right) \\ \frac{d}{dt} \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right| &= -\frac{9}{b_{i_\ell^*}^2} \hat{\lambda}_{i_\ell^*,j_\ell^*,5}(t) \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t) \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right| + \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right) \end{aligned}$$

where we applied the upper bound that $\left| \hat{\lambda}_{i_\ell^*,j_\ell^*,2}(t) \gamma_{i_\ell^*,j}^{(1)}(t) \right|, \left| \hat{\lambda}_{i_\ell^*,j_\ell^*,2}(t) \gamma_{i_\ell^*,j}^{(2)}(t) \right| \leq \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right)$ and $\left| \gamma_{i_\ell^*,j}^{(2)}(t) I_{i_\ell^*,i_\ell^*}^{(3)}(t) \right| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$ from Lemma 14. Since both $\hat{\lambda}_{i_\ell^*,j_\ell^*,1}(t) \gamma_{i_\ell^*,j_\ell^*}^{(1)}(t)$ and $\hat{\lambda}_{i_\ell^*,j_\ell^*,5}(t) \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t)$ are positive, we have that $\left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|$ and $\left| \zeta_{i_\ell^*,j}^{(2)}(t) \right|$ enjoys an exponential decay up to $\mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right)$. Recall that at the end of phase 1 we have $\left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} \right)$. Therefore, we can conclude that

$$\left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right)$$

Now, we focus on $\gamma_{i_\ell^*,j}^{(1)}(t)$ and $\gamma_{i_\ell^*,j}^{(2)}(t)$ for $j \neq j_\ell^*$. In particular, by Lemma 7, we have that

$$\begin{aligned} \frac{d}{dt} \gamma_{i_\ell^*,j}^{(1)}(t) &= -\frac{1}{a_{i_\ell^*}^2} \hat{\lambda}_{i_\ell^*,j_\ell^*,1}(t) \gamma_{i_\ell^*,j_\ell^*}^{(1)}(t) \left| \gamma_{i_\ell^*,j}^{(1)}(t) \right| + \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right) \\ \frac{d}{dt} \gamma_{i_\ell^*,j}^{(2)}(t) &= -\frac{9}{b_{i_\ell^*}^3} \hat{\lambda}_{i_\ell^*,j_\ell^*,5}(t) \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t) \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right| + \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right) \end{aligned}$$

Since both $\hat{\lambda}_{i_\ell^*,j_\ell^*,1}(t) \gamma_{i_\ell^*,j_\ell^*}^{(1)}(t)$ and $\hat{\lambda}_{i_\ell^*,j_\ell^*,5}(t) \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t)$ are positive, we have that $\left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|$ and $\left| \gamma_{i_\ell^*,j}^{(2)}(t) \right|$ for $j \neq j_\ell^*$ enjoys an exponential decay up to $\mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right)$. Thus, we can conclude that

$$\left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|, \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right| \leq \mathcal{O} \left(\frac{m^2}{\delta_{\mathbb{P}} \sqrt{d}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right); \forall j \neq j_\ell^*$$

□

Phase 2 Growth of $\gamma_{i_\ell^*,j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t)$. In this section, we show that $\gamma_{i_\ell^*,j_\ell^*}^{(2)}(t)$ continues growing up to at least $1 - \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right)$.

Lemma 16. *Suppose that the inductive hypothesis in Condition 2 and the initialization condition in Condition 1 hold. Then for all $t \geq T_1 + \mathcal{O}(\log d)$, where T_1 is defined in Lemma 13, we have that*

$$\gamma_{i_\ell^*,j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t) \geq 1 - \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} \right)$$

Proof. To start, we need to perform a more fine-grained analysis of the dynamic of $\gamma_{i_\ell^*,j_\ell^*}^{(2)}(t)$. Recall from the proof of Lemma 7 we have that for $(i, j) = (i_\ell^*, j_\ell^*)$

$$\begin{aligned} \sum_{r=1}^m \lambda_{i,r,5}(t) \gamma_{r,j}^{(2)}(t) &= \lambda_{i,i,5}(t) \gamma_{i,j}^{(5)}(t) \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \varepsilon_{\mathcal{A},\ell}^{(2)}(t) \right) \\ \sum_{r=1}^m \lambda_{i,r,2}(t) \zeta_{i,j}^{(1)}(t) &= \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \varepsilon_{5,\ell}(t) \right) \\ \sum_{r=1}^m \lambda_{i,r,3}(t) \zeta_{r,j}^{(1)}(t) &= \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \varepsilon_{5,\ell}(t) \right) \\ \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) \zeta_{i,j}^{(1)}(t) &= \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right) \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O} \left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \end{aligned}$$

and also

$$\begin{aligned}
 \sum_{r=1}^m \lambda_{i,r,2}(t) I_{i,i}^{(3)}(t) &= \pm \mathcal{O} \left(\varepsilon_{5,\ell}(t)^2 + m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 \sum_{r=1}^m \lambda_{i,r,3}(t) I_{r,i}^{(3)}(t) &= \pm \mathcal{O} \left(\varepsilon_{5,\ell}(t)^2 + m \varepsilon_{\mathcal{A},\ell}^{(2)}(t)^3 \right) \\
 \sum_{r=1}^m \lambda_{i,r,5}(t) I_{r,i}^{(2)}(t) &= \lambda_{i,i,5}(t) \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,2}(t) I_{i,i}^{(3)}(t) &= \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,3}(t) \zeta_{i,r}^{(2)}(t) &= \pm \mathcal{O} \left(\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i,r,5}(t) \gamma_{i,r}^{(2)}(t) &= \lambda_{i_\ell^*, j_\ell^*, 5}(t) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right)
 \end{aligned}$$

In this case, we need to refine the bound that

$$\begin{aligned}
 \sum_{r=1}^{m^*} \hat{\lambda}_{i_\ell^*, r, 2}(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i_\ell^*, r, 2}(t) I_{i_\ell^*, i_\ell^*}^{(3)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 2}(t) I_{i_\ell^*, i_\ell^*}^{(3)}(t) \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 \sum_{r=1}^{m^*} \hat{\lambda}_{i_\ell^*, r, 3}(t) \zeta_{i_\ell^*, r}^{(2)}(t) &= \hat{\lambda}_{i_\ell^*, j_\ell^*, 3}(t) \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \pm \mathcal{O} \left(m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right)
 \end{aligned}$$

Denote $\hat{\varepsilon}_\ell(t) = \max \left\{ \left| \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*, j_\ell^*}^{(2)}(t) \right|, \left| I_{i_\ell^*, i_\ell^*}^{(3)}(t) \right| \right\}$. Then we have that

$$\begin{aligned}
 \left| \sum_{r=1}^{m^*} \hat{\lambda}_{i_\ell^*, r, 2}(t) \zeta_{i_\ell^*, j_\ell^*}^{(1)}(t) \right| &\leq \mathcal{O} \left(\hat{\varepsilon}_\ell(t)^2 + m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 \left| \sum_{r=1}^{m^*} \hat{\lambda}_{i_\ell^*, r, 2}(t) I_{i_\ell^*, i_\ell^*}^{(3)}(t) \right| &\leq \mathcal{O} \left(\hat{\varepsilon}_\ell(t)^2 + m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 \left| \sum_{r=1}^{m^*} \hat{\lambda}_{i_\ell^*, r, 3}(t) \zeta_{i_\ell^*, r}^{(2)}(t) \right| &\leq \mathcal{O} \left(\hat{\varepsilon}_\ell(t)^2 + m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right)
 \end{aligned}$$

Noticing that $\varepsilon_{5,\ell}(t) \leq \hat{\varepsilon}_\ell(t)$, we have that

$$\frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \right) \hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) - \mathcal{O} \left(\hat{\varepsilon}_\ell(t)^2 + m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right)$$

By Lemma 21 we have that

$$\hat{\lambda}_{i_\ell^*, j_\ell^*, 5}(t) \geq 2 \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k - \mathcal{O} \left(\hat{\varepsilon}_\ell(t)^2 \right)$$

This gives that

$$\begin{aligned}
 \frac{d}{dt} \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) &\geq 2 \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \right) \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)^k - \mathcal{O} \left(\hat{\varepsilon}_\ell(t)^2 + m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right) \\
 &\geq c_0^2 \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2 \right) - \mathcal{O} \left(\hat{\varepsilon}_\ell(t)^2 + m \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 \right)
 \end{aligned}$$

Let \hat{T} be the first time when $\hat{\varepsilon}_\ell(t) \leq \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$. Recall that for all $t \leq \hat{T}$ we have $\hat{\varepsilon}_\ell(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$. Therefore, we have that either $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 1 - \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d}\right)$ or $\frac{d}{dt}\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq 0$. Thus, we must have that at \hat{T}

$$\gamma_{i_\ell^*, j_\ell^*}^{(2)}(\hat{T}) \geq \gamma_{i_\ell^*, j_\ell^*}^{(2)}(T_1) \geq 0.9$$

For all $t \geq \hat{T}$, we have that

$$\frac{d}{dt}\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq c_0^2 \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2\right) - \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$$

Before $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)$ first reaches $\mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right)$, we have that

$$\frac{d}{dt}\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq \frac{c_0^2}{2} \left(1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t)^2\right)$$

This gives that

$$\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t + \hat{T}) \geq \frac{\left(1 + \gamma_{i_\ell^*, j_\ell^*}^{(2)}(\hat{T})\right) \exp\left(\frac{c_0^2 t}{2}\right) - 1 + \gamma_{i_\ell^*, j_\ell^*}^{(2)}(\hat{T})}{\left(1 + \gamma_{i_\ell^*, j_\ell^*}^{(2)}(\hat{T})\right) \exp\left(\frac{c_0^2 t}{2}\right) + 1 - \gamma_{i_\ell^*, j_\ell^*}^{(2)}(\hat{T})} \geq \frac{\exp\left(\frac{c_0^2 t}{2}\right) - 0.1}{\exp\left(\frac{c_0^2 t}{2}\right) + 0.1}$$

Thus, for $t \geq \mathcal{O}(\log d)$ we have that $\gamma_{i_\ell^*, j_\ell^*}^{(2)}(t + \hat{T}) \geq 1 - \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3\right)$ and stays at that magnitude. We conclude the proof by noticing that $\hat{T} \leq T_1 + \mathcal{O}(\log d)$ by Lemma 14, and the same analysis holds for $\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t)$. \square

With all the preparation work, we are ready to state the lemma for phase 2 convergence.

Lemma 17. *Suppose that the inductive hypothesis in Condition 2, the initialization condition in Condition 1, and the condition on the sigmoid function in Assumption 3 holds. Let T_c be defined in (14). Then for all $t \geq T_1$, where T_1 is defined in Lemma 13, we have that*

$$\gamma_{i_\ell^*, j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*, j_\ell^*}^{(2)}(t) \geq \begin{cases} 0.9 & \text{for } t \geq T_1 \\ 1 - \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) & \text{for } t \geq T_1 + \mathcal{O}(\log d) \end{cases}$$

Moreover, for all $t \geq T_1$, we can upper bound $|\zeta_{i_\ell^*, j}^{(1)}(t)|, |\zeta_{i_\ell^*, j}^{(2)}(t)|$ for $j \in [m^*]$ and $|\gamma_{i_\ell^*, j}^{(1)}(t)|, |\gamma_{i_\ell^*, j}^{(2)}(t)|$ for $j \in [m^*] \setminus \{j_\ell^*\}$ by

$$\begin{aligned} |\zeta_{i_\ell^*, j}^{(1)}(t)|, |\zeta_{i_\ell^*, j}^{(2)}(t)| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right); \forall j \in [m^*] \\ |\gamma_{i_\ell^*, j}^{(1)}(t)|, |\gamma_{i_\ell^*, j}^{(2)}(t)| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right); \forall j \neq j_\ell^* \end{aligned}$$

Proof. The first part of the proof simply follows from a combination of Lemma 14, Lemma 15, and Lemma 16. \square

A.8 Formalizing the Proof of Theorem 1

Now we are ready to prove the theorem for gradient flow.

Proof of Theorem 1. By Lemma 17, we have that

$$\begin{aligned} |\zeta_{i_\ell^*, j}^{(1)}(t)|, |\zeta_{i_\ell^*, j}^{(2)}(t)| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right); \forall j \in [m^*] \\ |\gamma_{i_\ell^*, j}^{(1)}(t)|, |\gamma_{i_\ell^*, j}^{(2)}(t)| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right); \forall j \neq j_\ell^* \end{aligned}$$

By the definition of $\varepsilon_{\mathcal{A},\ell}^{(1)}(t)$, if $\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \geq \varepsilon_{\mathcal{A},\ell+1}^{(2)}(t)$, then we must have that $\left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right|$ and $\left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|, \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right|$ is greater than $\varepsilon_{\mathcal{A},\ell+1}^{(2)}(t)$. In this case, we must have that $\varepsilon_{\mathcal{A},\ell}^{(1)}(t)$ is dominated by $\left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right|$ and $\left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|, \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right|$. For $t \leq T_c$, we can thus conclude that $\left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right|$ and $\left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|, \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right|$ must stay below $\mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$. Otherwise, we will have that $\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \leq \varepsilon_{\mathcal{A},\ell+1}^{(2)}(t)$ and thus

$$\begin{aligned} \left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}} + m\varepsilon_{\mathcal{A},\ell+1}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell+1}^{(1)}(t)^2\right); \forall j \in [m^*] \\ \left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|, \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}} + m\varepsilon_{\mathcal{A},\ell+1}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell+1}^{(1)}(t)^2\right); \forall j \neq j_\ell^* \end{aligned}$$

which implies that for all t such that $\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$, it holds that

$$\begin{aligned} \left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right); \forall j \in [m^*] \\ \left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|, \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right); \forall j \neq j_\ell^* \end{aligned}$$

Combining with the inductive hypothesis that $\varepsilon_{\mathcal{F},\ell}(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ gives that $\varepsilon_{\mathcal{F},\ell+1}(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ when $\varepsilon_{\mathcal{A},\ell+1}^{(1)}(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$. This gives the third inductive hypothesis. For the same reasoning, by Lemma 17, we can also have the first inductive hypothesis. Lastly, the second inductive hypothesis follows from Lemma 9.

Based on the first and second inductive hypothesis, we can conclude bullet point 1-3 in Theorem 1. To see the last statement, we recall from Lemma 17 that

$$\gamma_{i_\ell^*,j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t) \geq 1 - \mathcal{O}\left(m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \text{ for } t \geq T_1 + \mathcal{O}(\log d)$$

Let $T^* \geq T_1 + \mathcal{O}(\log d)$ to be as small as possible for all ℓ . Then $T^* \leq \mathcal{O}(\sqrt{d} + \log d) \leq \mathcal{O}(\sqrt{d})$ by Lemma 13.

By definition of $\varepsilon_{\mathcal{A},\ell}^{(1)}(t)$, we have that $\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \leq \varepsilon_{\mathcal{A},\ell+1}^{(1)}(t)$ for all $\ell \in [m^* - 1]$. Thus, $\varepsilon_{\mathcal{A},\ell}^{(1)}(t) \leq \varepsilon_{\mathcal{A},m^*}^{(1)}(t)$ for all $\ell \in [m^*]$. This gives that

$$\gamma_{i_\ell^*,j_\ell^*}^{(1)}(t), \gamma_{i_\ell^*,j_\ell^*}^{(2)}(t) \geq 1 - \mathcal{O}\left(m\varepsilon_{\mathcal{A},m^*}^{(1)}(t)^3 + \frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right) \text{ for } t \geq T^*$$

Thus, it remains to bound $\varepsilon_{\mathcal{A},m^*}^{(1)}(t)$. By definition, $\varepsilon_{\mathcal{A},m^*}^{(1)}(t)$ depends on $\gamma_{i_\ell^*,j}^{(1)}(t), \gamma_{i_\ell^*,j}^{(2)}(t)$ for all $\ell \in [m^*]$ and $j \in [m^*] \setminus \{j_\ell^*\}$, $\zeta_{i,j}^{(1)}(t), \zeta_{i,j}^{(2)}(t)$ for all $i \in [m]$ and $j \in [m^*]$, and $I_{i,j}^{(1)}(t), I_{i,j}^{(2)}(t), I_{i,j}^{(3)}(t)$ for all $i, j \in [m]$ with $i \neq j$. Recall that by Lemma 17, we have that for all $t \geq T_1$, it holds that

$$\begin{aligned} \left| \zeta_{i_\ell^*,j}^{(1)}(t) \right|, \left| \zeta_{i_\ell^*,j}^{(2)}(t) \right| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right); \forall j \in [m^*] \\ \left| \gamma_{i_\ell^*,j}^{(1)}(t) \right|, \left| \gamma_{i_\ell^*,j}^{(2)}(t) \right| &\leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}} + m\varepsilon_{\mathcal{A},\ell}^{(1)}(t)^3 + \varepsilon_{\mathcal{A},\ell}^{(1)}(t)^2\right); \forall j \neq j_\ell^* \end{aligned}$$

For $\zeta_{i,j}^{(1)}(t), \zeta_{i,j}^{(2)}(t)$ and $I_{i,j}^{(1)}(t), I_{i,j}^{(2)}(t), I_{i,j}^{(3)}(t)$ with $i \in [m] \setminus \mathcal{R}_{m^*}$, we obtain from Lemma 11 that the above are bounded by $\mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ for all $t \leq \mathcal{O}\left(\frac{\delta_{\mathbb{P}}^2 d}{m^5}\right)$. Lastly, for $I_{i,j}^{(1)}(t), I_{i,j}^{(2)}(t), I_{i,j}^{(3)}(t)$ with $i \in \mathcal{R}_{m^*}$ or $j \in \mathcal{R}_{m^*}$, we can obtain from (17), (18), and (19) that

$$\max \left\{ \left| I_{i,j}^{(1)}(t) \right|, \left| I_{i,j}^{(2)}(t) \right|, \left| I_{i,j}^{(3)}(t) \right| \right\} \leq \max \left\{ \left| \gamma_{j,j_\ell^*}^{(1)}(t) \right|, \left| \gamma_{j,j_\ell^*}^{(2)}(t) \right|, \left| \zeta_{j,j_\ell^*}^{(1)}(t) \right|, \left| \zeta_{j,j_\ell^*}^{(2)}(t) \right| \right\} + \mathcal{O}\left(\frac{m^4}{\delta_{\mathbb{P}}^2 d^{\frac{3}{4}}}\right)$$

As $|\zeta_{j,j_{\ell'}^*}^{(1)}(t)|, |\zeta_{j,j_{\ell'}^*}^{(1)}(t)|$ are bounded above, we just need to look into $|\gamma_{j,j_{\ell'}^*}^{(1)}(t)|, |\gamma_{j,j_{\ell'}^*}^{(2)}(t)|$ which are bounded by $\mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ as shown in Lemma 10 when $t \leq \left\{T_{m^*}(\xi) + \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right), \mathcal{O}\left(\frac{\delta_{\mathbb{P}}^2 d^{\frac{3}{2}}}{m^5}\right)\right\}$. Combining all the bounds above we can conclude that $\varepsilon_{\mathcal{A},m^*}^{(1)}(t) \leq \mathcal{O}\left(\frac{m^2}{\delta_{\mathbb{P}}\sqrt{d}}\right)$ for all $t \leq \min\left\{T_{m^*}(\xi) + \mathcal{O}\left(\frac{\delta_{\mathbb{P}}\sqrt{d}}{m^2}\right)\right\}$. Thus, we can obtain that

$$\gamma_{i_{\ell}^*,j_{\ell}^*}^{(1)}(t), \gamma_{i_{\ell}^*,j_{\ell}^*}^{(2)}(t) \geq 1 - \mathcal{O}\left(\frac{m^7}{\delta_{\mathbb{P}}^3 d^{\frac{3}{2}}}\right); \forall T^* \leq t \leq T^* + \mathcal{O}\left(\frac{\delta_{\mathbb{P}}\sqrt{d}}{m^2}\right)$$

□

B Proof of Theorem 2

Proof of Theorem 2. We simply need to show that under the stopping criteria (8) the procedure in (7) satisfies that $r_{\tau} \in [m] \setminus [m^*]$ for all $\tau \in [\tau^*]$ and that $\tau^* = m - m^*$. This is done in **Part 1** and **Part 2** below. Before we start, we define the loss over the pruned model as

$$\mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}}\left[(f_{\mathcal{S}}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}))^2\right]$$

and we write the target pruned model as

$$\hat{f}(\boldsymbol{\theta}, \mathbf{x}) = \sum_{i=1}^{m^*} \pi(\bar{\mathbf{v}}_i^{\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^{\top} \mathbf{x})$$

For the convenience, we also denote $h_i(\mathbf{x}) := \pi(\bar{\mathbf{v}}_i^{\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^{\top} \mathbf{x})$ and $h_i^*(\mathbf{x}) := \pi(\bar{\mathbf{v}}_i^{*\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^{*\top} \mathbf{x})$.

Part 1. Assume that $\mathcal{S}_{\tau-1} \subseteq [m] \setminus [m^*]$ and $|\mathcal{S}_{\tau-1}| \leq m - m^*$. Let $r_{\tau} \in [m] \setminus [m^*]$ and $r'_{\tau} \in [m^*]$. Let $\mathcal{S}_{\tau} = \mathcal{S}_{\tau-1} \cup \{r_{\tau}\}$ and $\mathcal{S}'_{\tau} = \mathcal{S}_{\tau-1} \cup \{r'_{\tau}\}$. Moreover, let $\mathcal{S}^{\perp} = [m] \setminus [m^*] \setminus \mathcal{S}_{\tau}$. Then we have that

$$\begin{aligned} \mathcal{L}_{\mathcal{S}_{\tau}}(\boldsymbol{\theta}) - \mathcal{L}_{\mathcal{S}'_{\tau}}(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x}}\left[(f_{\mathcal{S}_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}))^2 - (f_{\mathcal{S}'_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}))^2\right] \\ &= \mathbb{E}_{\mathbf{x}}\left[(f_{\mathcal{S}_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) + f_{\mathcal{S}'_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) - 2f^*(\mathbf{x})) (f_{\mathcal{S}_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) - f_{\mathcal{S}'_{\tau}}(\boldsymbol{\theta}, \mathbf{x}))\right] \\ &= \mathbb{E}_{\mathbf{x}}\left[(f_{\mathcal{S}_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) + f_{\mathcal{S}'_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) - 2\hat{f}(\boldsymbol{\theta}, \mathbf{x})) (f_{\mathcal{S}_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) - f_{\mathcal{S}'_{\tau}}(\boldsymbol{\theta}, \mathbf{x}))\right] \\ &\quad - 2 \underbrace{\mathbb{E}_{\mathbf{x}}\left[(\hat{f}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) (f_{\mathcal{S}_{\tau}}(\boldsymbol{\theta}, \mathbf{x}) - f_{\mathcal{S}'_{\tau}}(\boldsymbol{\theta}, \mathbf{x}))\right]}_{\mathcal{T}_1} \\ &= \mathbb{E}_{\mathbf{x}}\left[\left(2 \sum_{i \in \mathcal{S}^{\perp}} h_i(\mathbf{x}) + h_{r_{\tau}}(\mathbf{x}) - h_{r'_{\tau}}(\mathbf{x})\right) (h_{r'_{\tau}}(\mathbf{x}) - h_{r_{\tau}}(\mathbf{x}))\right] - 2\mathcal{T}_1 \\ &= -\mathbb{E}_{\mathbf{x}}\left[(h_{r'_{\tau}}(\mathbf{x}) - h_{r_{\tau}}(\mathbf{x}))^2\right] - 2\mathcal{T}_1 + 2 \underbrace{\mathbb{E}_{\mathbf{x}}\left[\sum_{i \in \mathcal{S}^{\perp}} h_i(\mathbf{x}) (h_{r'_{\tau}}(\mathbf{x}) - h_{r_{\tau}}(\mathbf{x}))\right]}_{\mathcal{T}_2} \\ &\leq -\mathbb{E}_{\mathbf{x}}\left[(h_{r'_{\tau}}(\mathbf{x}) - h_{r_{\tau}}(\mathbf{x}))^2\right] + 2|\mathcal{T}_1| + 2|\mathcal{T}_2| \end{aligned}$$

It suffice to upper bound $|\mathcal{T}_1|$ and $|\mathcal{T}_2|$, and lower bound $\mathbb{E}_{\mathbf{x}}\left[(h_{r'_{\tau}}(\mathbf{x}) - h_{r_{\tau}}(\mathbf{x}))^2\right]$. To start, the lower bound can be derived as

$$\mathbb{E}_{\mathbf{x}}\left[(h_{r'_{\tau}}(\mathbf{x}) - h_{r_{\tau}}(\mathbf{x}))^2\right] = \mathbb{E}_{\mathbf{x}}\left[h_{r'_{\tau}}(\mathbf{x})^2\right] + \mathbb{E}_{\mathbf{x}}\left[h_{r_{\tau}}(\mathbf{x})^2\right] - 2\mathbb{E}_{\mathbf{x}}\left[h_{r_{\tau}}(\mathbf{x}) h_{r'_{\tau}}(\mathbf{x})\right] \geq 12 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} - \mathcal{O}(\varepsilon^2) \quad (20)$$

where the last inequality follows from Lemma 22. Next, for $|\mathcal{T}_1|$, we have that

$$\begin{aligned}
 \mathcal{T}_1 &= \mathbb{E}_{\mathbf{x}} \left[\left(\hat{f}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}) \right) \left(h_{r_\tau}(\mathbf{x}) - h_{r'_\tau}(\mathbf{x}) \right) \right] \\
 &= \mathbb{E}_{\mathbf{x}} \left[h_{r_\tau}(\mathbf{x}) \sum_{i=1}^{m^*} h_i(\mathbf{x}) \right] - \mathbb{E}_{\mathbf{x}} \left[h_{r_\tau}(\mathbf{x}) \sum_{i=1}^{m^*} h_i^*(\mathbf{x}) \right] - \mathbb{E}_{\mathbf{x}} \left[h'_{r_\tau}(\mathbf{x}) \sum_{i=1, i \neq r'_\tau}^{m^*} h_i(\mathbf{x}) \right] \\
 &\quad + \mathbb{E}_{\mathbf{x}} \left[h'_{r_\tau}(\mathbf{x}) \sum_{i=1, i \neq r'_\tau}^{m^*} h_i^*(\mathbf{x}) \right] + \mathbb{E}_{\mathbf{x}} \left[h_{r'_\tau}(\mathbf{x})^2 - h_{r'_\tau}(\mathbf{x}) h_{r'_\tau}^*(\mathbf{x}) \right] \\
 &= \mathbb{E}_{\mathbf{x}} \left[h_{r'_\tau}(\mathbf{x})^2 - h_{r'_\tau}(\mathbf{x}) h_{r'_\tau}^*(\mathbf{x}) \right] \pm \mathcal{O}(m^* \varepsilon^4)
 \end{aligned} \tag{21}$$

By Lemma 22, since $r'_\tau \in [m^*]$, we have that

$$\mathbb{E}_{\mathbf{x}} \left[h_{r'_\tau}(\mathbf{x}) h_{r'_\tau}^*(\mathbf{x}) \right] = 6 \mathbb{E}_{\mathbf{x}} \left[h_{r'_\tau}(\mathbf{x})^2 \right] \pm \mathcal{O}(\varepsilon) \tag{22}$$

Thus, we have that

$$|\mathcal{T}_1| \leq \mathcal{O}(\varepsilon + m^* \varepsilon^4) \tag{23}$$

For $|\mathcal{T}_2|$, we notice that $r_\tau \in \mathcal{S}_\tau$ and $r'_\tau \in [m^*]$. Therefore, $r_\tau, r'_\tau \notin \mathcal{S}^\perp$. Thus

$$|\mathcal{T}_2| \leq \sum_{i \in \mathcal{S}^\perp} |\mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x}) h_{r_\tau}(\mathbf{x})]| + \sum_{i \in \mathcal{S}^\perp} |\mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x}) h_{r'_\tau}(\mathbf{x})]| \leq \mathcal{O}(m \varepsilon^4) \tag{24}$$

Combining (20), (23), and (24) gives that

$$\mathcal{L}_{\mathcal{S}_\tau}(\boldsymbol{\theta}) - \mathcal{L}_{\mathcal{S}'_\tau}(\boldsymbol{\theta}) \leq -12 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} + \mathcal{O}(\varepsilon + m \varepsilon^4) \leq 0$$

when $\varepsilon \leq o\left(\frac{1}{\sqrt{m}}\right)$.

Part 2. Assume that $\tau^* < m - m^*$. We show that $\mathcal{L}_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}) \geq \mathcal{L}_{\mathcal{S}_{\tau^*+1}}(\boldsymbol{\theta} + 1)$ by letting $\mathcal{S}_{\tau^*+1} = \mathcal{S}_{\tau^*} \cup \{r^*\}$ where $r^* \in [m] \setminus [m^*]$. Similar to before, let $\mathcal{S}^\perp = [m] \setminus [m^*] \setminus \mathcal{S}_{\tau^*+1}$, we have that

$$\begin{aligned}
 \mathcal{L}_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}) - \mathcal{L}_{\mathcal{S}_{\tau^*+1}}(\boldsymbol{\theta} + 1) &= \mathbb{E}_{\mathbf{x}} \left[\left(f_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}) \right)^2 - \left(f_{\mathcal{S}_{\tau^*+1}}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}) \right)^2 \right] \\
 &= \mathbb{E}_{\mathbf{x}} \left[\left(f_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}, \mathbf{x}) + f_{\mathcal{S}_{\tau^*+1}}(\boldsymbol{\theta}, \mathbf{x}) - 2f^*(\mathbf{x}) \right) \left(f_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}, \mathbf{x}) - f_{\mathcal{S}_{\tau^*+1}}(\boldsymbol{\theta}, \mathbf{x}) \right) \right] \\
 &= \mathbb{E}_{\mathbf{x}} \left[h_{r^*}(\mathbf{x}) \left(f_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}, \mathbf{x}) + f_{\mathcal{S}_{\tau^*+1}}(\boldsymbol{\theta}, \mathbf{x}) - 2f^*(\mathbf{x}) \right) \right] \\
 &= \mathbb{E}_{\mathbf{x}} \left[h_{r^*}(\mathbf{x}) \left(f_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}, \mathbf{x}) + f_{\mathcal{S}_{\tau^*+1}}(\boldsymbol{\theta}, \mathbf{x}) - 2\hat{f}(\boldsymbol{\theta}, \mathbf{x}) \right) \right] \\
 &\quad - 2 \underbrace{\mathbb{E}_{\mathbf{x}} \left[h_{r^*}(\mathbf{x}) \left(\hat{f}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}) \right) \right]}_{\mathcal{T}_1} \\
 &= \mathbb{E}_{\mathbf{x}} \left[h_{r^*}(\mathbf{x})^2 \right] + 2 \underbrace{\mathbb{E}_{\mathbf{x}} \left[h_{r^*}(\mathbf{x}) \sum_{i \in \mathcal{S}^\perp} h_i(\mathbf{x}) \right]}_{\mathcal{T}_2} - 2\mathcal{T}_1 \\
 &\geq \mathbb{E}_{\mathbf{x}} \left[h_{r^*}(\mathbf{x})^2 \right] - 2|\mathcal{T}_1| - 2|\mathcal{T}_2|
 \end{aligned}$$

As before, we have that $\mathbb{E}_{\mathbf{x}} \left[h_{r^*}(\mathbf{x})^2 \right] \geq 6 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} - \mathcal{O}(\varepsilon^2)$ and $|\mathcal{T}_2| \leq \mathcal{O}(m \varepsilon^4)$. For \mathcal{T}_1 , we have that

$$\begin{aligned}
 |\mathcal{T}_1| &= \left| \mathbb{E}_{\mathbf{x}} \left[h_{r^*}(\mathbf{x}) \left(\hat{f}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x}) \right) \right] \right| \\
 &\leq \sum_{i=1}^{m^*} |\mathbb{E}_{\mathbf{x}}[h_{r^*}(\mathbf{x}) h_i(\mathbf{x})]| + \sum_{i=1}^{m^*} |\mathbb{E}_{\mathbf{x}}[h_{r^*}(\mathbf{x}) h_i^*(\mathbf{x})]| \\
 &\leq \mathcal{O}(m \varepsilon^4)
 \end{aligned}$$

Thus, we have that

$$\mathcal{L}_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta}) - \mathcal{L}_{\mathcal{S}_{\tau^*}}(\boldsymbol{\theta} + 1) \geq 6 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} - \mathcal{O}(\varepsilon^2 + m\varepsilon^4) \geq 0$$

when $\varepsilon \leq \mathcal{O}\left(\frac{1}{\sqrt{m}}\right)$. This shows that $\tau^* \geq m - m^*$. Next, we assume that $\tau^* > m - m^*$. Then $r_{m-m^*+1} \in [m^*]$.

We show that $\mathcal{L}_{\mathcal{S}_{m-m^*}}(\boldsymbol{\theta}) \leq \mathcal{L}_{\mathcal{S}_{m-m^*+1}}(\boldsymbol{\theta})$. As before, let \cdot . Notice that by **Part 1**, $f_{\mathcal{S}_{m-m^*}}(\boldsymbol{\theta}, \mathbf{x}) = \hat{f}(\boldsymbol{\theta}, \mathbf{x})$. Then we have that

$$\begin{aligned} \mathcal{L}_{\mathcal{S}_{m-m^*}}(\boldsymbol{\theta}) - \mathcal{L}_{\mathcal{S}_{m-m^*+1}}(\boldsymbol{\theta}) &= \mathbb{E}_{\mathbf{x}} \left[h_{r_{m-m^*+1}}(\mathbf{x}) (f_{\mathcal{S}_{m-m^*}}(\boldsymbol{\theta}, \mathbf{x}) + f_{\mathcal{S}_{m-m^*+1}}(\boldsymbol{\theta}, \mathbf{x}) - 2f^*(\mathbf{x})) \right] \\ &= 2\mathbb{E}_{\mathbf{x}} \left[h_{r_{m-m^*+1}}(\mathbf{x}) (\hat{f}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \right] - \mathbb{E}_{\mathbf{x}} \left[h_{r_{m-m^*+1}}(\mathbf{x})^2 \right] \end{aligned}$$

As before, we have that $\mathbb{E}_{\mathbf{x}} \left[h_{r_{m-m^*+1}}(\mathbf{x})^2 \right] \geq 6 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} - \mathcal{O}(\varepsilon^2)$. It remains to upper bound the first term. In particular, we have that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}} \left[h_{r_{m-m^*+1}}(\mathbf{x}) (\hat{f}(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \right] &\leq \mathbb{E}_{\mathbf{x}} \left[h_{r_{m-m^*+1}}(\mathbf{x})^2 - h_{r_{m-m^*+1}}(\mathbf{x}) h_{r_{m-m^*+1}}^*(\mathbf{x}) \right] \\ &\quad + \sum_{i \neq r_{m-m^*+1}} \mathbb{E}_{\mathbf{x}} \left[h_{r_{m-m^*+1}}(\mathbf{x}) (h_i(\mathbf{x}) - h_i^*(\mathbf{x})) \right] \\ &\leq \mathcal{O}(\varepsilon + m^* \varepsilon^4) \end{aligned}$$

Thus, we can conclude that

$$\mathcal{L}_{\mathcal{S}_{m-m^*}}(\boldsymbol{\theta}) - \mathcal{L}_{\mathcal{S}_{m-m^*+1}}(\boldsymbol{\theta}) \leq -6 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} + \mathcal{O}(\varepsilon + m\varepsilon^4) \leq 0$$

when $\varepsilon \leq \mathcal{O}\left(\frac{1}{\sqrt{m}}\right)$. This shows that $\tau^* = m - m^*$, which finishes the proof. \square

C Proof of Theorem 3

We will analyze the Hessian in a small region near the global minima $\boldsymbol{\theta}^* = \{(\bar{\mathbf{v}}_i^*, \bar{\mathbf{w}}_i^*)\}_{i=1}^{m^*}$. To do this, we utilize the following second-order Stein's lemma.

Lemma 18. *Let $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$ and $g, h : \mathbb{R} \rightarrow \mathbb{R}$. Then we have that*

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)} [g(\mathbf{v}^\top \mathbf{x}) h(\mathbf{w}^\top \mathbf{x}) \mathbf{x} \mathbf{x}^\top] &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)} [g(\mathbf{v}^\top \mathbf{x}) h(\mathbf{w}^\top \mathbf{x})] \mathbf{I}_d \\ &\quad + \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)} [g'(\mathbf{v}^\top \mathbf{x}) h'(\mathbf{w}^\top \mathbf{x})] (\mathbf{v} \mathbf{w}^\top + \mathbf{w} \mathbf{v}^\top) \\ &\quad + \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)} [g''(\mathbf{v}^\top \mathbf{x}) h(\mathbf{w}^\top \mathbf{x})] \mathbf{v} \mathbf{v}^\top \\ &\quad + \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)} [g(\mathbf{v}^\top \mathbf{x}) h''(\mathbf{w}^\top \mathbf{x})] \mathbf{w} \mathbf{w}^\top \end{aligned}$$

The proof of Lemma 18 follows by applying Stein's lemma twice. In particular, we shall prove the following result

Theorem 4. *Let $\boldsymbol{\theta} = \{(\mathbf{v}_i, \mathbf{w}_i)\}_{i=1}^{m^*}$ be the parameter of the MoE, let $\alpha_1, \dots, \alpha_{m^*}, \beta_1, \dots, \beta_{m^*} \geq \Omega(1)$, and let $\mathbf{u}_1, \dots, \mathbf{u}_{m^*}, \mathbf{q}_1, \dots, \mathbf{q}_{m^*} \in \mathbb{R}^d$ be any set of vectors such that $\|(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathbf{u}_i\|_2^2 \geq \mathcal{Q}^* \|\mathbf{u}_i\|_2^2$ and $\|(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathbf{q}_i\|_2^2 \geq \mathcal{Q}^* \|\mathbf{q}_i\|_2^2$ for some $\mathcal{Q}^* > 0$ for all $i \in [m^*]$. If $\boldsymbol{\theta}$ also satisfies that $\|\bar{\mathbf{v}}_i - \bar{\mathbf{v}}_i^*\|_2, \|\bar{\mathbf{w}}_i - \bar{\mathbf{w}}_i^*\|_2 \leq \frac{\varepsilon}{2}$ for some $\varepsilon \leq o\left(\frac{N_{\min} \mathcal{Q}^*}{m^{*2}}\right)$, and $\frac{C_{S,0}}{C_{S,1}} \geq \frac{N_{\max}^2 (1+\beta_i)^2}{N_{\min}^2 \mathcal{Q}^{*2} \alpha_i^2}$, then we have that*

$$\begin{bmatrix} \alpha_1 \mathbf{u}_1 \\ \vdots \\ \alpha_{m^*} \mathbf{u}_{m^*} \\ \beta_1 \mathbf{q}_1 \\ \vdots \\ \beta_{m^*} \mathbf{q}_{m^*} \end{bmatrix}^\top \nabla^2 \mathcal{L}(\boldsymbol{\theta}) \begin{bmatrix} \mathbf{u}_1 \\ \vdots \\ \mathbf{u}_{m^*} \\ \mathbf{q}_1 \\ \vdots \\ \mathbf{q}_{m^*} \end{bmatrix} \geq N_{\min} \mathcal{Q}^* \kappa \sum_{i=1}^{m^*} \left(\|\mathbf{u}_i\|_2^2 + \|\mathbf{q}_i\|_2^2 \right)$$

for some constant $\kappa > 0$.

Proof. Form of Hessian. Here we are going to compute $\frac{\partial^2}{\partial \mathbf{v}_i \partial \mathbf{v}_j} \mathcal{L}(\boldsymbol{\theta})$, $\frac{\partial^2}{\partial \mathbf{w}_i \partial \mathbf{w}_j} \mathcal{L}(\boldsymbol{\theta})$, and $\frac{\partial^2}{\partial \mathbf{v}_i \partial \mathbf{w}_j} \mathcal{L}(\boldsymbol{\theta})$. Recall that the gradient takes the form

$$\begin{aligned}\frac{\partial}{\partial \mathbf{v}_i} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{\|\mathbf{v}_i\|_2} (\mathbf{I}_d - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) \mathbf{x}] \\ \frac{\partial}{\partial \mathbf{w}_i} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{\|\mathbf{w}_i\|_2} (\mathbf{I}_d - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}) \mathbf{x}]\end{aligned}$$

Therefore

$$\begin{aligned}\frac{\partial^2}{\partial \mathbf{v}_i \partial \mathbf{v}_j} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{\|\mathbf{v}_i\|_2 \|\mathbf{v}_j\|_2} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \underbrace{\mathbb{E}_{\mathbf{x}}[\pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) \pi'(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^\top \mathbf{x}) \mathbf{x} \mathbf{x}^\top]}_{\mathcal{T}_{i,j,1}} (\mathbf{I} - \bar{\mathbf{v}}_j \bar{\mathbf{v}}_j^\top) \\ &\quad + \frac{\mathbb{I}\{i=j\}}{\|\mathbf{v}_i\|_2^2} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \underbrace{\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi''(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) \mathbf{x} \mathbf{x}^\top]}_{\mathcal{T}_{i,j,2}} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \\ &\quad - \frac{\mathbb{I}\{i=j\}}{\|\mathbf{v}_i\|_2^2} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \underbrace{\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) \bar{\mathbf{v}}_i^\top \mathbf{x}]}_{\mathcal{T}_{i,j,3}} \\ &\quad - \frac{\mathbb{I}\{i=j\}}{\|\mathbf{v}_i\|_2^2} \underbrace{\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) (\bar{\mathbf{v}}_i \mathbf{x}^\top + \mathbf{x} \bar{\mathbf{v}}_i^\top)]}_{\mathcal{T}_{i,j,4}} \\ \frac{\partial^2}{\partial \mathbf{w}_i \partial \mathbf{w}_j} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{\|\mathbf{w}_i\|_2 \|\mathbf{w}_j\|_2} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \underbrace{\mathbb{E}_{\mathbf{x}}[\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}) \pi(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_j^\top \mathbf{x}) \mathbf{x} \mathbf{x}^\top]}_{\mathcal{T}_{i,j,5}} (\mathbf{I} - \bar{\mathbf{w}}_j \bar{\mathbf{w}}_j^\top) \\ &\quad + \frac{\mathbb{I}\{i=j\}}{\|\mathbf{w}_i\|_2^2} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \underbrace{\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma''(\bar{\mathbf{w}}_i^\top \mathbf{x}) \mathbf{x} \mathbf{x}^\top]}_{\mathcal{T}_{i,j,6}} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \\ &\quad - \frac{\mathbb{I}\{i=j\}}{\|\mathbf{w}_i\|_2^2} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \underbrace{\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}) \bar{\mathbf{w}}_i^\top \mathbf{x}]}_{\mathcal{T}_{i,j,7}} \\ &\quad - \frac{\mathbb{I}\{i=j\}}{\|\mathbf{w}_i\|_2^2} \underbrace{\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}) (\bar{\mathbf{w}}_i \mathbf{x}^\top + \mathbf{x} \bar{\mathbf{w}}_i^\top)]}_{\mathcal{T}_{i,j,8}} \\ \frac{\partial^2}{\partial \mathbf{v}_i \partial \mathbf{w}_j} \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{\|\mathbf{v}_i\|_2 \|\mathbf{w}_j\|_2} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \underbrace{\mathbb{E}_{\mathbf{x}}[\pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) \pi(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_j^\top \mathbf{x}) \mathbf{x} \mathbf{x}^\top]}_{\mathcal{T}_{i,j,9}} (\mathbf{I} - \bar{\mathbf{w}}_j \bar{\mathbf{w}}_j^\top) \\ &\quad + \frac{\mathbb{I}\{i=j\}}{\|\mathbf{v}_i\|_2 \|\mathbf{w}_i\|_2} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \underbrace{\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}) \mathbf{x} \mathbf{x}^\top]}_{\mathcal{T}_{i,j,10}} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)\end{aligned}$$

For the convenience of the analysis, we define $C_{S,0} = 2 \sum_{k=0}^{\infty} \frac{c_k^2}{k!}$ and $C_{S,1} = 6 \sum_{k=0}^{\infty} \frac{c_{k+1}^2}{k!}$. Our next lemma controls the magnitudes of these blocks.

Bounding $\mathcal{T}_{i,j,2}$, $\mathcal{T}_{i,j,6}$ and $\mathcal{T}_{i,j,10}$. By Lemma 18, we have that

$$\begin{aligned}\mathcal{T}_{i,j,2} &= \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi''(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})] \mathbf{I}_d \\ &\quad + \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}^2 (\pi''(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \\ &\quad + \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}^2 (f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi''(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})] \\ &\quad + \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}} (f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}} (\pi''(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}))^\top]\end{aligned}$$

Similarly, we have that

$$\begin{aligned} \mathcal{T}_{i,j,6} &= \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma''(\bar{\mathbf{w}}_i^\top \mathbf{x})] \mathbf{I}_d \\ &\quad + \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}^2(\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma''(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \\ &\quad + \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}^2(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma''(\bar{\mathbf{w}}_i^\top \mathbf{x})] \\ &\quad + \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}(\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma''(\bar{\mathbf{w}}_i^\top \mathbf{x}))]^\top \end{aligned}$$

Also, we have that

$$\begin{aligned} \mathcal{T}_{i,j,10} &= \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x})] \mathbf{I}_d \\ &\quad + \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}^2(\pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \\ &\quad + \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}^2(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x})] \\ &\quad + \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}(\pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}))]^\top \end{aligned}$$

Thus, we can apply Lemma 23 to obtain that

$$\|\mathcal{T}_{i,j,2}\|_2, \|\mathcal{T}_{i,j,6}\|_2, \|\mathcal{T}_{i,j,10}\|_2 \leq \mathcal{O}(m^* \varepsilon)$$

Bounding $\mathcal{T}_{i,j,3}$ and $\mathcal{T}_{i,j,7}$ By Stein's Lemma, we have that

$$\begin{aligned} \mathcal{T}_{i,j,3} &= \mathbb{E}_{\mathbf{x}}[\bar{\mathbf{v}}_i^\top \nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})] \\ &\quad + \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \bar{\mathbf{v}}_i^\top \nabla_{\mathbf{x}}(\pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \\ \mathcal{T}_{i,j,7} &= \mathbb{E}_{\mathbf{x}}[\bar{\mathbf{w}}_i^\top \nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x})] \\ &\quad + \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \bar{\mathbf{w}}_i^\top \nabla_{\mathbf{x}}(\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \end{aligned}$$

Therefore, we have that

$$\begin{aligned} |\mathcal{T}_{i,j,3}| &\leq \left\| \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})] \right\|_2 \\ &\quad + \left\| \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}(\pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \right\|_2 \\ |\mathcal{T}_{i,j,7}| &\leq \left\| \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x})] \right\|_2 \\ &\quad + \left\| \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}(\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \right\|_2 \end{aligned}$$

Applying Lemma 23 gives that

$$|\mathcal{T}_{i,j,3}|, |\mathcal{T}_{i,j,7}| \leq \mathcal{O}(m^* \varepsilon)$$

Bounding $\mathcal{T}_{i,j,4}$ and $\mathcal{T}_{i,j,8}$. By the structure of $\mathcal{T}_{i,j,4}$ and $\mathcal{T}_{i,j,8}$, we have that

$$\begin{aligned} \|\mathcal{T}_{i,j,4}\|_2 &\leq 2 \left\| \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})] \mathbf{x} \right\|_2 \\ &\leq 2 \left\| \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})] \right\|_2 \\ &\quad + 2 \left\| \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}(\pi'(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \right\|_2 \end{aligned}$$

Similarly, for $\mathcal{T}_{i,j,8}$, we have that

$$\begin{aligned} \|\mathcal{T}_{i,j,8}\|_2 &\leq 2 \left\| \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x})] \mathbf{x} \right\|_2 \\ &\leq 2 \left\| \mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x})] \right\|_2 \\ &\quad + 2 \left\| \mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}(\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_i^\top \mathbf{x}))] \right\|_2 \end{aligned}$$

Thus, we have that

$$\begin{aligned}
 & (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathcal{T}_{i,i,1} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \\
 &= \mathbb{E}_{\mathbf{x}} \left[\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \\
 &\quad + 2 \mathbb{E}_{\mathbf{x}} \left[\left(\sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 + \sigma'' (\bar{\mathbf{w}}_i^\top \mathbf{x}) \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x}) \right) \pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \right] (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \\
 &= \mathbb{E}_{\mathbf{x}} \left[\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \\
 &\quad + 2 \mathbb{E}_{\mathbf{x}} \left[\sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \right] (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) + \hat{\mathcal{T}}_{i,i,1} \\
 & (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathcal{T}_{i,i,5} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \\
 &= \mathbb{E}_{\mathbf{x}} \left[\pi (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \\
 &\quad + 2 \mathbb{E}_{\mathbf{x}} \left[\left(\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 + \pi'' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi (\bar{\mathbf{v}}_i^\top \mathbf{x}) \right) \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top)
 \end{aligned}$$

with $\|\hat{\mathcal{T}}_{i,i,1}\|_2 \leq \mathcal{O}(\varepsilon)$. This gives that

$$\begin{aligned}
 \mathbf{u}^\top (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathcal{T}_{i,i,1} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathbf{u} &\geq \mathbb{E}_{\mathbf{x}} \left[\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] \|(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathbf{u}\|_2^2 - \mathcal{O}(\varepsilon) \\
 \mathbf{u}^\top (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathcal{T}_{i,i,5} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathbf{u} &\geq \mathbb{E}_{\mathbf{x}} \left[\pi (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] \|(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathbf{u}\|_2^2
 \end{aligned}$$

Imposing the condition that $\|(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathbf{u}\|_2^2 \geq \mathcal{Q}^* \|\mathbf{u}\|_2^2$ gives that

$$\begin{aligned}
 \mathbf{u}^\top (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathcal{T}_{i,i,1} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathbf{u} &\geq \mathcal{Q}^* \mathbb{E}_{\mathbf{x}} \left[\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] \|\mathbf{u}\|_2^2 - \mathcal{O}(\varepsilon) \\
 \mathbf{u}^\top (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathcal{T}_{i,i,5} (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathbf{u} &\geq \mathcal{Q}^* \mathbb{E}_{\mathbf{x}} \left[\pi (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] \|\mathbf{u}\|_2^2
 \end{aligned}$$

For $\mathcal{T}_{i,i,10}$, we first notice that

$$\begin{aligned}
 & \left| \mathbb{E}_{\mathbf{x}} \left[\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi (\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x}) \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon) \\
 & \left| \mathbb{E}_{\mathbf{x}} \left[\pi'' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi (\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma'' (\bar{\mathbf{w}}_i^\top \mathbf{x}) \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon) \\
 & \left| \mathbb{E}_{\mathbf{x}} \left[\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 \sigma'' (\bar{\mathbf{w}}_i^\top \mathbf{x}) \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon) \\
 & \left| \mathbb{E}_{\mathbf{x}} \left[\left(\pi''' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 + 2\pi'' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \right) \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x}) \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon) \\
 & \left| \mathbb{E}_{\mathbf{x}} \left[\left(2\sigma'' (\bar{\mathbf{w}}_i^\top \mathbf{x}) \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 + \sigma''' (\bar{\mathbf{w}}_i^\top \mathbf{x}) \sigma (\bar{\mathbf{w}}_i^\top \mathbf{x}) \right) \pi' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi (\bar{\mathbf{v}}_i^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon)
 \end{aligned}$$

Therefore, we have that

$$\begin{aligned}
 & (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathcal{T}_{i,i,10} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \\
 &= \mathbb{E}_{\mathbf{x}} \left[\left(\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 + \pi'' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi (\bar{\mathbf{v}}_i^\top \mathbf{x}) \right) \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \bar{\mathbf{v}}_i \bar{\mathbf{w}}_i^\top (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top)
 \end{aligned}$$

This implies that

$$\begin{aligned}
 \mathbf{u}_1^\top (\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) \mathcal{T}_{i,i,10} (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) \mathbf{u}_2 \\
 \leq \mathbb{E}_{\mathbf{x}} \left[\left(\pi' (\bar{\mathbf{v}}_i^\top \mathbf{x})^2 + \pi'' (\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi (\bar{\mathbf{v}}_i^\top \mathbf{x}) \right) \sigma' (\bar{\mathbf{w}}_i^\top \mathbf{x})^2 \right] \|\mathbf{u}_1\|_2 \|\mathbf{u}_2\|_2
 \end{aligned}$$

□

Now, we are ready to prove the fine-tuning convergence.

Proof of Theorem 3. By the mean-value theorem, we have that

$$\nabla \mathcal{L}(\boldsymbol{\theta}) = \nabla \mathcal{L}(\boldsymbol{\theta}^*) + \nabla^2 \mathcal{L}(\hat{\boldsymbol{\theta}})(\boldsymbol{\theta} - \boldsymbol{\theta}^*) = \nabla^2 \mathcal{L}(\hat{\boldsymbol{\theta}})(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$$

for some $\hat{\boldsymbol{\theta}} \in [\boldsymbol{\theta}, \boldsymbol{\theta}^*]$. The gradient flow dynamic implies that

$$\begin{aligned} \frac{d}{dt} \|\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*\|_2^2 &= \left\langle \bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*, \frac{d}{dt} \bar{\boldsymbol{\theta}}(t) \right\rangle \\ &= \left\langle \bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*, \mathbf{N}_{\boldsymbol{\theta}(t)} \frac{d}{dt} \boldsymbol{\theta}(t) \right\rangle \\ &= -\langle \mathbf{N}_{\boldsymbol{\theta}(t)} (\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*), \nabla \mathcal{L}(\boldsymbol{\theta}) \rangle \\ &= -\langle \mathbf{N}_{\boldsymbol{\theta}(t)} (\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*), \nabla^2 \mathcal{L}(\hat{\boldsymbol{\theta}}) (\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \rangle \end{aligned}$$

Notice that $\mathbf{N}_{\boldsymbol{\theta}(t)} (\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*)$ takes the form

$$\mathbf{N}_{\boldsymbol{\theta}(t)} (\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*) = \begin{bmatrix} \|\mathbf{v}_1\|_2^{-1} \bar{\mathbf{v}}_1 \\ \vdots \\ \|\mathbf{v}_{m^*}\|_2^{-1} \bar{\mathbf{v}}_{m^*} \\ \|\mathbf{w}_1\|_2^{-1} \bar{\mathbf{w}}_1 \\ \vdots \\ \|\mathbf{w}_{m^*}\|_2^{-1} \bar{\mathbf{w}}_{m^*} \end{bmatrix}$$

Thus, we are going to apply Theorem 4 with $N_{\min} = 1 - o(1) \frac{\delta_{\mathbb{P}}}{m^2}$, $N_{\max} = 1 + o(1) \frac{\delta_{\mathbb{P}}}{m^2}$, and $\alpha_i = \|\mathbf{v}_i\|_2^{-1} \geq \frac{m^2}{m^2 + o(1)\delta_{\mathbb{P}}}$, $\beta_i = \|\mathbf{w}_i\|_2^{-1} \leq \frac{m^2}{m^2 - o(1)\delta_{\mathbb{P}}}$. This leads to the condition that $\varepsilon \leq o\left(\frac{\mathcal{Q}^{*2}}{m^{*2}}\right)$ and $\frac{C_{S,0}}{C_{S,1}} \geq \frac{1.05}{\mathcal{Q}^{*2}}$. Under such condition, by Theorem 4, we have that

$$\frac{d}{dt} \|\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*\|_2^2 \leq -\left(1 - o(1) \frac{\delta_{\mathbb{P}}}{m^2}\right) \mathcal{Q}^* \kappa \|\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*\|_2^2$$

This shows that $\|\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*\|_2^2$ decreases monotonically. To find \mathcal{Q}^* , we notice that

$$\begin{aligned} \|(\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) (\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*)\|_2^2 &= \langle \bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*, (\mathbf{I} - \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top) (\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*) \rangle \\ &= \|\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*\|_2^2 - (\bar{\mathbf{v}}_i^\top (\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*))^2 \\ &= \|\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*\|_2^2 - (1 - \bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_i^*)^2 \\ &= \|\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*\|_2^2 - \frac{1}{4} \|\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*\|_2^4 \\ &= \left(1 - \frac{1}{4} \|\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*\|_2^2\right) \|\bar{\mathbf{v}}_i(t) - \bar{\mathbf{v}}_i^*\|_2^2 \end{aligned}$$

Similarly, we can obtain that

$$\|(\mathbf{I} - \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top) (\bar{\mathbf{w}}_i(t) - \bar{\mathbf{w}}_i^*)\|_2^2 = \left(1 - \frac{1}{4} \|\bar{\mathbf{w}}_i(t) - \bar{\mathbf{w}}_i^*\|_2^2\right) \|\bar{\mathbf{w}}_i(t) - \bar{\mathbf{w}}_i^*\|_2^2$$

This gives that $\mathcal{Q}^* = 1 - \frac{1}{4} \max_{i \in [m^*]} \max \left\{ \|\bar{\mathbf{v}}_i(0) - \bar{\mathbf{v}}_i^*\|_2^2, \|\bar{\mathbf{w}}_i(0) - \bar{\mathbf{w}}_i^*\|_2^2 \right\} = 1 - \mathcal{O}(\varepsilon)$. Thus, the condition that $\varepsilon \leq o\left(\frac{1}{m^{*2}}\right)$ and $C_{S,0} \geq 1.1C_{S,1}$ suffice. This gives us that

$$\frac{d}{dt} \|\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*\|_2^2 \leq -\left(1 - o(1) \frac{\delta_{\mathbb{P}}}{m^2}\right) (1 - \mathcal{O}(\varepsilon)) \kappa \|\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*\|_2^2 \leq \frac{\kappa}{2} \|\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*\|_2^2$$

Solving the ODE gives that

$$\|\bar{\boldsymbol{\theta}}(t) - \boldsymbol{\theta}^*\|_2^2 \leq \exp\left(-\frac{\kappa t}{2}\right) \|\bar{\boldsymbol{\theta}}(0) - \boldsymbol{\theta}^*\|_2^2$$

□

D Auxiliary Results

D.1 Hermite Polynomials

Lemma 19 (Restatement of Lemma 1). *Let $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$. For some multi-index $\mathbf{k} \in \mathbb{N}^n$, we define the multi-variate Hermite polynomial as*

$$He_{\mathbf{k}}(\mathbf{x}) = \prod_{i=1}^n He_{\mathbf{k}[i]}(\mathbf{x}[i])$$

Then we have that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)}[He_{\mathbf{k}}(\mathbf{x})] = \left(\prod_{i=1}^n \mathbf{k}[i]! \right) \sum_{\mathbf{M} \in \mathcal{S}} \prod_{i,j=1}^n \frac{\Sigma[i,j]^{\mathbf{M}[i,j]}}{\mathbf{M}[i,j]!}$$

where the set \mathcal{S} is defined by

$$\mathcal{S} = \left\{ \mathbf{M} \in \mathbb{N}^{n \times n} : \mathbf{M} = \mathbf{M}^\top, \sum_{j=1}^n \mathbf{M}[i,j] = \mathbf{k}[i], \mathbf{M}[i,i] = 0; \forall i \in [n], \right\}$$

Proof. Consider the generating function of Hermite polynomials

$$\exp\left(xt - \frac{t^2}{2}\right) = \sum_{k=0}^{\infty} \frac{He_k(x)}{k!} \cdot t^k$$

Let $x_1, \dots, x_n \sim \mathcal{N}(0, 1)$, then for all $\{\mathbf{t}_i\}_{i=1}^n$ we have

$$\begin{aligned} \exp\left(\sum_{i=1}^n \left(x_i t_i - \frac{t_i^2}{2}\right)\right) &= \prod_{i=1}^n \exp\left(x_i t_i - \frac{t_i^2}{2}\right) \\ &= \prod_{i=1}^n \left(\sum_{k=0}^{\infty} \frac{He_k(x_i)}{k!} \cdot t_i^k\right) \\ &= \sum_{\mathbf{k} \in \mathbb{N}^n} \left(\prod_{i=1}^n \frac{He_{\mathbf{k}[i]}(x_i)}{k[i]!}\right) \cdot \mathbf{t}^{\mathbf{k}} \end{aligned}$$

where $\mathbf{k} \in \mathbb{N}^n$ is the multi-index. On the other hand, if we write $x_i = \mathbf{u}_i^\top \mathbf{x}$ for some $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and \mathbf{u}_i satisfying $\|\mathbf{u}_i\|_2 = 1$, then we have

$$\begin{aligned} \mathbb{E}_{x_i \sim \mathcal{N}(0,1)} \left[\exp\left(\sum_{i=1}^n \left(x_i t_i - \frac{t_i^2}{2}\right)\right) \right] &= \mathbb{E}_{x_i \sim \mathcal{N}(0,1)} \left[\exp\left(\sum_{i=1}^n \left(\mathbf{u}_i^\top \mathbf{x} \cdot t_i - \frac{t_i^2}{2}\right)\right) \right] \\ &= \mathbb{E}_{x_i \sim \mathcal{N}(0,1)} \left[\exp\left(\mathbf{x}^\top \left(\sum_{i=1}^n t_i \mathbf{u}_i\right)\right) \right] \exp\left(-\frac{1}{2} \sum_{i=1}^n t_i^2\right) \end{aligned}$$

Since $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, we must have that $\mathbf{x}^\top (\sum_{i=1}^n t_i \mathbf{u}_i) \sim \mathcal{N}\left(\mathbf{0}, \|\sum_{i=1}^n t_i \mathbf{u}_i\|_2^2\right)$. By the moment-generating function of Gaussian random variable we have that

$$\mathbb{E}_{x_i \sim \mathcal{N}(0,1)} \left[\exp\left(\mathbf{x}^\top \left(\sum_{i=1}^n t_i \mathbf{u}_i\right)\right) \right] = \exp\left(\frac{1}{2} \left\| \sum_{i=1}^n t_i \mathbf{u}_i \right\|_2^2\right)$$

Thus, we have that

$$\mathbb{E}_{x_i \sim \mathcal{N}(0,1)} \left[\exp\left(\sum_{i=1}^n \left(x_i t_i - \frac{t_i^2}{2}\right)\right) \right] = \exp\left(\frac{1}{2} \left(\left\| \sum_{i=1}^n t_i \mathbf{u}_i \right\|_2^2 - \sum_{i=1}^n t_i^2 \right)\right) = \exp\left(\frac{1}{2} \sum_{i \neq j} \mathbf{u}_i^\top \mathbf{u}_j t_i t_j\right)$$

Applying Taylor's expansion gives

$$\exp\left(\frac{1}{2}\sum_{i \neq j} \mathbf{u}_i^\top \mathbf{u}_j t_i t_j\right) = \sum_{\ell=0}^{\infty} \frac{1}{\ell!} \left(\sum_{i < j} \mathbf{u}_i^\top \mathbf{u}_j t_i t_j\right)^\ell$$

Combining the results gives

$$\sum_{\mathbf{k} \in \mathbb{N}^d} \mathbb{E}_{x_i \sim \mathcal{N}(0,1)} \left[\prod_{i=1}^n He_{\mathbf{k}[i]}(x_i) \right] \prod_{i=1}^n \frac{t_i^{\mathbf{k}[i]}}{\mathbf{k}[i]!} = \sum_{\ell=0}^{\infty} \frac{1}{\ell!} \left(\sum_{i < j} \mathbf{u}_i^\top \mathbf{u}_j t_i t_j\right)^\ell$$

We intend to find out the coefficients of term $\prod_{i=1}^n t_i^{\mathbf{k}[i]}$ on the right-hand side. Notice that such term must only appears for term with ℓ satisfying $2\ell = \|\mathbf{k}\|_1$. By the multinomial theorem we have that

$$\left(\sum_{i \neq j} \mathbf{u}_i^\top \mathbf{u}_j t_i t_j\right)^\ell = \sum_{\mathbf{k}': \|\mathbf{k}'\|_1 = \ell} \ell! \cdot \prod_{i < j} \frac{(\mathbf{u}_i^\top \mathbf{u}_j)^{\mathbf{k}'[i,j]}}{\mathbf{k}'[i,j]!} t_i^{\mathbf{k}'[i,j]} t_j^{\mathbf{k}'[i,j]}$$

Therefore, we must have that

$$\mathbb{E}_{x_i \sim \mathcal{N}(0,1)} \left[\prod_{i=1}^n He_{\mathbf{k}[i]}(x_i) \right] = \left(\prod_{i=1}^n \mathbf{k}[i]! \right) \sum_{\mathbf{M} \in \mathcal{S}} \prod_{i < j} \frac{\mathbb{E}[x_i x_j]^{\mathbf{M}[i,j]}}{\mathbf{M}[i,j]!}$$

where \mathcal{S} is given by

$$\mathcal{S} = \left\{ \mathbf{M} \in \mathbb{N}^{n \times n} : \text{Tr}(\mathbf{M}) = 0; \forall i \in [n], \sum_{j=1}^n \mathbf{M}[i,j] = \mathbf{k}[i] \right\}$$

Using vector notations, we have that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)} [He_{\mathbf{k}}(\mathbf{x})] = \left(\prod_{i=1}^n \mathbf{k}[i]! \right) \sum_{\mathbf{M} \in \mathcal{S}} \prod_{i < j} \frac{\Sigma[i,j]^{\mathbf{M}[i,j]}}{\mathbf{M}[i,j]!}$$

□

Lemma 20 (Parseval's Identity). *Let $f(x)$ be given such that $\mathbb{E}_{x \sim \mathcal{N}(0,1)} [f^{(\ell)}(x)^2] \leq B$, and let c_k be the k th Hermite coefficient of $f(x)$. Then we have that*

$$c_{k+\ell}^2 \leq B \cdot k!; \quad \forall k \geq 0$$

Proof. Taking the Hermite expansion of $f^{(\ell)}$ gives

$$f^{(\ell)} = \sum_{k=0}^{\infty} \frac{c'_k}{k!} He_k(x); \quad c'_k = \mathbb{E}_{x \sim \mathcal{N}(0,1)} [f^{(\ell+k)}(x)] = c_{k+\ell}$$

Therefore, we have that

$$\mathbb{E}_{x \sim \mathcal{N}(0,1)} [f^{(\ell)}(x)^2] = \sum_{k=0}^{\infty} \frac{c'_k{}^2}{k!} = \sum_{k=0}^{\infty} \frac{c_{k+\ell}^2}{k!} \leq B$$

This implies that $c_{k+\ell}^2 \leq B \cdot k!$ since $c_{k+\ell}^2 \geq 0$ for all k . □

Lemma 21. *Let $\mathbf{v}_1, \mathbf{v}_2, \mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$ be vectors of unit norm such that*

$$\max\{|\mathbf{v}_1^\top \mathbf{w}_1|, |\mathbf{v}_1^\top \mathbf{w}_2|, |\mathbf{v}_2^\top \mathbf{w}_1|, |\mathbf{v}_2^\top \mathbf{w}_2|\} \leq \delta_r$$

with some $\delta_r \in (0, 1)$. Let $\{h_k\}_{k=0}^{\infty}, \{h'_k\}_{k=0}^{\infty}$ be two sequences of real numbers such that

$$\sum_{k=0}^{\infty} \frac{h_{k+a} h'_{k+b}}{k!} \leq \mathcal{O}(1); \quad \forall a + b \leq 6, a, b \in \mathbb{N} \cup \{0\} \quad (25)$$

Then we have that

$$\begin{aligned} & \sum_{k,\ell=0}^{\infty} \frac{h_k h'_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_3(\mathbf{w}_1^\top \mathbf{x}) He_3(\mathbf{w}_2^\top \mathbf{x})] \\ &= 6 \sum_{k=0}^{\infty} \frac{h_k h'_k}{k!} (\mathbf{v}_1^\top \mathbf{v}_2)^k (\mathbf{w}_1^\top \mathbf{w}_2)^3 \pm \mathcal{O}(\delta_r^2 (\mathbf{w}_1^\top \mathbf{w}_2)^2 + \delta_r^4) \end{aligned} \quad (26)$$

$$\begin{aligned} & \sum_{k,\ell=0}^{\infty} \frac{h_k h'_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_2(\mathbf{w}_2^\top \mathbf{x})] \\ &= 2 \sum_{k=0}^{\infty} \frac{h_k h'_k}{k!} (\mathbf{v}_1^\top \mathbf{v}_2)^k (\mathbf{w}_1^\top \mathbf{w}_2)^2 \pm \mathcal{O}(\delta_r^2 \cdot |\mathbf{w}_1^\top \mathbf{w}_2| + \delta_r^4) \end{aligned} \quad (27)$$

$$\begin{aligned} & \sum_{k,\ell=0}^{\infty} \frac{h_k h'_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_3(\mathbf{w}_2^\top \mathbf{x})] \\ &= 6 \sum_{k=0}^{\infty} \frac{h_{k+1} h'_k}{k!} (\mathbf{v}_1^\top \mathbf{v}_2)^k (\mathbf{w}_1^\top \mathbf{w}_2)^2 \mathbf{v}_1^\top \mathbf{w}_2 \\ & \quad + 6 \sum_{k=0}^{\infty} \frac{h_k h'_{k+1}}{k!} (\mathbf{v}_1^\top \mathbf{v}_2)^k (\mathbf{w}_1^\top \mathbf{w}_2)^2 \mathbf{v}_2^\top \mathbf{w}_2 \pm \mathcal{O}(\delta_r^3). \end{aligned} \quad (28)$$

Proof. The general idea of proving this lemma is to use Lemma 1. In particular, in our case we have that

$$\Sigma[i, j] = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{w}_1, \mathbf{w}_2]^\top [\mathbf{v}_1, \mathbf{v}_2, \mathbf{w}_1, \mathbf{w}_2] \in \mathbb{R}^{4 \times 4}$$

For the convenience of the analysis, we denote

$$\gamma_1 = \mathbf{v}_1^\top \mathbf{v}_2; \gamma_2 = \mathbf{w}_1^\top \mathbf{w}_2; \zeta_{ij} = \mathbf{v}_i^\top \mathbf{w}_j$$

By the assumption, we have that $|\zeta_{ij}| \leq \delta_r$.

Proof of (26). We start from the first equation. In particular we need to study

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_3(\mathbf{w}_1^\top \mathbf{x}) He_3(\mathbf{w}_2^\top \mathbf{x})]$$

By Lemma 1, we need to enumerate all \mathcal{S} , which is essentially the symmetric matrix \mathbf{M} with zero diagonal and non-negative entries whose row-sum equal to the vector $[k, \ell, 3, 3]$. This is equivalent to construct a weighted graph with four nodes and node degree $[k, \ell, 3, 3]$. Thus, it suffice to consider cases $k = \ell$, $|k - \ell| = 2$, $|k - \ell| = 4$, and $|k - \ell| = 6$. Due to symmetry between k and ℓ , we will first study the case $k \geq \ell$ and switch the indices to obtain all cases.

Case $k = \ell$. The node \mathbf{w}_1 and \mathbf{w}_2 has a total degree of 6, therefore, the pair of node \mathbf{v}_1 and \mathbf{v}_2 can have outgoing degree at most 6. Thus, the condition can be broken down into $\mathbf{M}[1, 2] \in \{k, k-1, k-2, k-3\}$. When $\mathbf{M}[1, 2] = k$, we have that $\mathbf{M}[3, 4] = 3$, and all other edges 0. When $\mathbf{M}[1, 2] = k-1$, we have that $\mathbf{M}[3, 4] = 2$ and either $(\mathbf{M}[1, 3], \mathbf{M}[2, 4]) = (1, 1)$ or $(\mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (1, 1)$. When $\mathbf{M}[1, 2] = k-2$, we have that $\mathbf{M}[3, 4] = 1$. Here we can have $(\mathbf{M}[1, 3], \mathbf{M}[2, 4]) = (2, 2)$, $(\mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (2, 2)$, or $(\mathbf{M}[1, 3], \mathbf{M}[2, 4], \mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (1, 1, 1, 1)$. When $\mathbf{M}[1, 2] = k-3$, then $\mathbf{M}[3, 4] = 0$. Thus we have $(\mathbf{M}[1, 3], \mathbf{M}[2, 4]) = (3, 3)$, $(\mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (3, 3)$ or $(\mathbf{M}[1, 3], \mathbf{M}[2, 4], \mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (1, 1, 2, 2)$, or $(\mathbf{M}[1, 3], \mathbf{M}[2, 4], \mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (2, 2, 1, 1)$. Plugging the possibilities into Lemma 1 gives that, under the case $k = \ell$, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_3(\mathbf{w}_1^\top \mathbf{x}) He_3(\mathbf{w}_2^\top \mathbf{x})] \\ &= 6k! \gamma_1^k \gamma_2^3 + 18P(k, 1)k! \gamma_1^{k-1} \gamma_2^2 (\zeta_{11} \zeta_{22} + \zeta_{12} \zeta_{21}) \\ & \quad + 9P(k, 2)k! \gamma_1^{k-2} \gamma_2 (\zeta_{11}^2 \zeta_{22}^2 + 4\zeta_{11} \zeta_{12} \zeta_{21} \zeta_{22} + \zeta_{12}^2 \zeta_{21}^2) \\ & \quad + P(k, 3)k! \gamma_1^{k-3} (\zeta_{11}^3 \zeta_{22}^3 + \zeta_{12}^3 \zeta_{21}^3 + 9\zeta_{11}^2 \zeta_{12} \zeta_{21} \zeta_{22}^2 + 9\zeta_{11} \zeta_{12}^2 \zeta_{21}^2 \zeta_{22}) \\ &= 6k! \gamma_1^k \gamma_2^3 \pm \mathcal{O}(\delta_r^2 \gamma_2^2) \cdot k! P(\ell, 1) \pm \mathcal{O}(\delta_r^4) \cdot k! (P(\ell, 2) + P(\ell, 3)) \end{aligned}$$

where $P(k, a) = \frac{k!}{(k-a)!}$ if $k \geq a$ and $P(k, a) = 0$ if $k < a$ represents the permutation number.

Case $k = \ell + 2$. In this case we have that $\mathbf{M}[1, 2] \in \{\ell, \ell - 1, \ell - 2\}$. If $\mathbf{M}[1, 2] = \ell$, then $\mathbf{M}[3, 4] = 2$. Here we have that $(\mathbf{M}[1, 3], \mathbf{M}[1, 4]) = (1, 1)$. If $\mathbf{M}[1, 2] = \ell - 1$, then $\mathbf{M}[3, 4] = 1$. Here we have $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 4]) = (2, 1, 1)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (1, 2, 1)$. If $\mathbf{M}[1, 2] = \ell - 2$, then $\mathbf{M}[3, 4] = 0$. Here we have $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 4]) = (3, 1, 2)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (1, 3, 2)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 3], \mathbf{M}[2, 4]) = (2, 2, 1, 1)$. Gathering all possibilities gives that, under the case $k = \ell + 2$, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_3(\mathbf{w}_1^\top \mathbf{x}) He_3(\mathbf{w}_2^\top \mathbf{x})] \\ &= 18k! \gamma_1^\ell \gamma_2^2 \zeta_{11} \zeta_{12} + 18P(\ell, 1)k! \gamma_1^{\ell-1} \gamma_2 \zeta_{11} \zeta_{12} (\zeta_{11} \zeta_{22} + \zeta_{12} \zeta_{21}) \\ &\quad + 3P(\ell, 2)k! \gamma_1^{\ell-2} \zeta_{11} \zeta_{12} (\zeta_{11}^2 \zeta_{22}^2 + \zeta_{12}^2 \zeta_{21}^2) \\ &\quad + 9P(\ell, 2)k! \gamma_1^{\ell-2} \zeta_{11}^2 \zeta_{12}^2 \zeta_{21} \zeta_{22} \\ &= \pm \mathcal{O}(\delta_r^2 \gamma_2^2) \cdot k! \pm \mathcal{O}(\delta_r^4) \cdot k! (P(\ell, 1) + P(\ell, 2)) \end{aligned}$$

Case $k = \ell + 4$. In this case we have that $\mathbf{M}[1, 2] \in \{\ell, \ell - 1\}$. If $\mathbf{M}[1, 2] = \ell$, then $\mathbf{M}[3, 4] = 1$. Here we have that $(\mathbf{M}[1, 3], \mathbf{M}[1, 4]) = (2, 2)$. If $\mathbf{M}[1, 2] = \ell - 1$, then $\mathbf{M}[3, 4] = 0$. Here we have that $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 4]) = (3, 2, 1)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (2, 3, 1)$. Gathering all possibilities gives that, under the case $k = \ell + 4$, we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_3(\mathbf{w}_1^\top \mathbf{x}) He_3(\mathbf{w}_2^\top \mathbf{x})] \\ &= 9k! \gamma_1^\ell \gamma_2 \zeta_{11}^2 \zeta_{12}^2 + 3P(\ell, 1)k! \gamma_1^{\ell-1} \zeta_{11}^2 \zeta_{12}^2 (\zeta_{11} \zeta_{22} + \zeta_{21} \zeta_{12}) \\ &= \pm \mathcal{O}(\delta_r^4) \cdot k! (1 + P(\ell, 1)) \end{aligned}$$

Case $k = \ell + 6$. In this case we have that $\mathbf{M}[1, 2] = \ell$, $\mathbf{M}[3, 4] = 0$ and $\mathbf{M}[1, 3] = \mathbf{M}[1, 4] = 3$. Thus, if $k = \ell + 6$, we have that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_3(\mathbf{w}_1^\top \mathbf{x}) He_3(\mathbf{w}_2^\top \mathbf{x})] = k! \gamma_1^\ell \gamma_2 \zeta_{11}^3 \zeta_{12}^3 = \pm \mathcal{O}(\delta_r^4) \cdot k!$$

Putting things together, we have that

$$\begin{aligned} & \sum_{k, \ell=0}^{\infty} \frac{h_k h'_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_3(\mathbf{w}_1^\top \mathbf{x}) He_3(\mathbf{w}_2^\top \mathbf{x})] \\ &= 6 \sum_{k=0}^{\infty} \frac{h_k h'_k}{k!} \gamma_1^k \gamma_2^3 \pm \mathcal{O}(\delta_r^2) \sum_{k=0}^{\infty} \frac{h_k h'_k}{k!} \pm \mathcal{O}(\delta_r^2 \gamma_2^2) \sum_{k=0}^{\infty} \frac{h_{k+1} h'_{k+1}}{k!} \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{c_{k+2} c'_{k+2}}{k!} \\ &\quad \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_{k+3} h'_{k+3}}{k!} \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_k h'_{k+2} + h'_k h_{k+2}}{k!} \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_{k+1} h'_{k+3} + h'_{k+1} h_{k+3}}{k!} \\ &\quad \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_{k+2} h'_{k+4} + h'_{k+2} h_{k+4}}{k!} \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_k h'_{k+4} + c'_k c_{k+4}}{k!} \\ &\quad \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_{k+1} h'_{k+5} + h'_{k+1} h_{k+5}}{k!} \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_k h'_{k+6} + h'_k h_{k+6}}{k!} \\ &= 6 \sum_{k=0}^{\infty} \frac{h_k h'_k}{k!} \gamma_1^k \gamma_2^3 \pm \mathcal{O}(\delta_r^2 \gamma_2^2 + \delta_r^4) \end{aligned}$$

Proof of (27). Similar to before, the combination of k, ℓ can be $k = \ell$, $|k - \ell| = 2$, or $|k - \ell| = 4$ due to the total degree of \mathbf{w}_1 and \mathbf{w}_2 is 4. We study the case $k \geq \ell$.

Case $k = \ell$. In this case, we have $\mathbf{M}[1, 2] \in \{\ell, \ell - 1, \ell - 2\}$. If $\mathbf{M}[1, 2] = \ell$, then $\mathbf{M}[3, 4] = 2$ and all other edges are 0. If $\mathbf{M}[1, 2] = \ell - 1$, then $\mathbf{M}[3, 4] = 1$, and either $(\mathbf{M}[1, 3], \mathbf{M}[2, 4]) = (1, 1)$ or $(\mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (1, 1)$. If $\mathbf{M}[1, 2] = \ell - 2$, then $\mathbf{M}[3, 4] = 0$. Here we can have $(\mathbf{M}[1, 3], \mathbf{M}[2, 4]) = (2, 2)$ or $(\mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (2, 2)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 3], \mathbf{M}[2, 4]) = (1, 1, 1, 1)$. Gathering all possibilities gives that, under the case $k = \ell$, we

have

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_2(\mathbf{w}_2^\top \mathbf{x}) \right] \\
 &= 2k! \gamma_1^k \gamma_2^2 + 4P(\ell, 1) k! \gamma_1^{k-1} \gamma_2 (\zeta_{11} \zeta_{22} + \zeta_{12} \zeta_{21}) \\
 &\quad + P(\ell, 2) k! \gamma_1^{k-2} (\zeta_{11}^2 \zeta_{22}^2 + \zeta_{12}^2 \zeta_{21}^2 + 4\zeta_{11} \zeta_{12} \zeta_{21} \zeta_{22}) \\
 &= 2k! \gamma_1^k \gamma_2^2 \pm \mathcal{O}(\delta_r^2 \gamma_2) \cdot k! P(\ell, 1) \pm \mathcal{O}(\delta_r^4) \cdot k! P(\ell, 2)
 \end{aligned}$$

In the case $\zeta_{22} = 0$, we denote $\hat{\zeta} = \max\{|\zeta_{21}|, |\zeta_{12}|\}$ we have that

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_2(\mathbf{w}_2^\top \mathbf{x}) \right] \\
 &= 2k! \gamma_1^k \gamma_2^2 + 4P(\ell, 1) k! \gamma_1^{k-1} \gamma_2 \zeta_{12} \zeta_{21} + P(\ell, 2) k! \gamma_1^{k-2} \zeta_{12}^2 \zeta_{21}^2 \\
 &= 2k! \gamma_1^k \gamma_2^2 \pm \mathcal{O}(\hat{\zeta}^2 \gamma_2) \cdot k! P(\ell, 1) \pm \mathcal{O}(\hat{\zeta}^4) \cdot k! P(\ell, 2)
 \end{aligned}$$

Case $k = \ell + 2$. In this case we have $\mathbf{M}[1, 2] \in \{\ell, \ell - 1\}$. If $\mathbf{M}[1, 2] = \ell$, then $\mathbf{M}[3, 4] = 1$ and $(\mathbf{M}[1, 3], \mathbf{M}[1, 4]) = (1, 1)$. If $\mathbf{M}[1, 2] = \ell - 1$, then $\mathbf{M}[3, 4] = 0$ and either $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], (\mathbf{M})[2, 3]) = (1, 2, 1)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], (\mathbf{M})[2, 4]) = (2, 1, 1)$. Gathering all possibilities gives that, under the case $k = \ell + 2$, we have

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_2(\mathbf{w}_2^\top \mathbf{x}) \right] \\
 &= 4k! \gamma_1^\ell \gamma_2 \zeta_{11} \zeta_{12} + 2P(\ell, 1) k! \gamma_1^{\ell-1} \zeta_{11} \zeta_{12} (\zeta_{11} \zeta_{22} + \zeta_{12} \zeta_{21}) \\
 &= \pm \mathcal{O}(\delta_r^2 \gamma_2) \cdot k! \pm \mathcal{O}(\delta_r^4) \cdot k! P(\ell, 1)
 \end{aligned}$$

In the case where $\zeta_{22} = 0$, we have that

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_2(\mathbf{w}_2^\top \mathbf{x}) \right] \\
 &= 4k! \gamma_1^\ell \gamma_2 \zeta_{11} \zeta_{12} + 2P(\ell, 1) k! \gamma_1^{\ell-1} \zeta_{11} \zeta_{12}^2 \zeta_{21} \\
 &= \pm \mathcal{O}(\delta_r \hat{\zeta} \gamma_2) \cdot k! \pm \mathcal{O}(\hat{\zeta}^3) \cdot k! P(\ell, 1)
 \end{aligned}$$

Case $k = \ell + 4$. In this case we can only have $\mathbf{M}[1, 2] = \ell$, $\mathbf{M}[3, 4] = 0$, and $\mathbf{M}[1, 3] = \mathbf{M}[1, 4] = 2$. Thus if $k = \ell + 4$, we have

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_2(\mathbf{w}_2^\top \mathbf{x}) \right] = k! \gamma_1^\ell \zeta_{11}^2 \zeta_{12}^2 = \pm \mathcal{O}(\delta_r^4) \cdot k!$$

In the case where $\zeta_{22} = 0$, we have that

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_2(\mathbf{w}_2^\top \mathbf{x}) \right] = k! \gamma_1^\ell \zeta_{11}^2 \zeta_{12}^2 = \pm \mathcal{O}(\delta_r^2 \hat{\zeta}^2) \cdot k!$$

Putting things together, we have that

$$\begin{aligned}
 & \sum_{k, \ell=0}^{\infty} \frac{h_k h'_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_2(\mathbf{w}_1^\top \mathbf{x}) He_2(\mathbf{w}_2^\top \mathbf{x}) \right] \\
 &= 2 \sum_{k=0}^{\infty} \frac{h_k h_{k'}}{k!} \gamma_1^k \gamma_2^2 \pm \mathcal{O}(\delta_r^2 \gamma_2) \sum_{k=0}^{\infty} \frac{h_{k+1} h'_{k+1}}{k!} \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_{k+2} h_{k+2}^2}{k!} \\
 &\quad \pm \mathcal{O}(\delta_r^2 \gamma_2) \sum_{k=0}^{\infty} \frac{h_k h_{k+2} + h'_k h_{k+2}}{k!} \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_{k+1} h_{k+3} + h'_{k+1} h'_{k+3}}{k!} \\
 &\quad \pm \mathcal{O}(\delta_r^4) \sum_{k=0}^{\infty} \frac{h_k h_{k+4} + h'_k h_{k+4}}{k!} \\
 &= 2 \sum_{k=0}^{\infty} \frac{h_k h_{k'}}{k!} \gamma_1^k \gamma_2^2 \pm \mathcal{O}(\delta_r^2 \gamma_2 + \delta_r^4)
 \end{aligned}$$

In the case where $\zeta_{22} = 0$, we have that

$$\begin{aligned}
 & \sum_{k,\ell=0}^{\infty} \frac{h_k h'_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [H e_k(\mathbf{v}_1^\top \mathbf{x}) H e_\ell(\mathbf{v}_2^\top \mathbf{x}) H e_2(\mathbf{w}_1^\top \mathbf{x}) H e_2(\mathbf{w}_2^\top \mathbf{x})] \\
 &= 2 \sum_{k=0}^{\infty} \frac{h_k h'_k}{k!} \gamma_1^k \gamma_2^2 \pm \mathcal{O}(\hat{\zeta}^2 \gamma_2) \sum_{k=0}^{\infty} \frac{h_{k+1} h'_{k+1}}{k!} \pm \mathcal{O}(\hat{\zeta}^4) \sum_{k=0}^{\infty} \frac{h_{k+2} h'_{k+2}}{k!} \\
 & \quad \pm \mathcal{O}(\delta_r \hat{\zeta} \gamma_2) \sum_{k=0}^{\infty} \frac{h_k h'_{k+2}}{k!} \pm \mathcal{O}(\hat{\zeta}^3) \sum_{k=0}^{\infty} \frac{h_{k+1} h'_{k+3}}{k!} \\
 & \quad \pm \mathcal{O}(\delta_r^2 \hat{\zeta}^2) \sum_{k=0}^{\infty} \frac{h_k h'_{k+4}}{k!} \\
 &= 2 \sum_{k=0}^{\infty} \frac{h_k h'_k}{k!} \gamma_1^k \gamma_2^2 \pm \mathcal{O}(\delta_r \hat{\zeta} \gamma_2 + \delta_r \hat{\zeta}^2)
 \end{aligned}$$

Proof of (28). We notice that in this case k, ℓ must satisfy $|k - \ell| \in \{1, 3, 5\}$. Similar to before, we assume that $k \geq \ell$.

Case $k = \ell + 1$. In this case $\mathbf{M}[1, 2] \in \{\ell, \ell - 1, \ell - 2\}$. If $\mathbf{M}[1, 2] = \ell$, then $\mathbf{M}[3, 4] = 2$ and $\mathbf{M}[1, 4] = 1$. If $\mathbf{M}[1, 2] = \ell - 1$, then $\mathbf{M}[3, 4] = 1$, and either $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 4]) = (1, 1, 1)$ or $(\mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (2, 1)$. If $\mathbf{M}[1, 2] = \ell - 2$, then $\mathbf{M}[3, 4] = 0$. Here we have $(\mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (3, 2)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 3], \mathbf{M}[2, 4]) = (1, 2, 1, 1)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 4]) = (2, 1, 2)$. Gathering all possibilities gives that, under the case $k = \ell + 1$, we have

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [H e_k(\mathbf{v}_1^\top \mathbf{x}) H e_\ell(\mathbf{v}_2^\top \mathbf{x}) H e_2(\mathbf{w}_1^\top \mathbf{x}) H e_3(\mathbf{w}_2^\top \mathbf{x})] \\
 &= 6k! \gamma_1^\ell \gamma_2^2 \zeta_{12} + 12P(\ell, 1) k! \gamma_1^{\ell-1} \gamma_2 \zeta_{12} (2\zeta_{11} \zeta_{22} + \zeta_{12} \zeta_{21}) \\
 & \quad + P(\ell, 2) k! \gamma_1^{\ell-2} \zeta_{12} (\zeta_{12}^2 \zeta_{21}^2 + 3\zeta_{11}^2 \zeta_{22}^2 + 6\zeta_{11} \zeta_{12} \zeta_{21} \zeta_{22}) \\
 &= 6k! \gamma_1^\ell \gamma_2^2 \zeta_{12} \pm \mathcal{O}(\delta_r^3) (P(\ell, 1) + P(\ell, 2))
 \end{aligned}$$

Case $k = \ell + 3$. In this case we have $\mathbf{M}[1, 2] \in \{\ell, \ell - 1\}$. If $\mathbf{M}[1, 2] = \ell$, then $\mathbf{M}[3, 4] = 1$, and $(\mathbf{M}[1, 3], \mathbf{M}[1, 4]) = (1, 2)$. If $\mathbf{M}[1, 2] = \ell - 1$, then $\mathbf{M}[3, 4] = 0$, and we have $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 3]) = (1, 3, 1)$ or $(\mathbf{M}[1, 3], \mathbf{M}[1, 4], \mathbf{M}[2, 4]) = (2, 2, 1)$. Gathering all possibilities gives that, under the case $k = \ell + 3$, we have

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [H e_k(\mathbf{v}_1^\top \mathbf{x}) H e_\ell(\mathbf{v}_2^\top \mathbf{x}) H e_2(\mathbf{w}_1^\top \mathbf{x}) H e_3(\mathbf{w}_2^\top \mathbf{x})] \\
 &= 6k! \gamma_1^\ell \gamma_2 \zeta_{11} \zeta_{12}^2 + 6P(\ell, 1) k! \gamma_1^{\ell-1} \zeta_{11} \zeta_{12}^2 (3\zeta_{21} + 2\zeta_{22}) \\
 &= \pm \mathcal{O}(\delta_r^3) (1 + P(\ell, 1))
 \end{aligned}$$

Case $k = \ell + 5$. In this case we must have that $\mathbf{M}[1, 2] = \ell$, $\mathbf{M}[3, 4] = 0$, and $\mathbf{M}[1, 3] = 2$, $\mathbf{M}[1, 4] = 3$. Thus

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [H e_k(\mathbf{v}_1^\top \mathbf{x}) H e_\ell(\mathbf{v}_2^\top \mathbf{x}) H e_2(\mathbf{w}_1^\top \mathbf{x}) H e_3(\mathbf{w}_2^\top \mathbf{x})] = k! \gamma_1^\ell \zeta_{11}^2 \zeta_{12}^3 = \pm \mathcal{O}(\delta_r^5)$$

Putting things together gives

$$\begin{aligned}
 & \sum_{k,\ell=0}^{\infty} \frac{h_k h'_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [H e_k(\mathbf{v}_1^\top \mathbf{x}) H e_\ell(\mathbf{v}_2^\top \mathbf{x}) H e_2(\mathbf{w}_1^\top \mathbf{x}) H e_3(\mathbf{w}_2^\top \mathbf{x})] \\
 &= 6 \sum_{k=0}^{\infty} \frac{h_{k+1} h'_k}{k!} \gamma_1^k \gamma_2^2 \zeta_{12} + 6 \sum_{k=0}^{\infty} \frac{h_k h'_{k+1}}{k!} \gamma_1^k \gamma_2^2 \zeta_{22} \pm \mathcal{O}(\delta_r^3) \sum_{k=0}^{\infty} \frac{h_{k+1} h'_{k+2} + h'_{k+1} h_{k+2}}{k!} \\
 & \quad \pm \mathcal{O}(\delta_r^3) \sum_{k=0}^{\infty} \frac{h_{k+2} h'_{k+3} + h'_{k+2} h_{k+3}}{k!} \pm \mathcal{O}(\delta_r^3) \sum_{k=0}^{\infty} \frac{h_k h'_{k+3} + h'_k h_{k+3}}{k!} \\
 & \quad \pm \mathcal{O}(\delta_r^3) \sum_{k=0}^{\infty} \frac{h_{k+1} h'_{k+4} + h'_{k+1} h_{k+4}}{k!} \pm \mathcal{O}(\delta_r^3) \sum_{k=0}^{\infty} \frac{h_k h'_{k+5} + h'_k h_{k+5}}{k!} \\
 &= 6 \sum_{k=0}^{\infty} \frac{h_{k+1} h'_k}{k!} \gamma_1^k \gamma_2^2 \zeta_{12} + 6 \sum_{k=0}^{\infty} \frac{h_k h'_{k+1}}{k!} \gamma_1^k \gamma_2^2 \zeta_{22} \pm \mathcal{O}(\delta_r^3)
 \end{aligned}$$

□

Lemma 22. Let $h_i(\mathbf{x}) = \pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x})$ and $h_j(\mathbf{x}) = \pi(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^\top \mathbf{x})$. Suppose that

$$\max \{ |\bar{\mathbf{v}}_i^\top \bar{\mathbf{w}}_i|, |\bar{\mathbf{v}}_j^\top \bar{\mathbf{w}}_j|, |\bar{\mathbf{v}}_i^\top \bar{\mathbf{w}}_j|, |\bar{\mathbf{v}}_j^\top \bar{\mathbf{w}}_i| \} \leq \delta_r$$

then we have that

$$\mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x})^2] = 6 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \pm \mathcal{O}(\delta_r^2)$$

If it holds that $|\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j| \leq \delta_r$, then we have that

$$\mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x}) h_j(\mathbf{x})] = \pm \mathcal{O}(\delta_r^3);$$

If it holds that $|\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j| \leq \delta_r$, then we have that

$$\mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x}) h_j(\mathbf{x})] = 6 \mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x})^2] \pm \mathcal{O}(\delta_r)$$

Proof. Adopting the Hermite expansion, by Lemma 21 we have that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x}) h_j(\mathbf{x})] &= \sum_{k,\ell=0}^{\infty} \frac{c_k c_\ell}{k! \ell!} \mathbb{E}_{\mathbf{x}}[He_k(\bar{\mathbf{v}}_i^\top \mathbf{x}) He_\ell(\bar{\mathbf{v}}_j^\top \mathbf{x}) He_3(\bar{\mathbf{w}}_i^\top \mathbf{x}) (\bar{\mathbf{w}}_j^\top \mathbf{x})] \\ &= 6 (\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j)^3 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} (\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j)^k \pm \mathcal{O}(\delta_r^2 (\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j)^2 + \delta_r^4) \end{aligned}$$

Thus, in the case where $i \neq j$ and $|\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j| \leq \delta_r$, we have that

$$\mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x}) h_j(\mathbf{x})] = \pm \mathcal{O}(\delta_r^3)$$

On the other hand, if $i = j$, then we have that

$$\mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x})^2] = 6 \sum_{k=0}^{\infty} \frac{c_k^2}{k!} \pm \mathcal{O}(\delta_r^2)$$

If $\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j \geq 1 - \delta_r$ and $\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j \geq 1 - \delta$, then we have that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x}) h_j(\mathbf{x})] &= 6 (1 - 3\delta_r) \sum_{k=0}^{\infty} \frac{c_k^2}{k!} (\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j)^k \pm \mathcal{O}(\delta_r^2) \\ &= 6 (1 - 3\delta_r) \mathbb{E}_{\mathbf{x}}[\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi(\bar{\mathbf{v}}_j^\top \mathbf{x})] \pm \mathcal{O}(\delta_r^2) \end{aligned}$$

Applying the Taylor's expansion gives that

$$\begin{aligned} \mathbb{E}_{\mathbf{x}}[\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \pi(\bar{\mathbf{v}}_j^\top \mathbf{x})] &= \mathbb{E}_{z_1, z_2 \sim \mathcal{N}(0,1), \text{Cov}(z_1, z_2) = 1 - \delta_r} [\pi(z_1) \pi(z_2)] \\ &= \mathbb{E}_{z \sim \mathcal{N}(0,1)} \left[\pi(z)^2 \right] - \frac{1}{2} \text{Var}(\pi'(z_1)) \delta_r \pm \mathcal{O}(\delta_r) \\ &\geq \mathbb{E}_{z \sim \mathcal{N}(0,1)} \left[\pi(z)^2 \right] \pm \mathcal{O}(\delta_r) \end{aligned}$$

Thus, we have that

$$\mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x}) h_j(\mathbf{x})] = 6 \mathbb{E}_{\mathbf{x}}[h_i(\mathbf{x})^2] \pm \mathcal{O}(\delta_r)$$

□

Lemma 23. Let θ satisfy that $\|\bar{\mathbf{v}}_i^\top - \bar{\mathbf{v}}_i^*\|_2^2, \|\bar{\mathbf{w}}_i^\top - \bar{\mathbf{w}}_i^*\|_2^2 \leq \varepsilon$, $|\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j|, |\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j|, |\bar{\mathbf{v}}_i^\top \bar{\mathbf{v}}_j^*|, |\bar{\mathbf{w}}_i^\top \bar{\mathbf{w}}_j^*| \leq \varepsilon$ for all $i \neq j$, and $|\bar{\mathbf{v}}_i^\top \bar{\mathbf{w}}_j|, |\bar{\mathbf{w}}_i^\top \bar{\mathbf{v}}_j|, |\bar{\mathbf{v}}_i^\top \bar{\mathbf{w}}_j^*|, |\bar{\mathbf{w}}_i^\top \bar{\mathbf{v}}_j^*| \leq \varepsilon$ for all i, j . Then the following holds:

- $|\mathbb{E}_{\mathbf{x}}[(f(\theta, \mathbf{x}) - f^*(\mathbf{x})) \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x})]| \leq \mathcal{O}(m^* \varepsilon)$

- $\|\mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}^2(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x})]\|_2 \leq \mathcal{O}(m^* \varepsilon)$
- $\|\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}^2(\pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}))]\|_2 \leq \mathcal{O}(m^* \varepsilon)$
- $\|\mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}(\pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}))]\|_2 \leq \mathcal{O}(m^* \varepsilon)$
- $\|\mathbb{E}_{\mathbf{x}}[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}}(\pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}))]\|_2 \leq \mathcal{O}(m^* \varepsilon)$
- $\|\mathbb{E}_{\mathbf{x}}[\nabla_{\mathbf{x}}(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x})]\|_2 \leq \mathcal{O}(m^* \varepsilon)$

Proof. We first write out the gradient with respect to \mathbf{x}

$$\begin{aligned}
 \nabla_{\mathbf{x}} f(\boldsymbol{\theta}, \mathbf{x}) &= \sum_{j=1}^{m^*} \pi'(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^\top \mathbf{x}) \bar{\mathbf{v}}_j + \sum_{j=1}^{m^*} \pi(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_j^\top \mathbf{x}) \bar{\mathbf{w}}_j \\
 \nabla_{\mathbf{x}} f^*(\mathbf{x}) &= \sum_{j=1}^{m^*} \pi'(\bar{\mathbf{v}}_j^{*\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^{*\top} \mathbf{x}) \bar{\mathbf{v}}_j^* + \sum_{j=1}^{m^*} \pi(\bar{\mathbf{v}}_j^{*\top} \mathbf{x}) \sigma'(\bar{\mathbf{w}}_j^{*\top} \mathbf{x}) \bar{\mathbf{w}}_j^* \\
 \nabla_{\mathbf{x}}^2 f(\boldsymbol{\theta}, \mathbf{x}) &= \sum_{j=1}^{m^*} \pi'(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma'(\bar{\mathbf{w}}_j^\top \mathbf{x}) (\bar{\mathbf{v}}_j \bar{\mathbf{w}}_j^\top + \bar{\mathbf{w}}_j \bar{\mathbf{v}}_j^\top) \\
 &\quad + \sum_{j=1}^{m^*} \pi''(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^\top \mathbf{x}) \bar{\mathbf{v}}_j \bar{\mathbf{v}}_j^\top + \sum_{j=1}^{m^*} \pi(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma''(\bar{\mathbf{w}}_j^\top \mathbf{x}) \bar{\mathbf{w}}_j \bar{\mathbf{w}}_j^\top \\
 \nabla_{\mathbf{x}}^2 f^*(\mathbf{x}) &= \sum_{j=1}^{m^*} \pi'(\bar{\mathbf{v}}_j^{*\top} \mathbf{x}) \sigma'(\bar{\mathbf{w}}_j^{*\top} \mathbf{x}) (\bar{\mathbf{v}}_j^* \bar{\mathbf{w}}_j^{*\top} + \bar{\mathbf{w}}_j^* \bar{\mathbf{v}}_j^{*\top}) \\
 &\quad + \sum_{j=1}^{m^*} \pi''(\bar{\mathbf{v}}_j^{*\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^{*\top} \mathbf{x}) \bar{\mathbf{v}}_j^* \bar{\mathbf{v}}_j^{*\top} + \sum_{j=1}^{m^*} \pi(\bar{\mathbf{v}}_j^{*\top} \mathbf{x}) \sigma''(\bar{\mathbf{w}}_j^{*\top} \mathbf{x}) \bar{\mathbf{w}}_j^* \bar{\mathbf{w}}_j^{*\top}
 \end{aligned}$$

Moreover, we also have that

$$\begin{aligned}
 \nabla_{\mathbf{x}} \left(\pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right) &= \pi^{(a+1)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \bar{\mathbf{v}}_i + \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b+1)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \bar{\mathbf{w}}_i \\
 \nabla_{\mathbf{x}}^2 \left(\pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right) &= \pi^{(a+1)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b+1)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) (\bar{\mathbf{v}}_i \bar{\mathbf{w}}_i^\top + \bar{\mathbf{w}}_i \bar{\mathbf{v}}_i^\top) \\
 &\quad + \pi^{(a+2)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \bar{\mathbf{v}}_i \bar{\mathbf{v}}_i^\top \\
 &\quad + \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b+2)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \bar{\mathbf{w}}_i \bar{\mathbf{w}}_i^\top
 \end{aligned}$$

Therefore, for the last two bounds, we have

$$\begin{aligned}
 \left\| \mathbb{E}_{\mathbf{x}} \left[(f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \nabla_{\mathbf{x}} \left(\pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right) \right] \right\|_2 &= \mathcal{E}_{1,1} + \mathcal{E}_{2,1} \\
 \left\| \mathbb{E}_{\mathbf{x}} \left[\nabla_{\mathbf{x}} (f(\boldsymbol{\theta}, \mathbf{x}) - f^*(\mathbf{x})) \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right\|_2 &= \mathcal{E}_{1,2} + \mathcal{E}_{2,2}
 \end{aligned}$$

where $\mathcal{E}_{1,1}$ and $\mathcal{E}_{1,2}$ are a summation of two terms in the form

$$\left| \mathbb{E}_{\mathbf{x}} \left[\left(\pi(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^\top \mathbf{x}) - \pi(\bar{\mathbf{v}}_i^{*\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^{*\top} \mathbf{x}) \right) \pi^{(a')}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b')}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right|$$

Thus, by Taylor expansion, we obtain that $\mathcal{E}_{1,1}, \mathcal{E}_{1,2} \leq \mathcal{O}(\varepsilon)$. Moreover, $\mathcal{E}_{1,2}$ and $\mathcal{E}_{2,2}$ are a summation of $4m^*$ terms of the form

$$\begin{aligned}
 &\left| \mathbb{E}_{\mathbf{x}} \left[\pi(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^\top \mathbf{x}) \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b')}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right| \\
 &\left| \mathbb{E}_{\mathbf{x}} \left[\pi(\bar{\mathbf{v}}_j^{*\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^{*\top} \mathbf{x}) \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b')}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right|
 \end{aligned}$$

By Lemma 1, we have that $\mathcal{E}_{2,1}, \mathcal{E}_{2,2} \leq \mathcal{O}(m^* \varepsilon)$. This gives the last two property. For the rest of the property, we can apply similar strategy to decompose the objective in terms of

$$\left| \mathbb{E}_{\mathbf{x}} \left[\left(\pi^{(a_0)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b_0)}(\bar{\mathbf{w}}_i^\top \mathbf{x}) - \pi(\bar{\mathbf{v}}_i^{*\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_i^{*\top} \mathbf{x}) \right) \pi^{(a')}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b')}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right|$$

which can be upper bounded by Taylor expansion, and

$$\begin{aligned} & \left| \mathbb{E}_{\mathbf{x}} \left[\pi(\bar{\mathbf{v}}_j^\top \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^\top \mathbf{x}) \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b')}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right| \\ & \left| \mathbb{E}_{\mathbf{x}} \left[\pi(\bar{\mathbf{v}}_j^{*\top} \mathbf{x}) \sigma(\bar{\mathbf{w}}_j^{*\top} \mathbf{x}) \pi^{(a)}(\bar{\mathbf{v}}_i^\top \mathbf{x}) \sigma^{(b')}(\bar{\mathbf{w}}_i^\top \mathbf{x}) \right] \right| \end{aligned}$$

which can be upper bounded by lemma 25. Since there are in total $\mathcal{O}(m^*)$ terms for each quantity, we can conclude the desired result. \square

Lemma 24. *Let $\mathbf{v}_1, \mathbf{v}_2, \mathbf{w}_1, \mathbf{w}_2$ be unit vectors satisfying that any two of the four have an inner product with magnitude less than ε . Then for $a_1, a_2, b_1, b_2 > 0$ with $b_1 + b_2 \leq 3$, the following holds*

- $\left| \mathbb{E}_{\mathbf{x}} \left[\pi^{(a_1)}(\mathbf{v}_1^\top \mathbf{x}) \sigma^{(b_1)}(\mathbf{w}_1^\top \mathbf{x}) \pi^{(a_2)}(\mathbf{v}_2^\top \mathbf{x}) \sigma^{(b_2)}(\mathbf{w}_2^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon)$
- $\left\| \mathbb{E}_{\mathbf{x}} \left[\nabla_{\mathbf{x}}^2 \left(\pi^{(a_1)}(\mathbf{v}_1^\top \mathbf{x}) \sigma^{(b_1)}(\mathbf{w}_1^\top \mathbf{x}) \right) \pi^{(a_2)}(\mathbf{v}_2^\top \mathbf{x}) \sigma^{(b_2)}(\mathbf{w}_2^\top \mathbf{x}) \right] \right\|_2 \leq \mathcal{O}(\varepsilon)$
- $\left\| \mathbb{E}_{\mathbf{x}} \left[\pi^{(a_1)}(\mathbf{v}_1^\top \mathbf{x}) \sigma^{(b_1)}(\mathbf{w}_1^\top \mathbf{x}) \nabla_{\mathbf{x}}^2 \left(\pi^{(a_2)}(\mathbf{v}_2^\top \mathbf{x}) \sigma^{(b_2)}(\mathbf{w}_2^\top \mathbf{x}) \right) \right] \right\|_2 \leq \mathcal{O}(\varepsilon)$
- $\left\| \mathbb{E}_{\mathbf{x}} \left[\nabla_{\mathbf{x}} \left(\pi^{(a_1)}(\mathbf{v}_1^\top \mathbf{x}) \sigma^{(b_1)}(\mathbf{w}_1^\top \mathbf{x}) \right) \nabla_{\mathbf{x}} \left(\pi^{(a_2)}(\mathbf{v}_2^\top \mathbf{x}) \sigma^{(b_2)}(\mathbf{w}_2^\top \mathbf{x}) \right) \right] \right\|_2 \leq \mathcal{O}(\varepsilon)$

Proof. The first quantity is directly bounded by applying Lemma 25. For the rest, we write out the form of the gradients with respect to \mathbf{x} as

$$\begin{aligned} \nabla_{\mathbf{x}} \left(\pi^{(a)}(\mathbf{v}^\top \mathbf{x}) \sigma^{(b)}(\mathbf{w}^\top \mathbf{x}) \right) &= \pi^{(a+1)}(\mathbf{v}^\top \mathbf{x}) \sigma^{(b)}(\mathbf{w}^\top \mathbf{x}) \mathbf{v} + \pi^{(a)}(\mathbf{v}^\top \mathbf{x}) \sigma^{(b+1)}(\mathbf{w}^\top \mathbf{x}) \mathbf{w} \\ \nabla_{\mathbf{x}}^2 \left(\pi^{(a)}(\mathbf{v}^\top \mathbf{x}) \sigma^{(b)}(\mathbf{w}^\top \mathbf{x}) \right) &= \pi^{(a+2)}(\mathbf{v}^\top \mathbf{x}) \sigma^{(b)}(\mathbf{w}^\top \mathbf{x}) \mathbf{v} \mathbf{v}^\top + \pi^{(a)}(\mathbf{v}^\top \mathbf{x}) \sigma^{(b+2)}(\mathbf{w}^\top \mathbf{x}) \mathbf{w} \mathbf{w}^\top \\ &\quad + \pi^{(a+1)}(\mathbf{v}^\top \mathbf{x}) \sigma^{(b+1)}(\mathbf{w}^\top \mathbf{x}) (\mathbf{v} \mathbf{w}^\top + \mathbf{w} \mathbf{v}^\top) \end{aligned}$$

Therefore, for each of the rest property, it can be written in terms of a summation of terms of the form

$$\mathbb{E}_{\mathbf{x}} \left[\pi^{(a'_1)}(\mathbf{v}_1^\top \mathbf{x}) \sigma^{(b'_1)}(\mathbf{w}_1^\top \mathbf{x}) \pi^{(a'_2)}(\mathbf{v}_2^\top \mathbf{x}) \sigma^{(b'_2)}(\mathbf{w}_2^\top \mathbf{x}) \right]$$

Since $b_1 + b_2 \leq 3$, taking twice derivative gives $b'_1 + b'_2 \leq 5$. Therefore, applying Lemma 25 gives that all the rest terms are upper bounded by $\mathcal{O}(\varepsilon)$. \square

Lemma 25. *Let $\mathbf{v}_1, \mathbf{v}_2, \mathbf{w}_1, \mathbf{w}_2$ be unit vectors such that any two of the four have an inner product with magnitude upper bounded by ε . Then we have that for any a_1, a_2, b_1, b_2 such that $b_1 + b_2 \leq 5$*

$$\left| \mathbb{E}_{\mathbf{x}} \left[\pi^{(a_1)}(\mathbf{v}_1^\top \mathbf{x}) \pi^{(a_2)}(\mathbf{v}_2^\top \mathbf{x}) \sigma^{(b_1)}(\mathbf{w}_1^\top \mathbf{x}) \sigma^{(b_2)}(\mathbf{w}_2^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon)$$

Proof. Taking the Hermite expansion

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}} \left[\pi^{(a_1)}(\mathbf{v}_1^\top \mathbf{x}) \pi^{(a_2)}(\mathbf{v}_2^\top \mathbf{x}) \sigma^{(b_1)}(\mathbf{w}_1^\top \mathbf{x}) \sigma^{(b_2)}(\mathbf{w}_2^\top \mathbf{x}) \right] \\ &= \sum_{k, \ell=0}^{\infty} \frac{c_{k+a_1} c_{\ell+a_2}}{k! \ell!} \mathbb{E}_{\mathbf{x}} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_{3-b_1}(\mathbf{w}_1^\top \mathbf{x}) He_{3-b_2}(\mathbf{w}_2^\top \mathbf{x}) \right] \end{aligned}$$

We could observe that at least one of $3 - b_1$ and $3 - b_2$ is nonzero. Therefore, by Lemma 1, we have that $\mathbb{E}_{\mathbf{x}} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_{3-b_1}(\mathbf{w}_1^\top \mathbf{x}) He_{3-b_2}(\mathbf{w}_2^\top \mathbf{x}) \right]$ is a polynomial with lowest degree at most 1. Therefore, we have that

$$\left| \mathbb{E}_{\mathbf{x}} \left[He_k(\mathbf{v}_1^\top \mathbf{x}) He_\ell(\mathbf{v}_2^\top \mathbf{x}) He_{3-b_1}(\mathbf{w}_1^\top \mathbf{x}) He_{3-b_2}(\mathbf{w}_2^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon)$$

Moreover, this quantity is nonzero only when $k - \ell \leq 6 - b_1 - b_2$. Thus, by the boundedness of the Hermite coefficients, we can conclude that

$$\left| \mathbb{E}_{\mathbf{x}} \left[\pi^{(a_1)}(\mathbf{v}_1^\top \mathbf{x}) \pi^{(a_2)}(\mathbf{v}_2^\top \mathbf{x}) \sigma^{(b_1)}(\mathbf{w}_1^\top \mathbf{x}) \sigma^{(b_2)}(\mathbf{w}_2^\top \mathbf{x}) \right] \right| \leq \mathcal{O}(\varepsilon)$$

□

D.2 Other Auxiliary Results

Lemma 26. Let $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$, and define $f(\mathbf{v}) = \frac{1}{\|\mathbf{v}\|_2} \left(\mathbf{I} - \frac{\mathbf{v}\mathbf{v}^\top}{\|\mathbf{v}\|_2^2} \right) \mathbf{w}$. Then we have that

$$\mathcal{J}f(\mathbf{v}) = -\frac{\mathbf{v}^\top \mathbf{w}}{\|\mathbf{v}\|_2^3} \left(\mathbf{I} - \frac{\mathbf{v}\mathbf{v}^\top}{\|\mathbf{v}\|_2^2} \right) - \frac{1}{\|\mathbf{v}\|_2^3} (\mathbf{v}\mathbf{w}^\top + \mathbf{w}\mathbf{v}^\top)$$

Lemma 27. Let $f(x) = \frac{1}{1+e^{-x}}$ be the sigmoid function. Then we have that

- $f'''(x)^2 \leq 1$ for all $x \in \mathbb{R}$
- $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[f(x)] = \frac{1}{2}$
- $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[f^{(2k)}(x)] = 0$ for all $k \geq 1$

Proof. Using simple calculations, we can obtain that

$$f'(x) = f(x)(1 - f(x))$$

This gives that

$$f''(x) = f'(x)(1 - f(x)) - f(x)f'(x) = f'(x)(1 - 2f(x))$$

Thus, $f'''(x)$ can be written as

$$\begin{aligned} f'''(x) &= f''(x)(1 - 2f(x)) - 2f'(x)^2 \\ &= f'(x)(1 - 2f(x))^2 - 2f'(x)^2 \\ &= f'(x)(1 - 4f(x) + 4f(x)^2 - 2f'(x)) \\ &= f'(x)(1 - 6f(x) + 6f(x)^2) \end{aligned}$$

On the range $[0, 1]$, the function $1 - 6y + 6y^2$ takes extremes at $y = 0, \frac{1}{2}, 1$. At $y = 0$ and $y = 1$, we have $1 - 6y + 6y^2 = 1$. At $y = \frac{1}{2}$, we have that $1 - 6y + 6y^2 = -\frac{1}{2}$. Thus, we can conclude that $|1 - 6y + 6y^2| \leq 1$ for all $y \in [0, 1]$. Moreover, we have that $f'(x) = f(x)(1 - f(x)) \in [0, 1]$, since $f(x) \in [0, 1]$. Therefore, we can conclude that $|f'''(x)| \leq 1$, which implies the first property. To prove the second, we notice that

$$f(-x) = \frac{1}{1+e^x} = \frac{e^{-x}}{1+e^{-x}} = 1 - f(x)$$

Therefore, due to the symmetry of Gaussian distribution, we have that

$$\mathbb{E}_{x \sim \mathcal{N}(0,1)}[f(x)] = \mathbb{E}_{x \sim \mathcal{N}(0,1)}[f(-x)] = 1 - \mathbb{E}_{x \sim \mathcal{N}(0,1)}[f(x)]$$

This gives that $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[f(x)] = \frac{1}{2}$. To prove the third property, we notice that

$$f'(-x) = f(-x)(1 - f(-x)) = (1 - f(x))f(x) = f'(x)$$

which shows that $f'(-x)$ is even. Therefore, $f^{(2k)}(x)$ are odd functions for all $k \geq 1$. This implies the third property. □

Lemma 28. Consider function $f(x), g(x), h(x)$ given by the ODE system

$$\begin{aligned} f'(x) &= -a_1 f(x) + b_1 g(x) + c_1 h(x) + p \\ g'(x) &= -a_2 f(x) - b_2 g(x) + c_2 h(x) + p \\ h'(x) &= -a_3 f(x) + b_3 g(x) - c_3 h(x) + p \end{aligned}$$

for some $a_1, a_2, a_3, b_1, b_2, b_3, c_1, c_2, c_3 > 0$. If $b_2 c_3 \geq b_3 c_2$ and

$$a_1^2 b_2 + a_1^2 c_3 + a_1 a_2 b_1 + a_1 a_3 c_1 + a_1 b_2^2 + a_2 b_1 b_2 + a_1 c_3^2 + a_3 c_1 c_3 \geq a_2 b_3 c_1 + a_3 b_1 c_2$$

then we have that

$$\max \{|f(x)|, |g(x)|, |h(x)|\} \leq e^{-\Omega(x)} (|f(0)| + |g(0)| + |h(0)|) + \mathcal{O}(p)$$

for all $x \geq 0$

Proof. Let $\mathbf{A} \in \mathbb{R}^{3 \times 3}$ be given by

$$\mathbf{A} = \begin{bmatrix} -a_1 & b_1 & c_1 \\ -a_2 & -b_2 & c_2 \\ -a_3 & b_3 & -c_3 \end{bmatrix}$$

Then we have that

$$\begin{bmatrix} f'(x) \\ g'(x) \\ h'(x) \end{bmatrix} = \mathbf{A} \begin{bmatrix} f(x) \\ g(x) \\ h(x) \end{bmatrix} + p \mathbf{1}$$

Solving the system gives

$$\begin{bmatrix} f(x) \\ g(x) \\ h(x) \end{bmatrix} = e^{\mathbf{A}x} \begin{bmatrix} f(0) \\ g(0) \\ h(0) \end{bmatrix} + p \mathbf{A}^{-1} (e^{\mathbf{A}x} - \mathbf{I}) \mathbf{1}$$

Let λ be the eigenvalue of A with the largest real part. Then we have that

$$\left\| \begin{bmatrix} f(x) \\ g(x) \\ h(x) \end{bmatrix} \right\|_2 \leq e^{\lambda x} \left\| \begin{bmatrix} f(0) \\ g(0) \\ h(0) \end{bmatrix} \right\| + \mathcal{O}(p)$$

Thus, it suffice to show that all eigenvalues of \mathbf{A} has negative real parts. To do this, we write out the characteristic polynomial of \mathbf{A} as

$$P(y) = y^3 - \text{Tr}(\mathbf{A}) y^2 + \frac{1}{2} \left(\text{Tr}(\mathbf{A})^2 - \text{Tr}(\mathbf{A}^2) \right) y - \det(\mathbf{A})$$

By the Routh-Hurwitz criteria, it suffice to show that

$$\text{Tr}(\mathbf{A}) < 0, \det(\mathbf{A}) < 0, \frac{1}{2} \text{Tr}(\mathbf{A}) \left(\text{Tr}(\mathbf{A})^2 - \text{Tr}(\mathbf{A}^2) \right) < \det(\mathbf{A})$$

With the form of \mathbf{A} , we obtain that

$$\begin{aligned} \text{Tr}(\mathbf{A}) &= -(a_1 + b_2 + c_3) \\ \frac{1}{2} \text{Tr}(\mathbf{A}) \left(\text{Tr}(\mathbf{A})^2 - \text{Tr}(\mathbf{A}^2) \right) &= a_1 b_2 + a_1 c_3 + b_2 c_3 + a_2 b_1 + a_3 c_1 - b_3 c_2 \\ \det(\mathbf{A}) &= -a_1 b_2 c_3 - a_2 b_3 c_1 - a_3 b_1 c_2 - a_3 b_2 c_1 - a_2 b_1 c_3 - a_1 b_3 c_2 \end{aligned}$$

Thus, it is easy to see that $\text{Tr}(\mathbf{A}) < 0, \det(\mathbf{A}) < 0$. It remains to show that $\frac{1}{2} \text{Tr}(\mathbf{A}) \left(\text{Tr}(\mathbf{A})^2 - \text{Tr}(\mathbf{A}^2) \right) <$

$\det(\mathbf{A})$. This is equivalent to show that $S \geq 0$ with

$$\begin{aligned}
 S &= \det(\mathbf{A}) - \frac{1}{2} \text{Tr}(\mathbf{A}) \left(\text{Tr}(\mathbf{A})^2 - \text{Tr}(\mathbf{A}^2) \right) \\
 &= (a_1 + b_2 + c_3) (a_1 b_2 + a_1 c_3 + b_2 c_3 + a_2 b_1 + a_3 c_1 - b_3 c_2) \\
 &\quad - (a_1 b_2 c_3 + a_2 b_3 c_1 + a_3 b_1 c_2 + a_3 b_2 c_1 + a_2 b_1 c_3 + a_1 b_3 c_2) \\
 &= a_1^2 b_2 + a_1^2 c_3 + a_1 b_2 c_3 + a_1 a_2 b_1 + a_1 a_3 c_1 - a_1 b_3 c_2 + a_1 b_2^2 + a_1 b_2 c_3 + b_2^2 c_3 + a_2 b_1 b_2 \\
 &\quad + a_3 b_2 c_1 - b_2 b_3 c_2 + a_1 b_2 c_3 + a_1 c_3^2 + b_2 c_3^2 + a_2 b_1 c_3 + a_3 c_1 c_3 - b_3 c_2 c_3 \\
 &\quad - a_1 b_2 c_3 - a_2 b_3 c_1 - a_3 b_1 c_2 - a_3 b_2 c_1 - a_2 b_1 c_3 - a_1 b_3 c_2 \\
 &= a_1^2 b_2 + a_1^2 c_3 + 2a_1 b_2 c_3 + a_1 a_2 b_1 + a_1 a_3 c_1 - 2a_1 b_3 c_2 + a_1 b_2^2 + b_2^2 c_3 + a_2 b_1 b_2 \\
 &\quad - b_2 b_3 c_2 + a_1 c_3^2 + b_2 c_3^2 + a_3 c_1 c_3 - b_3 c_2 c_3 - a_2 b_3 c_1 - a_3 b_1 c_2
 \end{aligned}$$

If $b_2 c_3 \geq b_3 c_2$, then we have that

$$a_1 b_2 c_3 \geq a_1 b_3 c_2; \quad b_2 c_3^2 \geq b_3 c_2 c_3; \quad b_2^2 c_3 \geq b_2 b_3 c_2$$

This gives that

$$S \geq a_1^2 b_2 + a_1^2 c_3 + a_1 a_2 b_1 + a_1 a_3 c_1 + a_1 b_2^2 + a_2 b_1 b_2 + a_1 c_3^2 + a_3 c_1 c_3 - a_2 b_3 c_1 - a_3 b_1 c_2$$

□

Lemma 29. Let c_k be the k th order Hermite coefficient of $\pi(x)$. Let $z_1, z_2 \sim \mathcal{N}(0, 1)$ with $\text{Cov}(z_1, z_2) = \rho$. If $\mathbb{E}_{z_1, z_2}[\pi'(z_1) \pi'''(z_2)] \leq 0$ for all ρ , then we have that

$$\sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!} \gamma^k \leq 0; \quad \forall \gamma > 0$$

Proof. Notice that, by taking the Hermite expansion of $\pi(x)$ and $\pi''(x)$, we have that for $z_1, z_2 \sim \mathcal{N}(0, 1)$ with $\text{Cov}(z_1, z_2) = \gamma$

$$\mathbb{E}_{z_1, z_2}[\pi(z_1) \pi''(z_2)] = \sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!} \gamma^k$$

Let $f(\gamma) = \mathbb{E}_{z_1, z_2}[\pi(z_1) \pi''(z_2)]$. Then by Price's Theorem we have that

$$f'(\gamma) = \mathbb{E}_{z_1, z_2}[\pi'(z_1) \pi'''(z_2)] \leq 0$$

Moreover, at $\gamma = 0$, we have that

$$f(0) = \mathbb{E}_{z_1}[\pi(z_1)] \mathbb{E}_{z_2}[\pi''(z_2)] = 0$$

where the last equality is due to Lemma 27. Therefore, we can conclude that

$$\sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!} \gamma^k = f(\gamma) \leq f(0) = 0$$

□

Lemma 30. Let $x \in (-1, 1)$, and consider $S = \sum_{k=0}^{\infty} \frac{c_k^2}{k!} x^k$. If $c_k \leq \sqrt{B \cdot k!}$, then we have that

$$c_0^2 - \frac{|x|}{1 - |x|} \leq S \leq c_0^2 + \frac{|x|}{1 - |x|}$$

Proof. We write S as

$$S = c_0^2 + \sum_{k=1}^{\infty} \frac{c_k^2}{k!} x^k$$

Notice that

$$\left| \sum_{k=1}^{\infty} \frac{c_k^2}{k!} x^k \right| \leq \sum_{k=1}^{\infty} \frac{c_k^2}{k!} |x|^k \leq B \sum_{k=1}^{\infty} |x|^k = \frac{B|x|}{1-|x|}$$

Therefore, we can conclude that

$$c_0^2 - \frac{B|x|}{1-|x|} \leq S \leq c_0^2 + \frac{B|x|}{1-|x|}$$

□

Lemma 31. *Let c_k denote the k th order Hermite coefficient of $\pi(\cdot)$ such that $c_2 = 0$ and $c_{k+3} \leq \sqrt{B \cdot k!}$ for all $k \geq 0$. Let $\gamma \in (b, 1]$. If $\mathbb{E}_{z_1, z_2 \sim \mathcal{N}(0,1), \text{Cov}(z_1, z_2) = \rho} [\pi'(z_1) \pi^{(3)}(z_2)] \leq 0$ for all $\rho \in [-1, 1]$, then we have that*

$$\sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!} \gamma^k \leq \frac{|b|}{1-|b|}$$

Proof. Notice that

$$\mathbb{E}_{z_1, z_2} [\pi(z_1) \pi''(z_2)] = \mathbb{E}_{z_1, z_2} \left[\left(\sum_{k=0}^{\infty} \frac{c_k}{k!} \text{Herm}_k(z_1) \right) \left(\sum_{k=0}^{\infty} \frac{c_{k+2}}{k!} \text{Herm}_k(z_1) \right) \right] = \sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!}$$

Therefore, we can define

$$h(\rho) = \sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!} \rho^k = \mathbb{E}_{z_1, z_2} [\pi(z_1) \pi''(z_2)] \rho^k$$

By Price's Theorem, we have that

$$h'(\rho) = \mathbb{E}_{z_1, z_2} [\pi'(z_1) \pi^{(3)}(z_2)] \leq 0$$

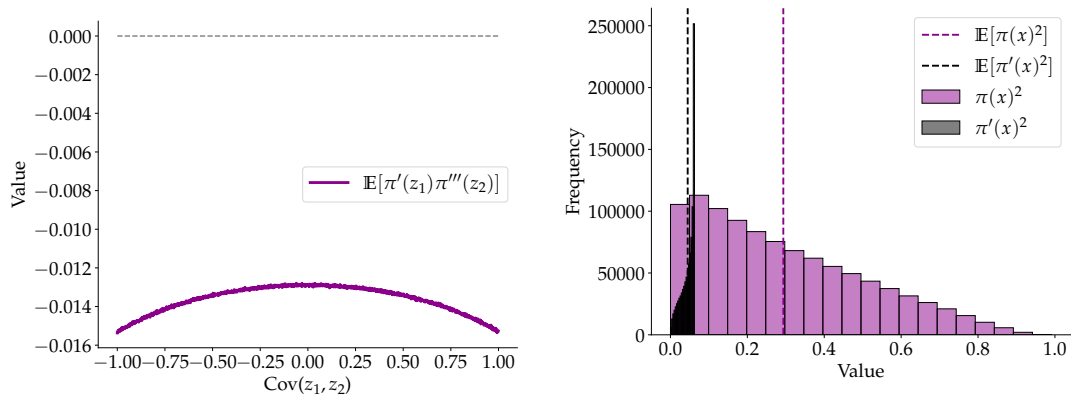
Therefore, $h(\rho)$ decreases monotonically, which implies that

$$h(\rho) = \sum_{k=0}^{\infty} \frac{c_k c_{k+2}}{k!} \rho^k \leq h(b) = b \sum_{k=0}^{\infty} \frac{c_{k+1} c_{k+3}}{k!} b^k \leq \sum_{k=1}^{\infty} |b|^k \leq \frac{|b|}{1-|b|}$$

□

E Plotting Sigmoid Property

In this section, we plot the simulation result of the properties of the sigmoid function. Figure 3a is generated by taking 10^5 samples of correlated Gaussian random variables for each covariance value $\rho \in [-1, 1]$. Figure 3b is generated by taking 10^6 samples of standard Gaussian random variables.



(a) Value of $\mathbb{E}_{z_1, z_2 \sim \mathcal{N}(0,1)}[\pi'(z_1)\pi'''(z_2)]$ for different values of $\text{Cov}(z_1, z_2)$. As in the figure, all values are negative.

(b) Values of $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[\pi(x)^2]$ and $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[\pi'(x)^2]$. Figure shows that $\mathbb{E}_{x \sim \mathcal{N}(0,1)}[\pi(x)^2] \geq 1.1\mathbb{E}_{x \sim \mathcal{N}(0,1)}[\pi'(x)^2]$.

Figure 3: Plot of the property of $\pi(x)$