

# BENIGN OVERFITTING IN ADVERSARIAL TRAINING FOR VISION TRANSFORMERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Despite the remarkable success of Vision Transformers (ViTs) across a wide range of vision tasks, recent studies have revealed that they remain vulnerable to adversarial examples, much like Convolutional Neural Networks (CNNs). A common empirical defense strategy is adversarial training, yet the theoretical underpinnings of its robustness in ViTs remain largely unexplored. In this work, we present the first theoretical analysis of adversarial training under simplified ViT architectures. We show that, when trained under a signal-to-noise ratio that satisfies a certain condition and within a moderate  $\ell_2$  perturbation budget, adversarial training enables ViTs to achieve nearly zero robust training loss and robust generalization error under certain regimes. Remarkably, this leads to strong generalization even in the presence of overfitting, a phenomenon known as *benign overfitting*, previously only observed in CNNs (with adversarial training). Experiments on both synthetic and real-world datasets further validate our theoretical findings.

## 1 INTRODUCTION

Vision Transformers (ViTs) have emerged as a powerful alternative to convolutional neural networks (CNNs) (Krizhevsky et al., 2012) for a wide range of computer vision tasks, including image classification, object detection, semantic segmentation, and vision-language modeling (Dosovitskiy et al., 2020; Liu et al., 2021; Hu et al., 2025a). Despite their strong performance, recent studies (Hu et al., 2023; Shao et al., 2021; Hu et al., 2025b; 2024) have revealed that ViTs, similar to CNNs, can still be vulnerable to small, carefully crafted perturbations. These perturbations, referring to adversarial examples (Szegedy et al., 2013), can often cause significant performance degradation.

A widely studied defense mechanism against such vulnerabilities is *adversarial training* (Goodfellow et al., 2014; 2016), which augments the training process with adversarially perturbed samples to improve model robustness. While adversarial training has proven effective in enhancing model robustness, it is frequently accompanied by a noticeable degradation in generalization performance on clean data (Raghunathan et al., 2019; Xu et al., 2023). Such a robustness-generalization trade-off (Fu & Wang, 2023; Zhang et al., 2021; Xiao et al., 2022) has been extensively studied, raising the question of whether it is possible to preserve robustness without sacrificing clean-data accuracy.

*Benign Overfitting* (Bartlett et al., 2020) refers to the phenomenon that overparameterized models can interpolate the training data (achieving near-zero empirical risk) yet still generalize well. In the standard (non-adversarial) setting, this behavior has been studied across various model architectures, including linear regression (Bartlett et al., 2020; Zhou & Ge, 2023), logistic regression Wang et al. (2021); Cao et al. (2021), ridge regression (Tsigler & Bartlett, 2023), kernel methods (Mei & Montanari, 2022), and neural networks (Li et al., 2021; Cao et al., 2022; Frei et al., 2023; Jiang et al., 2024; Frei & Vardi, 2024). In the adversarial setting, Chen et al. (2023) provides the first theoretical analysis that benign overfitting can arise in adversarial training for linear regression under sub-Gaussian mixture models. Subsequently, Wang et al. (2024) shows that adversarial training can still generalize well in the presence of inference-time attacks for two-layer neural networks under appropriate distributional assumptions. However, to the best of our knowledge, it remains unclear whether analogous behavior occurs in more advanced architectures such as ViTs.

Compared with linear and two-layer models, analyzing ViTs poses distinctive challenges for analyzing robust benign overfitting. Unlike CNNs or two-layer neural networks with activation functions followed by a linear model, ViTs incorporate attention heads with query, key, and value projection

matrices, leading to substantially more complex training dynamics. The effect of perturbations on attention heads differs significantly from their effect on linear layers or activation-plus-linear structures, and varying perturbation magnitudes can markedly influence the training dynamics of ViTs (A more detailed discussion is provided in Section 4). Thus, the benign overfitting behavior or its conditions for ViTs may be significantly different from previous models.

To fill the gap, in this work, we provide the first comprehensive theoretical analysis of robust benign overfitting for a simplified ViT model. Specifically,

1. We validate that *benign overfitting* can also arise in adversarially trained Vision Transformers when the signal-to-noise ratio and the perturbation magnitude satisfy certain conditions, similar to linear and two-layer neural network models. That is, the adversarially trained interpolator attains near-zero robust training loss while maintaining small robust test error.
2. By analyzing the adversarial training dynamics of ViTs, we identify three key regimes: 1) Small perturbations yield trajectories close to clean training; 2) Moderate perturbations cause the attention mechanism to fail, such that the ViT collapses to a linear model; 3) Large perturbations lead to significant generalization error beyond benign overfitting. In all cases, we provide explicit upper bounds on the clean and robust test error.
3. We empirically verify our theoretical findings on both synthetic data and MNIST, demonstrating agreement between the derived bound and observed conditions for the occurrence of benign overfitting.

## 2 RELATED WORK

**Benign Overfitting.** Benign overfitting refers to the phenomenon where a predictor perfectly fits (interpolates) noisy training data yet still achieves strong generalization performance on unseen data. Bartlett et al. (2020) analyzed benign overfitting in linear regression with Gaussian noise, showing that in high-dimensional (overparameterized) regimes, the excess risk of interpolation can be asymptotically optimal. Foundational studies have further established benign overfitting in various linear settings, including regression, sparse regression, logistic regression, ridge regression, and kernel methods (Belkin et al., 2018; Bartlett et al., 2020; Hastie et al., 2022; Ding et al., 2024b;a). In neural networks, Frei et al. (2023) studied benign overfitting in two-layer networks without relying on the lazy training assumption in finite-width regimes. Cao et al. (2022); Kou et al. (2023) examined benign overfitting in convolutional networks from a feature learning perspective. Recently, Jiang et al. (2024) investigated learning dynamics in Vision Transformers (ViTs), delineating the boundary between benign and harmful overfitting. Frei & Vardi (2024) also analyzed benign overfitting in trained Transformer classifiers within in-context learning setups.

**Benign Overfitting for Adversarial Training.** Chen et al. (2023) initiated the study of benign overfitting under adversarial training in the context of linear classifiers with sub-Gaussian mixture data, proving that under moderate perturbations, linear classifiers can achieve near-optimal standard and adversarial risks. Building on this line of work, Wang et al. (2024) extended the analysis to two-layer networks and demonstrated that, regardless of whether the activation function is smooth or non-smooth, adversarial training can achieve near-optimal robust generalization error. In contrast, Hao & Zhang (2024) showed that under non-negligible noise, linear regression and NTK regression models tend to overfit the training data, yielding estimators with inflated Lipschitz norms and consequently elevated adversarial risk. However, none of them considered and their conclusions do not hold for transformer architectures.

**Adversarial Robustness in Transformer.** Adversarial training has been widely used to improve Transformer robustness. For Vision Transformers (ViTs), Herrmann et al. (2022) proposed PyramidAT, combining consistent dropout and stochastic depth to alleviate performance degradation, while Wu et al. (2022) reduced computational cost via attention-guided dropping of patch embeddings. Gopal et al. (2025) introduced SAFER, a layer-selective fine-tuning method with sharpness-aware minimization to mitigate adversarial overfitting, and Islam et al. (2025) analyzed layer-wise perturbation propagation, proposing a neuron-level suppression mechanism. However, these works are purely empirical and lack theoretical guarantees. Recent studies have begun to analyze robustness from a theoretical perspective in linear Transformer in-context learning, showing that robustness can be enhanced through adversarial training against hijacking attacks (Anwar et al.), short-suffix

training to defend long-suffix jailbreaks (Fu et al., 2025), and multi-task adversarial pretraining without downstream AT (Kumano et al., 2025). Yet, these analyses are restricted to the in-context learning setting with linear Transformers, and theoretical guarantees for adversarial robustness in general Transformer architectures remain an open problem.

**Differences in Adversarial Robustness Between Transformer and CNN.** Several studies have compared Transformers and CNNs under adversarial attacks. Bai et al. (2021) observe that under unified training settings, Transformers are not inherently more robust, with their OOD generalization mainly attributed to self-attention. Mo et al. (2022) show that under standard training, Transformers do not necessarily outperform CNNs under adversarial attack, and propose training strategies to improve ViT robustness. Dingeto & Kim (2024) propose a regularization method that enables ViTs to exhibit stronger adversarial robustness than CNNs. These studies are largely empirical, focusing on performance differences, and do not analyze the learning dynamics or the role of attention in adversarial training.

### 3 PROBLEM SETUP

In this section, we introduce the necessary definitions and formally describe the Gradient Descent-based Adversarial Training under the multi-patch data distribution and the two-layer Vision Transformer model.

**Notations.** For two sequences  $\{x_n\}$  and  $\{y_n\}$ , we write  $x_n = O(y_n)$  if there exist absolute constants  $C > 0$  and  $N > 0$  such that  $|x_n| \leq C|y_n|$  for all  $n \geq N$ . Similarly, we write  $x_n = \Omega(y_n)$  if  $y_n = O(x_n)$ . We say  $x_n = \Theta(y_n)$  if both  $x_n = O(y_n)$  and  $x_n = \Omega(y_n)$ . Moreover,  $x_n = o(y_n)$  if  $\lim_{n \rightarrow \infty} |x_n/y_n| = 0$ . Finally, we use  $\tilde{O}(\cdot)$ ,  $\tilde{\Omega}(\cdot)$ , and  $\tilde{\Theta}(\cdot)$  to denote the corresponding notations with logarithmic factors suppressed.

**Definition 1 (Data Generation Model).** Let  $\boldsymbol{\mu}_+, \boldsymbol{\mu}_- \in \mathbb{R}^d$  be fixed vectors representing the signals contained in data points, where  $\|\boldsymbol{\mu}_+\|_2 = \|\boldsymbol{\mu}_-\|_2 = \|\boldsymbol{\mu}\|_2$  and  $\langle \boldsymbol{\mu}_+, \boldsymbol{\mu}_- \rangle = 0$ . Then each data point  $(\mathbf{X}, y)$  with  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M)^\top \in \mathbb{R}^{M \times d}$  and  $y \in \{-1, 1\}$  is generated from the following distribution  $D$ :

- (1) The label  $y$  is generated as a Rademacher random variable, i.e.,  $\mathbb{P}(y = 1) = \mathbb{P}(y = -1) = \frac{1}{2}$ .
- (2) If  $y = 1$  then  $\mathbf{x}_1$  is given as  $\boldsymbol{\mu}_+$ , if  $y = -1$  then  $\mathbf{x}_1$  is given as  $\boldsymbol{\mu}_-$ , which represents signals.
- (3)  $\mathbf{x}_2, \dots, \mathbf{x}_M$  are given by noise vectors  $\boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_M$ , generated i.i.d from the Gaussian distribution  $\mathcal{N}(0, \sigma_p^2 \cdot (\mathbf{I} - \boldsymbol{\mu}_+ \boldsymbol{\mu}_+^\top \cdot \|\boldsymbol{\mu}\|_2^{-2} - \boldsymbol{\mu}_- \boldsymbol{\mu}_-^\top \cdot \|\boldsymbol{\mu}\|_2^{-2}))$ , which represent noises.

Our data generation model is motivated by the patch-level structure of real image data, where some patches encode class-relevant signals (e.g. semantic information) while others capture irrelevant noise (e.g. background artifacts). Similar constructions have been widely employed in the feature learning literature to analyze the generalization behavior of overparameterized classifiers (Allen-Zhu & Li, 2020; Cao et al., 2022; Jelassi & Li, 2022; Kou et al., 2023; Zou et al., 2023; Jiang et al., 2024; Han et al., 2024; Ding et al., 2025). In our setup, the noise component  $\boldsymbol{\xi}$  is modeled as a Gaussian variable, with its covariance structure designed to remain orthogonal to the signal component  $\boldsymbol{\mu}$ , ensuring that the data noise is independent of and unrelated to the feature.

**Two-layer Transformer.** Following the architecture introduced by Jiang et al. (2024), we consider a simplified two-layer Transformer consisting of a self-attention layer followed by a fixed linear layer, defined as the following, where  $\theta = (\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V)$ .

$$f(\mathbf{X}, \theta) = \frac{1}{M} \sum_{l=1}^M \varphi(\mathbf{x}_l^\top \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top) \mathbf{X} \mathbf{W}_V \mathbf{w}_O. \quad (1)$$

Here,  $\varphi(\cdot) : \mathbb{R}^M \rightarrow \mathbb{R}^M$  denotes the softmax function;  $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d_n}$  and  $\mathbf{W}_V \in \mathbb{R}^{d \times d_v}$  represent the query, key, and value matrices, respectively; and  $\mathbf{w}_O \in \mathbb{R}^{d_v}$  represents the weight vector of the linear layer. We use  $\theta$  to denote the collection of all the model weights. This model is not reduced to a linear or single-layer attention architecture; instead, it more closely resembles the structure of a real Transformer, with a correspondingly more complex parameter update process.

**Loss Function.** Let  $S = \{(\mathbf{X}_n, y_n)\}_{n=1}^N$  denote the training dataset drawn from the distribution  $D$  defined in Definition 1, where  $n$  indexes the samples (so  $(\mathbf{X}_n, y_n)$  represents the  $n$ -th sample). In

**Algorithm 1** Gradient Descent-based Adversarial Training

---

**Require:** Step size  $\eta$ , perturbation budget per token  $\tau$ , number of iterations  $T$ , standard deviation  $\sigma_V, \sigma_h$

- 1: Initialize network weights  $(\mathbf{W}_Q^0)_{ij}, (\mathbf{W}_K^0)_{ij} \sim \mathcal{N}(0, \sigma_h^2), (\mathbf{W}_V^0)_{ij} \sim \mathcal{N}(0, \sigma_V^2)$  i.i.d.
- 2: **for**  $t = 0, \dots, T - 1$  **do**
- 3:     **for**  $n = 1, \dots, N$  **do**
- 4:         Generate adversarial example for input  $\mathbf{X}_n$ :  
 $\tilde{\mathbf{X}}_n^t \leftarrow \arg \max_{\tilde{\mathbf{X}}_n \in B(\mathbf{X}_n, \tau)} \ell(y_n f(\tilde{\mathbf{X}}_n; \mathbf{W}_Q^t, \mathbf{W}_K^t, \mathbf{W}_V^t))$
- 5:     **end for**
- 6:     Update all weight matrices simultaneously:  $(\mathbf{W}_Q^{t+1}, \mathbf{W}_K^{t+1}, \mathbf{W}_V^{t+1}) \leftarrow (\mathbf{W}_Q^t, \mathbf{W}_K^t, \mathbf{W}_V^t) - \frac{\eta}{N} \sum_{n=1}^N (\nabla_{\mathbf{W}_Q}, \nabla_{\mathbf{W}_K}, \nabla_{\mathbf{W}_V}) \ell(y_n f(\tilde{\mathbf{X}}_n^t; \mathbf{W}_Q^t, \mathbf{W}_K^t, \mathbf{W}_V^t))$
- 7: **end for**
- 8: **return**  $\theta(T) = (\mathbf{W}_Q^T, \mathbf{W}_K^T, \mathbf{W}_V^T)$

---

this work, we adopt the empirical cross-entropy loss as a surrogate for the non-differentiable 0/1 loss, and train the two-layer Transformer by minimizing this loss:

$$L_S(\theta) = \frac{1}{N} \sum_{n=1}^N \ell(y_n f(\mathbf{X}_n, \theta)),$$

where  $\ell(z) = \log(1 + \exp(-z))$  and  $f(\mathbf{X}, \theta)$  is the two-layer Transformer. We measure the generalization ability of the two-layer Transformer using the *test error*, defined as the expected 0/1 loss over the data distribution  $D$ :

$$L_D(\theta) = \mathbb{E}_{(\mathbf{x}, y) \sim D} \mathbb{1}(yf(\mathbf{X}, \theta) \leq 0).$$

**Robust Loss.** We consider  $\ell_2$  norm-bounded adversarial perturbations applied to each component of the input sequence  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_M] \in \mathcal{X}^M$ , where each  $\mathbf{x}_m \in \mathbb{R}^d$  denotes a token (or patch) embedding. For a perturbation budget  $\tau > 0$ , the admissible perturbation set is  $B(\mathbf{X}, \tau) := \{ \tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_M] \mid \|\tilde{\mathbf{x}}_m - \mathbf{x}_m\|_2 \leq \tau, \forall m \in [M] \}$ . Under this threat model, the robust 0/1 loss is defined as  $\ell_{\text{rob}}^{0/1}(yf(\mathbf{X}, \theta)) := \max_{\tilde{\mathbf{X}} \in B(\mathbf{X}, \tau)} \mathbb{1}(yf(\tilde{\mathbf{X}}, \theta) \leq 0)$ , and the robust loss is defined as  $\ell_{\text{rob}}(y_n f(\mathbf{X}_n, \theta)) := \max_{\tilde{\mathbf{X}}_n \in B(\mathbf{X}_n, \tau)} \ell(y_n f(\tilde{\mathbf{X}}_n, \theta))$ . The robust test error and robust test loss are:

$$L_D^{\text{rob}}(\theta) := \mathbb{E}_{(\mathbf{x}, y) \sim D} \ell_{\text{rob}}^{0/1}(yf(\mathbf{X}, \theta)), \quad L_S^{\text{rob}}(\theta) := \frac{1}{N} \sum_{n=1}^N \ell_{\text{rob}}(y_n f(\mathbf{X}_n, \theta))$$

**Adversarial Training.** We adopt a Gradient Descent-based Adversarial Training algorithm to update the network parameters, as summarized in Algorithm 1. Note that we initialize the network weights  $\mathbf{W}_Q, \mathbf{W}_K$ , and  $\mathbf{W}_V$  with Gaussian distributions, where each entry of  $\mathbf{W}_Q$  and  $\mathbf{W}_K$  is drawn from  $\mathcal{N}(0, \sigma_h^2)$ , and each entry of  $\mathbf{W}_V$  is drawn from  $\mathcal{N}(0, \sigma_V^2)$ . The algorithm iteratively constructs adversarial training examples by maximizing the training loss with respect to the input and updates the model parameters based on them. In our setting, only the attention projection matrices  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$  are updated during training, while the output vector  $\mathbf{w}_O$  remains fixed.

## 4 MAIN RESULTS

In this section, we present our main theoretical results on the convergence and generalization of the ViT model, demonstrating how the signal-to-noise ratio  $\text{SNR} = \|\boldsymbol{\mu}\|_2 / (\sigma_p \sqrt{d})$  and the sample size  $N$  influence its adversarial training dynamics. We first introduce the following conditions.

**Condition 1.** Given a sufficiently small failure probability  $\delta > 0$  and a target training loss  $\epsilon > 0$ , suppose that:

- (1) The dimension  $d$  and  $d_h$  are sufficiently large satisfying  $d = \tilde{\Omega}(\epsilon^{-2} N^2 d_h)$  and  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$ .

- 216 (2) The training sample size  $N$  is large enough such that  $N = \Omega(\text{poly log}(d))$ .  
 217 (3) The number of input tokens is bounded as  $M = \Theta(1)$ , and the  $\ell_2$ -norm of linear layer weights  
 218 satisfies  $\|\mathbf{w}_O\|_2 = \Theta(1)$ .  
 219 (4) The learning rate  $\eta$  is chosen sufficiently small so that  $\eta \lesssim \tilde{O}(\min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}})$ .  
 220 (5) The Gaussian initialization is appropriately chosen such that the standard deviation  $\sigma_V$  sat-  
 221 isfies  $\sigma_V \leq \tilde{O}(\|\mathbf{w}_O\|_2^{-1} \cdot \min\{\|\boldsymbol{\mu}\|_2^{-1}, (\sigma_p \sqrt{d})^{-1}\} \cdot d_h^{-\frac{1}{4}} d^{-\frac{1}{2}})$ , and the variance  $\sigma_h^2$  satisfies  
 222  $\min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}} \cdot (\log(6N^2 M^2 / \delta))^{-2} \leq \sigma_h^2 \leq \min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}} \cdot$   
 223  $(\log(6N^2 M^2 / \delta))^{-\frac{3}{2}}$ .  
 224 (6) The target training loss is satisfying  $\epsilon \leq O(1/\text{poly log}(d))$ .  
 225  
 226  
 227

228 Conditions (1) and (2) ensure that the learning problem is set in a sufficiently over-parameterized  
 229 regime, allowing the model to fully capture the feature signal described in Definition 1. Con-  
 230 dition (3) guarantees that each class contains enough samples with high probability. Condi-  
 231 tions (4) and (5) simplify the analysis, though they can be generalized to the settings  $M = \Omega(1)$ ,  
 232  $\|\mathbf{w}_O\|_2 = o(1)$ , or  $\|\mathbf{w}_O\|_2 = \omega(1)$ . Together, Conditions (4) and (5) further ensure that the Trans-  
 233 former can be effectively trained. Finally, Condition (6) ensures that the Transformer sufficiently  
 234 overfits the training data. Similar conditions are widely made in the theoretical analysis of benign  
 235 overfitting in neural networks (Allen-Zhu & Li, 2020; Cao et al., 2022; Frei et al., 2022; Jelassi &  
 236 Li, 2022; Chatterji & Long, 2023; Zou et al., 2023; Kou et al., 2023; Frei & Vardi, 2024; Jiang et al.,  
 2024).

237 **Theorem 2** (Benign Overfitting under Adversarial Training). *Under Condition 1, we distinguish*  
 238 *two cases:*

239 **Case 1.** *If we have  $N \cdot \text{SNR}^2 = \Omega(1)$  and  $\tau \leq O(\frac{\|\boldsymbol{\mu}\|_2}{\log d_h})$ , ViT’s attention head is effectively trainable*  
 240 *. In this case, let  $T = \Theta(\eta^{-1} \epsilon^{-1} \|\boldsymbol{\mu}\|_2^{-2} \|\mathbf{w}_O\|_2^{-2})$ .*

242 **Case 2.** *If we have  $N \cdot \text{SNR}^2 = \Omega(\frac{1}{\epsilon})$  and  $\omega(\frac{\|\boldsymbol{\mu}\|_2}{\log d_h}) \leq \tau \leq O(\boldsymbol{\mu})$ , the attention head parameters*  
 243 *barely update, causing the attention weights to remain nearly uniform and the ViT degenerates into*  
 244 *a linear model . In this case, let  $T = M \cdot \Theta(\eta^{-1} \epsilon^{-1} \|\boldsymbol{\mu}\|_2^{-2} \|\mathbf{w}_O\|_2^{-2})$ .*

245 *In both cases, with probability at least  $1 - d^{-1}$ , the following holds:*

- 246 1. *The robust training loss converges to  $\epsilon$ :*

$$247 L_S^{\text{rob}}(\theta(T)) \leq \epsilon.$$

- 248 2. *The clean test error satisfies:*

$$249 L_D(\theta(T)) \leq \exp(-C \cdot d \text{SNR}^2).$$

- 250 3. *The robust test error satisfies:*

$$251 L_D^{\text{rob}}(\theta(T)) \leq \exp\left(-C \cdot d \text{SNR}^2 \left(1 - \frac{\tau}{\|\boldsymbol{\mu}\|_2}\right)^2\right).$$

252 Theorem 2 demonstrates that , under the assumption on  $\tau$ , the model attains adversarial robustness,  
 253 while the SNR condition guarantees that it prioritizes the signal over the noise, thereby leading to  
 254 benign overfitting. The adversarially-trained two-layer Transformer model exhibits three key ob-  
 255 servations. 1) The model fits the training data well, with training loss converging to  $\epsilon$ . 2) For  
 256 generalization guarantee, the clean test error decays rapidly with increasing  $d$  and SNR, and obvi-  
 257 ously less than  $\epsilon$  by Condition (2). This is consistent with classical benign overfitting phenomena  
 258 established in prior work (Jiang et al., 2024). 3) For robustness guarantees, Theorem 2 provides an  
 259 explicit upper bound on the robust test error as a function of the perturbation radius  $\tau$  and the sig-  
 260 nal strength  $\|\boldsymbol{\mu}\|_2$ , showing that the bound increases with increasing  $\tau$ , which is consist with prior  
 261 empirical observations(Madry et al., 2017; Schmidt et al., 2018).

262 **The impact of perturbation  $\tau$  radius.** Here, we give more discussion on the perturbation radius  
 263  $\tau$  to illustrate how it affects the different dynamics of ViTs, leading to the various results stated  
 264 in Theorem 2. First, it can be observed that the softmax structure in attention is highly sensitive  
 265 to perturbations according to Lemma 4. When  $\tau \leq O(\frac{\|\boldsymbol{\mu}\|_2}{\log d_h})$ , adversarial perturbations do not  
 266  
 267

dominate the learning dynamics, and the adversarial training trajectory of ViT remains close to standard clean training, corresponding to Case 1 in Theorem 2. In contrast, when  $\omega(\frac{\|\mu\|_2}{\log d_h}) \leq \tau \leq O(\mu)$ , the updates between signal and noise in attention are effectively canceled out by the perturbation, forcing attention weights to remain close to their initialized uniform distribution, under which ViT degenerates into a linear model, corresponding to Case 2 in Theorem 2. The relevant lemmas will be provided in Section 5.2 later.

**The difference between ViT and degenerated linear model.** Although both cases can achieve benign overfitting, the attention mechanism in ViTs in Case 1 enables the model to learn the signal more rapidly, leading to faster convergence, and allows it to extract useful information even from sparser signals. According to Theorem 2, for the degenerated linear model, the convergence time  $T$  is  $M$  times slower than that of the ViT, and achieving benign overfitting requires a higher signal-to-noise ratio,  $N \cdot \text{SNR}^2 = \Omega(1/\epsilon)$ , highlighting the advantages of the ViT architecture.

**Comparison with prior work.** The most closely related works to ours are Chen et al. (2023); Wang et al. (2024), which also investigated the phenomenon of benign overfitting under adversarial training. However, Chen et al. (2023) focused on linear regression models with moderate perturbations, and Wang et al. (2024) studied simplified neural networks, leaving more complex architectures unexplored. Our results on ViTs therefore complement this line of research. In addition, even in the more complex setting, our analysis does not rely on the implicit assumption of a large  $\|\mu\|_2$ , i.e.,  $\|\mu\|_2 = \Theta(d^r)$  for some  $r \in (1/4, 1/2]$  in (Chen et al., 2023). Moreover, our results require a minimum convergence time of  $T \sim O(N \max\{\text{SNR}^2, \text{SNR}^{-4}\})$ , which is significantly smaller than the convergence time  $T \sim O(\frac{d^2}{\|\mu\|^2 \epsilon^2})$  reported in prior studies on CNN models by (Wang et al., 2024) in the over-parameterized regime.

Next, we will show that once the perturbation radius  $\tau$  exceeds the signal strength  $\|\mu\|_2$ , no classifier can achieve nontrivial robust accuracy. This implies that excessive adversarial training leads to poor model performance, consistent with the findings of Wang et al. (2024). Combining with Theorem 2, we can see that the assumption of the relation between  $\tau$  and  $\|\mu\|_2$  is essential to understand benign overfitting. Also, our radius for benign overfitting is almost tight.

**Theorem 3.** *For any given classifier  $f(\cdot; \theta)$ , when  $\tau \geq \|\mu\|_2$ , the robust test error satisfies  $L_D^{\text{rob}}(\theta) \geq 0.25$ .*

## 5 PROOF SKETCH

In this section, we provide a proof sketch of the different adversarial training dynamics due to different perturbation radii. Based on the ViT formulation in Eq.(1), when the perturbations are relatively small, we show its impact on the learning of the  $\mathbf{W}_Q$  and  $\mathbf{W}_K$  matrices is limited. If more attention is allocated to signal tokens, the value matrix  $\mathbf{W}_V$  tends to align more closely with a perturbed signal vector  $\tilde{\mu}$ , while its direction remains dominated by the true signal  $\mu$ . Consequently, the gradients propagated to the attention heads associated with the signal are larger than those to the noise, forming a positive feedback loop that facilitates better generalization, similar to the learning dynamics observed in the clean training (Jiang et al., 2024). When the perturbations are moderate, adversarial training suppresses the learning of token-to-signal attention, keeping the attention weights near their initialization. When  $\mathbf{W}_V$  and  $\mathbf{W}_K$  are initialized with small Gaussian variance, the attention distribution remains nearly uniform, and the ViT degenerates into a linear model. In this regime, the model requires a larger SNR to align  $\mathbf{W}_V \mathbf{w}_O$  with the signal  $\mu$ , thereby achieving benign overfitting.

### 5.1 VECTORIZED Q & K AND SCALARIZED V WITH TIME-INDEPENDENT PERTURBED INPUTS

First, unlike CNNs that treat convolutional kernels as vectors for signal-noise decomposition, our setting requires handling the more complex interactions among the matrices  $\mathbf{W}_V$ ,  $\mathbf{W}_K$ , and  $\mathbf{W}_Q$ . Therefore, we consider vectorizing  $\mathbf{W}_K$ , and  $\mathbf{W}_Q$  and scalarizing  $\mathbf{W}_V$ . We fix a universal perturbed input  $\tilde{\mathbf{X}} = [\tilde{\mu}, \tilde{\xi}_{n,2}, \dots, \tilde{\xi}_{n,M}]$ , where  $\tilde{\mu}_+ \in B(\mu_+, \tau)$ ,  $\tilde{\mu}_- \in B(\mu_-, \tau)$ , and  $\tilde{\xi}_{n,i} \in B(\xi_{n,i}, \tau)$  are chosen once and remain fixed for all iterations  $t$ . These perturbations are universal and do not correspond to iteration-specific adversarial examples. For example, for  $\mathbf{W}_V$  we have the following.

**Definition 2** (Scalarized V). Let  $\mathbf{W}_V^{(t)}$  denote the V matrix of the ViT at the  $t$ -th iteration of adversarial training. For the fixed perturbed vectors above, there exist scalars  $\gamma_{V,+}^{(t)}$ ,  $\gamma_{V,-}^{(t)}$ , and  $\rho_{V,n,i}^{(t)}$  such that

$$\begin{aligned}\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_V^{(t)} \mathbf{w}_O &= \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_V^{(0)} \mathbf{w}_O + \gamma_{V,+}^{(t)} \|\mathbf{w}_O\|_2^2, \\ \tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_V^{(t)} \mathbf{w}_O &= \tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_V^{(0)} \mathbf{w}_O + \gamma_{V,-}^{(t)} \|\mathbf{w}_O\|_2^2, \\ \tilde{\boldsymbol{\xi}}_{n,i}^\top \mathbf{W}_V^{(t)} \mathbf{w}_O &= \tilde{\boldsymbol{\xi}}_{n,i}^\top \mathbf{W}_V^{(0)} \mathbf{w}_O + \rho_{V,n,i}^{(t)} \|\mathbf{w}_O\|_2^2,\end{aligned}$$

for  $i \in [M] \setminus \{1\}$  and  $n \in [N]$ .

We further denote the  $V_+^{(t)} := \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$ ,  $V_-^{(t)} := \tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$  and  $V_{n,i}^{(t)} := \tilde{\boldsymbol{\xi}}_{n,i}^\top \mathbf{W}_V^{(t)}$ .

With scalarized V in Definition 2, we can provide the dynamics of matrix  $\mathbf{W}_V^{(t)}$  by analyzing the update of coefficients  $\gamma^{(t)}$  as follows:

$$\begin{aligned}\gamma_{V,+}^{(t+1)} &= \gamma_{V,+}^{(t)} - \frac{\eta}{NM} \sum_{n \in S_+} \tilde{\ell}_n^{(t)} \left[ \sum_{l=1}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \varphi(\tilde{\mathbf{x}}_{n,l} \mathbf{W}_Q \mathbf{W}_K^\top (\tilde{\mathbf{X}}_n)^\top)_1 \right. \\ &\quad \left. + \sum_{i=2}^M \sum_{l=1}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle \varphi(\tilde{\mathbf{x}}_{n,l} \mathbf{W}_Q \mathbf{W}_K^\top (\tilde{\mathbf{X}}_n)^\top)_i \right].\end{aligned}$$

First, the perturbed signal and noise components break the orthogonality between signal and noise, resulting in the appearance of terms of the form  $\langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle$ . The upper bound of  $\langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle$  is given by  $(\|\boldsymbol{\mu}\|_2 + \tau)^2$ , while the upper bound of  $\langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle$  is  $(\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2)$ . Therefore, when the perturbation magnitude is small,  $\langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle$  dominates, and the training dynamics under adversarial training closely resemble those under clean training.

When bounding  $V_+^{(t)}$  (similar for  $V_-^{(t)}$ ), we can consider a special case  $\tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\mu}}^{(t)}$  and derive the bound via cumulative summation, yielding  $|V_+^{(t)}| \leq |V_+^{(0)}| + \sum_{s=0}^{t-1} |\gamma_{V,+}^{(s+1)} - \gamma_{V,+}^{(s)}| \cdot \|\mathbf{w}_O\|_2^2$  by Definition 2. In fact, we establish that for any  $\tilde{\boldsymbol{\mu}} \in B(\boldsymbol{\mu}, \tau)$ , the quantity  $\tilde{\boldsymbol{\mu}}^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$  admits a uniform upper bound.

Similarly, we can define vectorized  $\mathbf{W}_Q$  and  $\mathbf{W}_K$  as follows.

**Definition 3** (Vectorized Q & K). Let  $\mathbf{W}_Q^{(t)}$  and  $\mathbf{W}_K^{(t)}$  be the QK matrices of the ViT at the  $t$ -th iteration of gradient descent. Then we define the vectorized Q and vectorized K as follows

$$\begin{aligned}\mathbf{q}_+^{(t)} &= \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q^{(t)}, & \mathbf{q}_-^{(t)} &= \tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_Q^{(t)}, & \mathbf{q}_{n,i}^{(t)} &= \tilde{\boldsymbol{\xi}}_{n,i}^\top \mathbf{W}_Q^{(t)}, \\ \mathbf{k}_+^{(t)} &= \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_K^{(t)}, & \mathbf{k}_-^{(t)} &= \tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_K^{(t)}, & \mathbf{k}_{n,i}^{(t)} &= \tilde{\boldsymbol{\xi}}_{n,i}^\top \mathbf{W}_K^{(t)}\end{aligned}$$

for  $i \in [M] \setminus \{1\}$ ,  $n \in [N]$ .

With Definition 2 and 3, we can analyze the learning dynamics of the transformed coefficients rather than the original matrices.

## 5.2 EFFECTS OF PERTURBATION MAGNITUDE ON ATTENTION

Our second key technique is to analyze how the attention mechanism behaves under perturbations of varying magnitudes. First, we present the following lemma, which shows that the attention mechanism is robust to small perturbations.

**Lemma 4** (Informal). Under Condition 1, suppose the perturbation satisfies  $\tau \leq O(\frac{\|\boldsymbol{\mu}\|_2}{\log(d_n)})$  and  $t \geq \Omega\left(\frac{1}{\eta \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \log(6N^2M^2/\delta)}\right)$ . Then, there exists a universal constant  $C \leq e/2$  such that

$$\max_{\tilde{\mathbf{X}} \in B(\mathbf{X}, \tau)} \text{softmax}\left(\langle \mathbf{q}^{(t)}, \mathbf{k}^{(t)} \rangle\right) / \min_{\tilde{\mathbf{X}} \in B(\mathbf{X}, \tau)} \text{softmax}\left(\langle \mathbf{q}^{(t)}, \mathbf{k}^{(t)} \rangle\right) \leq C.$$

This lemma implies that the relative change in attention weights under perturbations is uniformly bounded; hence, attention computed on perturbed inputs closely matches that on the clean inputs. Thus, when considering the attention from the perturbed patch to another, it can be well-approximated by the attention from the clean signal to the clean noise.

Next, we present the following lemma, which characterizes the behavior of the attention mechanism under moderate perturbations.

**Lemma 5 (Informal).** *Under Condition 1, supposing the perturbation satisfies  $\omega(\frac{\|\boldsymbol{\mu}\|_2}{\log d_h}) \leq \tau \leq O(\boldsymbol{\mu})$ , for any  $t \geq 0$ , we have:*

$$\frac{1}{M} - o(1) \leq \text{softmax}(\langle \mathbf{q}^{(t)}, \mathbf{k}^{(t)} \rangle) \leq \frac{1}{M} + o(1)$$

This lemma demonstrates that, under moderate perturbations, the attention distribution remains invariant and stays close to its initial uniform form.

### 5.3 GENERALIZATION GUARANTEE

For a new data point  $(\mathbf{X}, y)$  generated from the distribution defined in Definition 1, we interpret the test error as the probability that the noise component dominates the model output. For adversarial test data with added perturbations, we need to bound the maximal distance between them and the corresponding clean test data. First, from the analysis in Section 5.1, we know that bounding  $V^{(t)} = \tilde{\mathbf{x}}^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$  implies that its deviation from the unperturbed counterpart  $\mathbf{x}^\top \mathbf{W}_V \mathbf{w}_O$  is at most  $|\langle \tilde{\mathbf{x}} - \mathbf{x}, \mathbf{W}_V \mathbf{w}_O \rangle| \leq \tau \|\mathbf{W}_V \mathbf{w}_O\|_2$ . Second, from the analysis in Section 5.2, we know that the attention component  $\text{softmax}(\langle \mathbf{q}, \mathbf{k} \rangle)$  exhibits robustness under perturbations of bounded magnitude. Combining these insights, we state the following lemma.

**Lemma 6 (Informal).** *Under Condition 1, if  $t \geq \Omega\left(\eta^{-1} \epsilon^{-1} \|\boldsymbol{\mu}\|_2^{-2} \|\mathbf{w}_O\|_2^{-2} \log^{-1}\left(\frac{6N^2 M^2}{\delta}\right)\right)$  and  $\tau \leq \frac{\|\boldsymbol{\mu}\|_2}{\log(d_h)}$ , we have:*

$$yf(\mathbf{X}, \theta(t)) - \min_{\tilde{\mathbf{X}} \in B(\mathbf{X}, \tau)} yf(\tilde{\mathbf{X}}, \theta(t)) \lesssim M \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2 \tau.$$

As a result, the robust test error can be interpreted as the probability of misclassification when the output of the ViT under clean test data is perturbed by its maximal adversarial deviation.

**Lemma 7 (Informal).** *Under the same conditions of Lemma 6, with high probability, we have for some constant  $c > 0$ :*

$$\begin{aligned} P(\exists \tilde{\mathbf{X}} \in B(\mathbf{X}, \tau) : yf(\tilde{\mathbf{X}}, \theta(t)) \leq 0) &= P(yf(\mathbf{X}, \theta(t)) + (yf(\mathbf{X}, \theta(t)) - \min_{\tilde{\mathbf{X}} \in B(\mathbf{X}, \tau)} yf(\tilde{\mathbf{X}}, \theta(t))) \leq 0) \\ &\leq \exp\left(-c \left(\frac{(V_+^{(t)} - V_-^{(t)}) - \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2 \tau}{\sigma_p \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2}\right)^2\right). \end{aligned}$$

Consequently, the robust test error differs from the clean test error only by an additive factor that scales with both the perturbation radius and the model complexity. More details are in Appendix.

## 6 EXPERIMENTS

### 6.1 EXPERIMENTAL SETUP

**Datasets.** For the experiments on synthetic data, we synthesize each sample following the distribution in Definition 1. Specifically, we first formalize the two signal vectors in Definition 1 as  $\boldsymbol{\mu}_+ = \|\boldsymbol{\mu}\|_2 \cdot [1, 0, \dots, 0]^\top$  and  $\boldsymbol{\mu}_- = \|\boldsymbol{\mu}\|_2 \cdot [0, 1, 0, \dots, 0]^\top$ , where the signal dimension  $d$  is set as 1024. Then, for each generated sample, the number  $M$  of tokens in this sample is set as 16 and every noise vector in this sample is drawn independently from the Gaussian distribution  $\boldsymbol{\xi}_i \sim \mathcal{N}(0, 0.4 \cdot \mathbb{I}_d)$ . We generate up to 22 samples for training and 100 samples for evaluation.

For the experiments on real-world data, we use a subset of MNIST where only samples with label ‘0’ or ‘1’ are used. To better simulate the feature learning setting in our theory, we transform each image to a vector with the following procedures: (1) Vectorize the image to a “signal” vector and

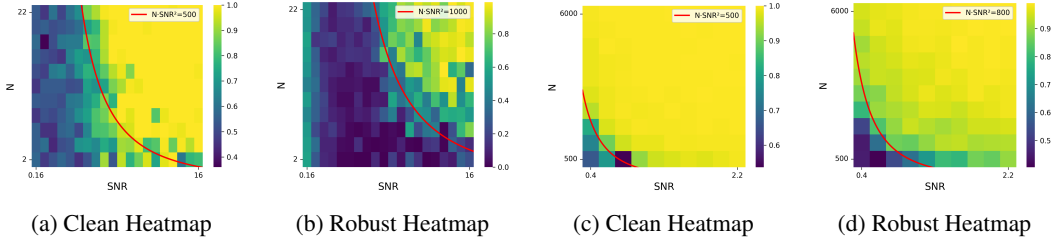


Figure 1: Clean and robust test accuracy under adversarial training across various signal-to-noise ratios (SNR) and sample sizes ( $N$ ). (a)&(b): results on synthetic data; (c)&(d): results on real-world data. High test accuracy is colored in yellow, whereas low test accuracy is colored in purple.

normalize its  $\ell_2$ -norm to 5. (2) Draw a noise vector following the Gaussian distribution  $\mathcal{N}(0, \mathbb{I}_d)$ , where  $d = 784$ . (3) Concatenate the signal and noise vectors to form a new vector as the sample adopted in our experiments.

**Model architectures.** For synthetic data experiments, we adopt the two-layer Transformer defined in (1), where both of the hidden dimension  $d_h$  and the value dimension  $d_v$  are set as 128. For real-world data experiments, We implement a ViT model consisting of two attention layers, each with four self-attention heads, followed by an MLP layer with ReLU activation. The hidden dimension of this ViT is set as 128. Model parameters in all experiments are initialized using PyTorch’s default method, followed by an additional scaling factor of  $1/16$ , matching the requirement of Condition (5).

**Model training & evaluation.** We use full-batch gradient descent to train all models in our experiments, with a learning rate of 0.1. Each model will be trained until its training loss falls below a target threshold 0.01. Besides, in each adversarial training step, we leverage projected gradient descent (PGD; Madry et al. 2017) to search adversarial examples with attack strength  $\tau/\|\mu\| = 0.05$  for 20 steps with a step size of  $0.2\tau$ , using per-token  $L_2$ -normalized gradient updates with projection onto the  $L_2$  ball of radius  $\tau$ . To assess the performance of trained models, we report both their clean and robust classification errors calculated on test datasets.

## 6.2 RESULTS ANALYSIS

**Phase Transition in benign overfitting.** We perform adversarial training with different number  $N$  of training data ranging from 2 to 22 and different SNR ranging from 0.16 to 16 on synthetic data. The clean and robust test accuracies of these models are collected and presented as heatmaps in Figures 1a and 1b. From the figure, we observe a clear decision boundary in the form  $N \cdot \text{SNR}^2 = \Omega(1)$  separates the high-accuracy and low-accuracy regions. This observation aligns well with our Theorem 2 that  $N \cdot \text{SNR}^2 = \Omega(1)$  is necessary for models in adversarial training to produce benign overfitting. We also conduct experiments on the real-world dataset MNIST, where similar conclusions can still be observed as shown in Figures 1c and 1d.

**Effects of signal-to-noise ratio and dataset size.** We then fix the number  $N$  of training data and plot curves of the training loss/robust test accuracy versus the training iteration under different SNR in Figures 2a and 2b. From them, we observe that as the training iteration increases, the model overfits to training data, but its robustness does not consistently increase unless the SNR is large. We also fix the SNR and plot similar curves with different number of training data in Figures 2c and 2d, where we find that even the model is overfitting, its robustness improves only when the training data number  $N$  is large. All these results indicate that benign overfitting can emerge in adversarial training only when the training data number  $N$  and the data SNR are both not too small, which coincides with the requirements that  $N \cdot \text{SNR}^2 = \Omega(1)$  or  $N \cdot \text{SNR}^2 = \Omega(\frac{1}{\epsilon})$  in our Theorem 2.

We provide the following additional results in the appendix: Empirical validation of different learning dynamics under varying perturbation radii in Appendix B.1; Additional experiments on MNIST, CIFAR-10, and Tiny-ImageNet with the state-of-the-art APGD attack (Croce & Hein, 2020) in Appendix B.2; Additional experiments on multi-norm attacks in Appendix B.3; and Additional experiments on realistic ViT-base models(Dosovitskiy et al., 2020) in Appendix B.4.

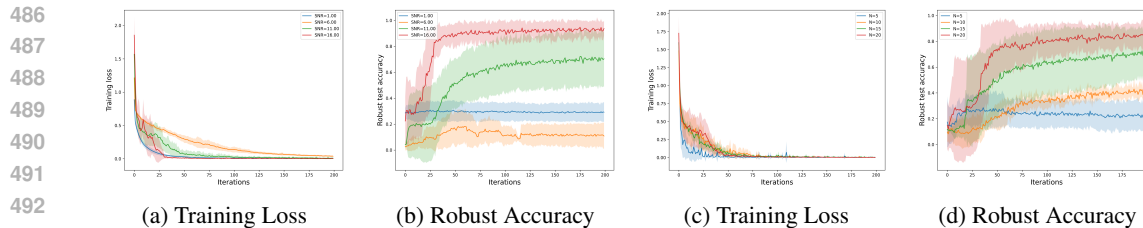


Figure 2: (a)&(b): Curves of robust training loss and robust test accuracy versus training iteration under fixed training data number  $N = 22$ . (c)&(d): Curves of robust training loss and robust test accuracy versus training iteration under fixed data SNR = 12.

## 7 CONCLUSION

Our paper presents the first comprehensive theoretical analysis of the generalization behavior after adversarial training on a two-layer Vision Transformer. We demonstrate that, under appropriate relationships between signal-to-noise ratio and perturbation magnitude, adversarially trained ViTs can interpolate the training data with vanishing robust loss while still achieving small robust test error. Our analysis reveals three perturbation regimes that clarify how adversarial training shapes the learning dynamics of attention heads, from clean-like behavior to linear collapse and eventual failure beyond the benign regime. Empirical results on synthetic data and MNIST corroborate the theory, aligning closely with the predicted conditions under which robust benign overfitting emerges. Empirically, experiments on synthetic data and MNIST corroborate the theory, matching the predicted conditions under which robust benign overfitting arises.

## REFERENCES

- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- Usman Anwar, Johannes von Oswald, Louis Kirsch, David Krueger, and Spencer Frei. Understanding in-context learning of linear models in transformers through an adversarial lens. *Transactions on Machine Learning Research*.
- Yutong Bai, Jieru Mei, Alan L Yuille, and Cihang Xie. Are transformers more robust than cnns? *Advances in neural information processing systems*, 34:26831–26843, 2021.
- Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. In *International conference on machine learning*, pp. 541–549. PMLR, 2018.
- Yuan Cao, Quanquan Gu, and Mikhail Belkin. Risk bounds for over-parameterized maximum margin classification on sub-gaussian mixtures. *arXiv preprint arXiv:2104.13628*, 2021.
- Yuan Cao, Zixiang Chen, Misha Belkin, and Quanquan Gu. Benign overfitting in two-layer convolutional neural networks. *Advances in neural information processing systems*, 35:25237–25250, 2022.
- Niladri S Chatterji and Philip M Long. Deep linear networks can benignly overfit when shallow ones do. *Journal of Machine Learning Research*, 24(117):1–39, 2023.
- Jinghui Chen, Yuan Cao, and Quanquan Gu. Benign overfitting in adversarially robust linear classification. In *Uncertainty in Artificial Intelligence*, pp. 313–323. PMLR, 2023.
- Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International conference on machine learning*, pp. 2206–2216. PMLR, 2020.

- 540 Meng Ding, Kaiyi Ji, Di Wang, and Jinhui Xu. Understanding forgetting in continual learning with  
541 linear regression. *arXiv preprint arXiv:2405.17583*, 2024a.
- 542
- 543 Meng Ding, Mingxi Lei, Liyang Zhu, Shaowei Wang, Di Wang, and Jinhui Xu. Revisiting differ-  
544 entially private relu regression. *Advances in Neural Information Processing Systems*, 37:55470–  
545 55506, 2024b.
- 546 Meng Ding, Mingxi Lei, Shaopeng Fu, Di Wang, and Jinhui Xu. Understanding private learning  
547 from feature perspective. In *ICLR 2025 Workshop on Navigating and Addressing Data Problems*  
548 *for Foundation Models*, 2025.
- 549
- 550 Hiskias Dingeto and Juntae Kim. Comparative study of adversarial defenses: Adversarial training  
551 and regularization in vision transformers and cnns. *Electronics*, 13(13):2534, 2024.
- 552
- 553 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
554 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An  
555 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*  
556 *arXiv:2010.11929*, 2020.
- 557 Spencer Frei and Gal Vardi. Trained transformer classifiers generalize and exhibit benign overfitting  
558 in-context. *arXiv preprint arXiv:2410.01774*, 2024.
- 559
- 560 Spencer Frei, Niladri S Chatterji, and Peter Bartlett. Benign overfitting without linearity: Neural  
561 network classifiers trained by gradient descent for noisy linear data. In *Conference on Learning*  
562 *Theory*, pp. 2668–2703. PMLR, 2022.
- 563
- 564 Spencer Frei, Gal Vardi, Peter Bartlett, and Nathan Srebro. Benign overfitting in linear classifiers  
565 and leaky relu networks from kkt conditions for margin maximization. In *The Thirty Sixth Annual*  
566 *Conference on Learning Theory*, pp. 3173–3228. PMLR, 2023.
- 567
- 568 Shaopeng Fu and Di Wang. Theoretical analysis of robust overfitting for wide dnns: An ntk ap-  
569 proach. *arXiv preprint arXiv:2310.06112*, 2023.
- 570
- 571 Shaopeng Fu, Liang Ding, Jingfeng Zhang, and Di Wang. Short-length adversarial training helps  
572 llms defend long-length jailbreak attacks: Theoretical and empirical evidence. *arXiv preprint*  
573 *arXiv:2502.04204*, 2025.
- 574
- 575 Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1.  
576 MIT Press, 2016.
- 577
- 578 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial  
579 examples. *arXiv preprint arXiv:1412.6572*, 2014.
- 580
- 581 Bhavna Gopal, Huanrui Yang, Mark Horton, and Yiran Chen. Safer: Sharpness aware layer-selective  
582 finetuning for enhanced robustness in vision transformers. *arXiv preprint arXiv:2501.01529*,  
583 2025.
- 584
- 585 Andi Han, Wei Huang, Yuan Cao, and Difan Zou. On the feature learning in diffusion models. *arXiv*  
586 *preprint arXiv:2412.01021*, 2024.
- 587
- 588 Yifan Hao and Tong Zhang. The surprising harmfulness of benign overfitting for adversarial robust-  
589 ness. *arXiv preprint arXiv:2401.12236*, 2024.
- 590
- 591 Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-  
592 dimensional ridgeless least squares interpolation. *Annals of statistics*, 50(2):949, 2022.
- 593
- 594 Charles Herrmann, Kyle Sargent, Lu Jiang, Ramin Zabih, Huiwen Chang, Ce Liu, Dilip Krishnan,  
595 and Deqing Sun. Pyramid adversarial training improves vit performance. In *Proceedings of the*  
596 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 13419–13429, 2022.
- 597
- 598 Lijie Hu, Yixin Liu, Ninghao Liu, Mengdi Huai, Lichao Sun, and Di Wang. Seat: stable and ex-  
599 plainable attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,  
600 pp. 12907–12915, 2023.

- 594 Lijie Hu, Yixin Liu, Ninghao Liu, Mengdi Huai, Lichao Sun, and Di Wang. Improving interpretation  
595 faithfulness for vision transformers. In *Forty-first International Conference on Machine Learning*,  
596 2024.
- 597 Lijie Hu, Songning Lai, Yuan Hua, Shu Yang, Jingfeng Zhang, and Di Wang. Stable vision concept  
598 transformers for medical diagnosis. In *Joint European Conference on Machine Learning and*  
599 *Knowledge Discovery in Databases*, pp. 317–332. Springer, 2025a.
- 600 Lijie Hu, Xinhai Wang, Yixin Liu, Ninghao Liu, Mengdi Huai, Lichao Sun, and Di Wang. To-  
601 wards stable and explainable attention mechanisms. *IEEE Transactions on Knowledge and Data*  
602 *Engineering*, 2025b.
- 603 Chashi Mahiul Islam, Samuel Jacob Chacko, Mao Nishino, and Xiuwen Liu. Mechanistic un-  
604 derstandings of representation vulnerabilities and engineering robust vision transformers. *arXiv*  
605 *preprint arXiv:2502.04679*, 2025.
- 606 Samy Jelassi and Yuanzhi Li. Towards understanding how momentum improves generalization in  
607 deep learning. In *International Conference on Machine Learning*, pp. 9965–10040. PMLR, 2022.
- 608 Jiarui Jiang, Wei Huang, Miao Zhang, Taiji Suzuki, and Liqiang Nie. Unveil benign overfitting for  
609 transformer in vision: Training dynamics, convergence, and generalization. *Advances in Neural*  
610 *Information Processing Systems*, 37:135464–135625, 2024.
- 611 Yiwen Kou, Zixiang Chen, Yuanzhou Chen, and Quanquan Gu. Benign overfitting in two-layer  
612 relu convolutional neural networks. In *International conference on machine learning*, pp. 17615–  
613 17659. PMLR, 2023.
- 614 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-  
615 lutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 616 Soichiro Kumano, Hiroshi Kera, and Toshihiko Yamasaki. Adversarially pretrained transformers  
617 may be universally robust in-context learners. *arXiv preprint arXiv:2505.14042*, 2025.
- 618 Zhu Li, Zhi-Hua Zhou, and Arthur Gretton. Towards an understanding of benign overfitting in neural  
619 networks. *arXiv preprint arXiv:2106.03212*, 2021.
- 620 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.  
621 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*  
622 *IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 623 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
624 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,  
625 2017.
- 626 Song Mei and Andrea Montanari. The generalization error of random features regression: Precise  
627 asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*,  
628 75(4):667–766, 2022.
- 629 Yichuan Mo, Dongxian Wu, Yifei Wang, Yiwen Guo, and Yisen Wang. When adversarial training  
630 meets vision transformers: Recipes from training to architecture. *Advances in Neural Information*  
631 *Processing Systems*, 35:18599–18611, 2022.
- 632 Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John C Duchi, and Percy Liang. Adversarial  
633 training can hurt generalization. *arXiv preprint arXiv:1906.06032*, 2019.
- 634 Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Ad-  
635 versarially robust generalization requires more data. *Advances in neural information processing*  
636 *systems*, 31, 2018.
- 637 Rulin Shao, Zhouxing Shi, Jinfeng Yi, Pin-Yu Chen, and Cho-Jui Hsieh. On the adversarial robust-  
638 ness of vision transformers. *arXiv preprint arXiv:2103.15670*, 2021.
- 639 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow,  
640 and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

648 Alexander Tsigler and Peter L. Bartlett. Benign overfitting in ridge regression. *Journal of Machine*  
649 *Learning Research*, 24(123):1–76, 2023.  
650

651 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,  
652 volume 47. Cambridge university press, 2018.

653 Ke Wang, Vidya Muthukumar, and Christos Thrampoulidis. Benign overfitting in multiclass classi-  
654 fication: All roads lead to interpolation. In *Advances in Neural Information Processing Systems*  
655 (*NeurIPS*), volume 34, pp. 14777–14790, 2021.  
656

657 Yunjuan Wang, Kaibo Zhang, and Raman Arora. Benign overfitting in adversarial training of neural  
658 networks. In *Forty-first International Conference on Machine Learning*, 2024.

659 Boxi Wu, Jindong Gu, Zhifeng Li, Deng Cai, Xiaofei He, and Wei Liu. Towards efficient adversarial  
660 training on vision transformers. In *European Conference on Computer Vision*, pp. 307–325.  
661 Springer, 2022.  
662

663 Jiancong Xiao, Yanbo Fan, Ruoyu Sun, Jue Wang, and Zhi-Quan Luo. Stability analysis and gen-  
664 eralization bounds of adversarial training. *Advances in Neural Information Processing Systems*,  
665 35:15446–15459, 2022.

666 Xilie Xu, Keyi Kong, Ning Liu, Lizhen Cui, Di Wang, Jingfeng Zhang, and Mohan Kankanhalli. An  
667 llm can fool itself: A prompt-based adversarial attack. *arXiv preprint arXiv:2310.13345*, 2023.  
668

669 Shufei Zhang, Zhuang Qian, Kaizhu Huang, Qiufeng Wang, Rui Zhang, and Xinpeng Yi. Towards  
670 better robust generalization with shift consistency regularization. In *International conference on*  
671 *machine learning*, pp. 12524–12534. PMLR, 2021.

672 Mo Zhou and Rong Ge. Implicit regularization leads to benign overfitting for sparse linear regres-  
673 sion. In *International Conference on Machine Learning*, pp. 42543–42573. PMLR, 2023.  
674

675 Difan Zou, Yuan Cao, Yuanzhi Li, and Quanquan Gu. The benefits of mixup for feature learning. In  
676 *International Conference on Machine Learning*, pp. 43423–43479. PMLR, 2023.  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

## A LLMs USAGE IN THE PAPER

LLMs were used only occasionally to help polish the writing (propose new words, grammar and spelling correction). All technical ideas, experimental designs, analyses, conclusions, writing were developed and carried out entirely by the authors. The authors have full responsibility for the final text.

## B ADDITIONAL EXPERIMENTS

### B.1 EMPIRICAL VALIDATION OF THEORETICAL REGIMES

In this section, we follow a similar setting as the experiments on synthetic data in Section 6, and add more experiments about training dynamics with different perturbation  $\tau$  radius.

#### Experiments setting.

We focus on tracking the dynamics of attention entropy, training loss, and  $W_V$  norm under benign overfitting regime with different perturbation  $\tau$  radius, and the key parameters are as follows:

- $N = 25$
- $SNR = 16$
- $M = 2$
- $d = 1024$
- $d_h = 512$
- $d_v = 512$
- $\sigma_p = 0.05$
- $\frac{\tau}{\|\mu\|_2} = 0.02, 0.1, 0.5$

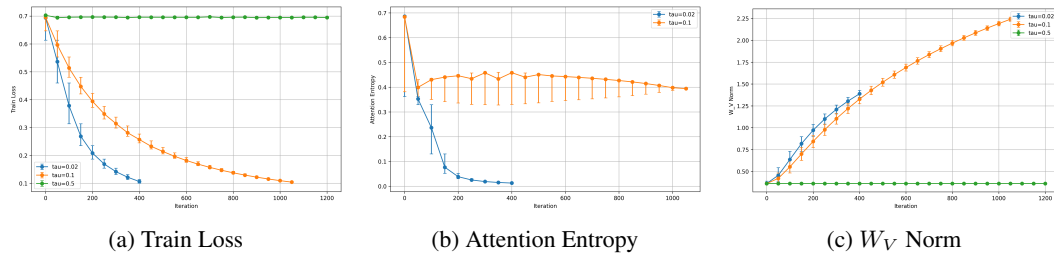


Figure 3: Training dynamics with different perturbation  $\tau$  radius: attention entropy, training loss, and  $W_V$  norm.

#### Experiments results.

1. In Figure 3a, when  $\frac{\tau}{\|\mu\|_2} = 0.02$  or  $0.1$ , the adversarial training loss converges to zero, indicating that the model successfully interpolates all noise-corrupted training samples, consistent with the benign overfitting behavior in Theorem 2. In contrast, when  $\frac{\tau}{\|\mu\|_2} = 0.5$ , the adversarial training loss fails to decrease, aligning with the non-convergence regime characterized in Theorem 3. Moreover, the case  $\frac{\tau}{\|\mu\|_2} = 0.02$  exhibits a faster convergence rate than the  $\frac{\tau}{\|\mu\|_2} = 0.1$  setting, highlighting the role of the attention mechanism in accelerating convergence, which is in agreement with our theoretical predictions in Theorem 2.

2. In Figure 3b, when  $\frac{\tau}{\|\mu\|_2} = 0.02$ , the attention entropy decreases to nearly zero, indicating that the attention mechanism correctly concentrates on the signal patch. This behavior is consistent with Case 1 of Theorem 2, where the perturbation level is sufficiently small for the model to recover the underlying signal structure.

In contrast, when  $\frac{\tau}{\|\mu\|_2} = 0.1$ , the attention entropy fails to decrease and instead remains high, demonstrating that moderate perturbations hinder the learning of attention weights. As a result, the

attention distribution remains nearly uniform rather than focusing on the signal patch. This phenomenon aligns with Case 2 of Theorem 2, where the perturbation magnitude prevents the attention mechanism from identifying the true signal.

3. In Figure 3c, when  $\frac{\tau}{\|\mu\|_2} = 0.1$ , the  $\|W_V\|_2$  norm exhibits the largest growth. This behavior is consistent with Case 2 of Theorem 2, where the ViT effectively collapses into a linear model and the value projection  $W_V$  becomes the dominant component driving the learning dynamics.

In contrast, for  $\frac{\tau}{\|\mu\|_2} = 0.02$ , the attention mechanism remains effective, so only mild updates to  $W_V$  are required for the model to fit the noisy training data and achieve benign overfitting. When  $\frac{\tau}{\|\mu\|_2} = 0.5$ , the perturbation is too large for the model to learn meaningful structure, resulting in  $W_V$  failing to make progress during training.

### B.2 ADDITIONAL EXPERIMENTS ON MNIST, CIFAR-10 AND TINY-IMAGENET WITH APGD

In this section, we follow the same experimental setup as in Section 6 for the MNIST dataset and further extend our evaluation by conducting additional experiments on both MNIST, CIFAR-10 and Tiny-ImageNet under APGD attacks. These results demonstrate that our theoretical insights continue to hold for larger and more complex datasets, as well as under stronger adversarial attacks.

**Experiments setting.** We conduct adversarial training using the state-of-the-art APGD attack model. For the MNIST dataset, we consider an attack strength of  $\frac{\tau}{\|\mu\|_2} = 0.05$ , and vary the number of training samples  $N$  from 1000 to 6000 and SNR from 0.4 to 2. For the CIFAR-10 dataset, we set a weaker attack strength of  $\frac{\tau}{\|\mu\|_2} = 0.01$ , and vary the number of training samples  $N$  from 1000 to 10000 and the SNR from 0.4 to 10. For the Tiny-ImageNet dataset, we set a weaker attack strength of  $\frac{\tau}{\|\mu\|_2} = 0.10$ , and vary the number of training samples  $N$  from 100 to 1000 and the SNR from 0.4 to 10.

**Experiments results.** The clean and robust test accuracies on MNIST, CIFAR-10 and Tiny-ImageNet are collected and presented as heatmaps in Figures 4, 5 and 6. In both figures, we observe a clear phase transition phenomenon. Moreover, as both the sample size  $N$  and the SNR increase, the clean and robust test error consistently decreases. This behavior is fully aligned with our theoretical analysis as well as the empirical findings reported in previous experiments.

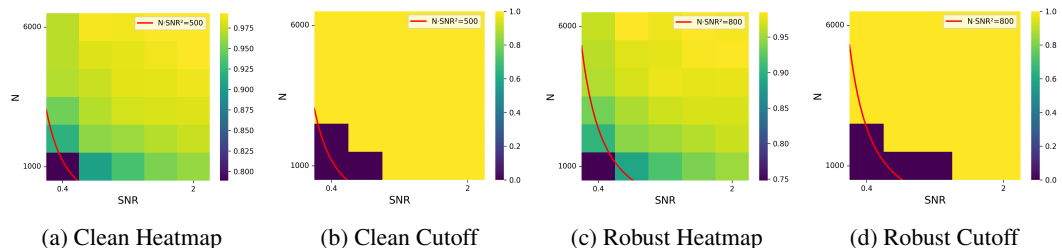


Figure 4: Clean and robust test accuracy heatmaps on MNIST with APGD attack across various signal-to-noise ratios (SNR) and sample sizes ( $N$ ). (b)&(d) are a heatmap that applies a cutoff value 0.93.

### B.3 ADDITIONAL EXPERIMENTS ON MULTI-NORM ATTACKS

In this section, we follow the same experimental setup as in Section 6 for the MNIST dataset and further extend our evaluation by conducting additional experiments on both MNIST and CIFAR-10 under multi-norm attacks. These results demonstrate that although the model’s robust test accuracy decreases under multi-norm attacks, our theoretical insights continue to hold.

**Experiments setting.** We perform adversarial training under a multi-norm PGD attack model spanning  $l_1, l_2$  and  $l_\infty$  perturbations.

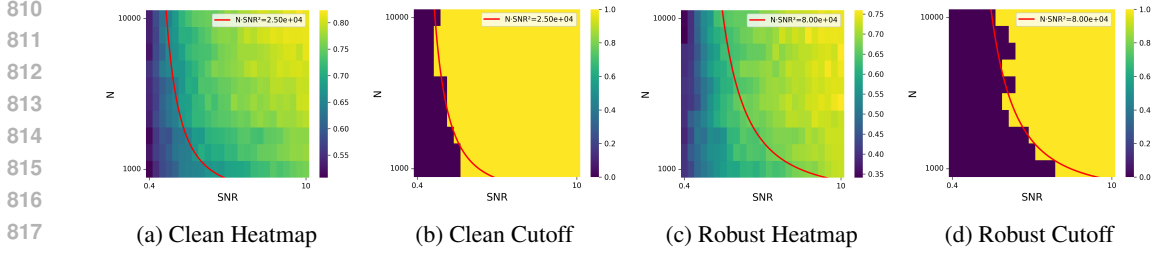


Figure 5: Clean and robust test accuracy heatmaps on CIFAR-10 with APGD attack across various signal-to-noise ratios (SNR) and sample sizes ( $N$ ). (b)&(d) are a heatmap that applies a cutoff value 0.65.

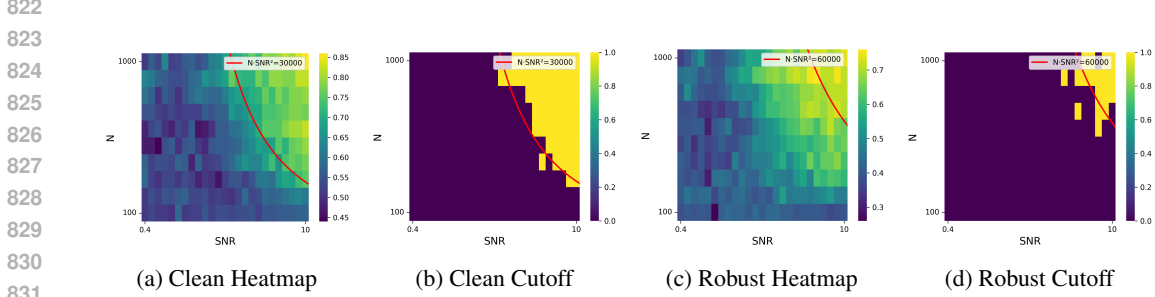


Figure 6: Clean and robust test accuracy heatmaps on Tiny Imagenet with APGD attack across various signal-to-noise ratios (SNR) and sample sizes ( $N$ ). (b)&(d) are a heatmap that applies a cutoff value 0.70.

For the **MNIST** dataset, we consider a base attack strength of  $\text{eps} = \frac{\tau_2}{\|\mu\|_2} = 0.05$ , and set  $(\frac{\tau_1}{\|\mu\|_2}, \frac{\tau_2}{\|\mu\|_2}, \frac{\tau_\infty}{\|\mu\|_2}) = (\text{eps} * 20, \text{eps}, \text{eps}/30)$ . We vary the number of training samples  $N$  from 1000 to 6000 and SNR from 0.4 to 2.

For the **CIFAR-10** dataset, we set a weaker base attack strength of  $\text{eps} = \frac{\tau_2}{\|\mu\|_2} = 0.01$ , and set  $(\frac{\tau_1}{\|\mu\|_2}, \frac{\tau_2}{\|\mu\|_2}, \frac{\tau_\infty}{\|\mu\|_2}) = (\text{eps} * 20, \text{eps}, \text{eps}/30)$ . We vary the number of training samples  $N$  from 1000 to 10000 and the SNR from 0.4 to 10.

**Experiments results.** The clean and robust test accuracies on MNIST and CIFAR-10 are collected and presented as heatmaps in Figures 7. We observe that although the robust test accuracy decreases under multi-norm attacks, a phase transition phenomenon still persists, and increasing either the sample size  $N$  or the SNR further reduces both clean and robust test error in a manner fully aligned with our theoretical results.

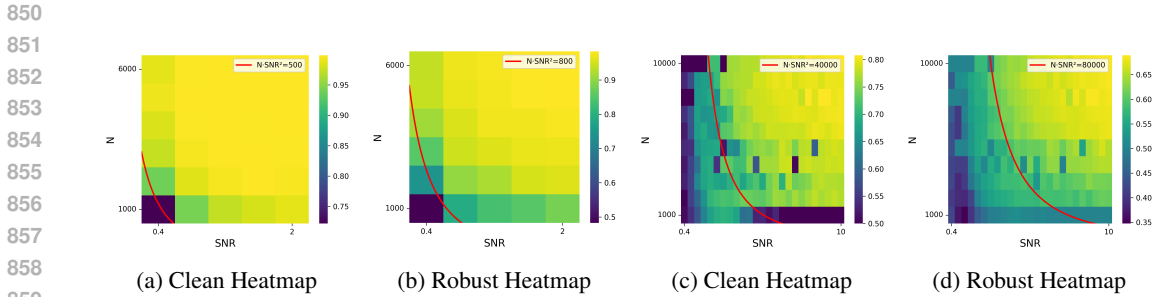


Figure 7: Clean and robust test accuracy under multi-norm attack adversarial training across various signal-to-noise ratios (SNR) and sample sizes ( $N$ ). (a)&(b): results on MNIST data; (c)&(d): results on CIFAR-10 data.

B.4 ADDITIONAL EXPERIMENTS ON REALISTIC ViT

In this section, we conduct real-world experiments on image classification benchmarks, including MNIST, CIFAR-10, and Tiny-ImageNet, using a realistic ViT architecture (Dosovitskiy et al., 2020). The results show that when our model is scaled from the simplified two-layer ViT to a full-fledged ViT model, our theoretical insights continue to hold.

**Experiments setting.** We adopt google/vit-base-patch16-224-in21k (Dosovitskiy et al., 2020) as the backbone model to extend our analysis to real-world scenarios. To align with our theoretical setting, we freeze all parameters except for the QKV matrices in the final attention layer. This setup effectively treats all preceding layers as a fixed feature-extraction encoder, whose output serves as the input to the last Transformer layer.

We conduct adversarial training on MNIST, CIFAR-10, and Tiny-ImageNet, using PGD as the threat model with a perturbation radius of  $\frac{\|\mu\|_2}{20}$  and 5 attack steps.

**Experiments results.** The clean and robust test accuracies of ViT-base on MNIST, CIFAR-10 and Tiny-ImageNet are collected and presented as heatmaps in Figures 8. We observe that when the model is scaled from the simplified two-layer ViT to a realistic ViT architecture, a phase transition phenomenon still persists, and increasing either the sample size  $N$  or the SNR further reduces both clean and robust test error in a manner fully aligned with our theoretical results.

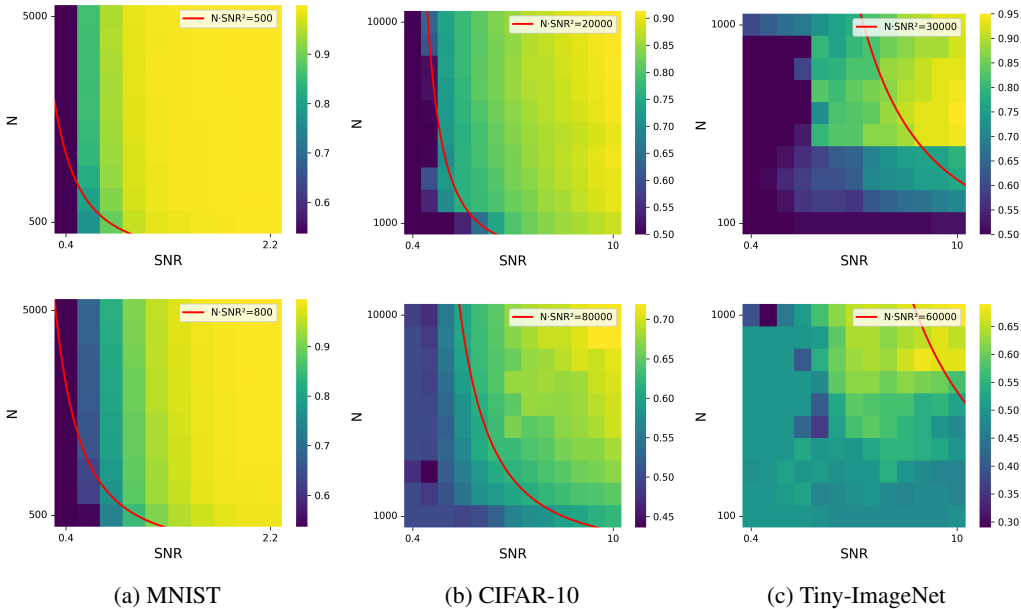


Figure 8: Clean and robust test accuracy of ViT-base under adversarial training across various signal-to-noise ratios (SNR) and sample sizes (N). Top row: clean test accuracy. Bottom row: robust test accuracy.

C EXTENSION TO MHA

We let the parameters be  $\theta := \{(W_{Q,h}, W_{K,h}, W_{V,h})\}_{h=1}^H$ , where  $W_{Q,h}, W_{K,h} \in \mathbb{R}^{d \times d_h}$  and  $W_{V,h} \in \mathbb{R}^{d \times d_v}$  for each  $h \in [H]$ . Here  $H$  denotes the number of attention heads, which we treat as a fixed constant. Under this parameterization, the network can be written as:

$$f(\mathbf{X}, \theta) = \sum_{h=1}^H f_h(\mathbf{X}, \theta)$$

where,

$$f_h(\mathbf{X}, \theta) = \frac{1}{M} \sum_{l=1}^M \varphi(\mathbf{x}_l^\top \mathbf{W}_{Q,h} \mathbf{W}_{K,h}^\top \mathbf{X}^\top) \mathbf{X} \mathbf{W}_{V,h} \mathbf{w}_O.$$

The gradients in the multi-head attention module,  $\frac{\partial f_h}{\partial W_{K,h}}$ ,  $\frac{\partial f_h}{\partial W_{Q,h}}$ ,  $\frac{\partial f_h}{\partial W_{V,h}}$ , remain unchanged. However, the gradient of the loss with respect to the output of each single head, i.e.,  $\frac{\partial \ell}{\partial f_h}$ , does change.

Intuitively, the model output increases by approximately an  $H$ -fold factor, which causes the scale of the loss  $\ell'$  to decrease accordingly.

More concretely, following our analysis of the signal attention head,  $\ell'^{(t)} = \frac{1}{M} \pm o(1)$  stay when  $t \leq T_2 = \Theta\left(\frac{1}{\eta(\|\mu\|_2 + \tau)^2 \|w_O\|_2^2}\right)$ . Thus, this implies  $f_h^{(T_2)}(\mathbf{X}, \theta) = o(1)$ . The  $H$ -fold increase in the multi-head model outputs does not alter this result, so the effect of the changes in  $\frac{\partial \ell}{\partial f_h}$  can be ignored.

Therefore, under the MHA setting, the training dynamics of the model still follow those of the single-head attention case, and our conclusions remain unchanged.

## D BASIC CALCULATION

### D.1 NOTION

### D.2 UPDATE RULES

We fix a universal perturbed input  $\widetilde{\mathbf{X}} = [\widetilde{\mu}, \widetilde{\xi}_{n,2}, \dots, \widetilde{\xi}_{n,M}]$ , where  $\widetilde{\mu}_+ \in B(\mu_+, \tau)$ ,  $\widetilde{\mu}_- \in B(\mu_-, \tau)$ , and  $\widetilde{\xi}_{n,i} \in B(\xi_{n,i}, \tau)$  are chosen once and remain fixed for all iterations  $t$ . These perturbations are universal and do not correspond to iteration-specific adversarial examples.

**Definition 4** (Scalarized V). *Let  $\mathbf{W}_V^{(t)}$  be the V matrix of the ViT at the  $t$ -th iteration of adversarial training. Then there exist coefficients  $\gamma_{V,+}^{(t)}, \gamma_{V,-}^{(t)}, \rho_{V,n,i}^{(t)}$  such that*

$$\begin{aligned} \widetilde{\mu}_+^\top \mathbf{W}_V^{(t)} \mathbf{w}_O &= \widetilde{\mu}_+^\top \mathbf{W}_V^{(0)} \mathbf{w}_O + \gamma_{V,+}^{(t)} \|\mathbf{w}_O\|_2^2, \\ \widetilde{\mu}_-^\top \mathbf{W}_V^{(t)} \mathbf{w}_O &= \widetilde{\mu}_-^\top \mathbf{W}_V^{(0)} \mathbf{w}_O + \gamma_{V,-}^{(t)} \|\mathbf{w}_O\|_2^2, \\ \widetilde{\xi}_{n,i}^\top \mathbf{W}_V^{(t)} \mathbf{w}_O &= \widetilde{\xi}_{n,i}^\top \mathbf{W}_V^{(0)} \mathbf{w}_O + \rho_{V,n,i}^{(t)} \|\mathbf{w}_O\|_2^2 \end{aligned}$$

for  $i \in [M] \setminus \{1\}, n \in [N]$ .

We further denote the  $V_+^{(t)} := \widetilde{\mu}_+^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$ ,  $V_-^{(t)} := \widetilde{\mu}_-^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$  and  $V_{n,i}^{(t)} := \widetilde{\xi}_{n,i}^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$  with time-independent perturbations, and  $\widetilde{V}_+^{(t)} := \widetilde{\mu}_+^{(t)\top} \mathbf{W}_V^{(t)} \mathbf{w}_O$ ,  $\widetilde{V}_-^{(t)} := \widetilde{\mu}_-^{(t)\top} \mathbf{W}_V^{(t)} \mathbf{w}_O$  and  $\widetilde{V}_{n,i}^{(t)} := \widetilde{\xi}_{n,i}^{(t)\top} \mathbf{W}_V^{(t)} \mathbf{w}_O$  with time-dependent perturbations, where  $\widetilde{\mu}_+^{(t)}, \widetilde{\mu}_-^{(t)}, \widetilde{\xi}_{n,i}^{(t)}$  is the adversarial sample at  $t$ -th iteration. We refer to it as scalarized V.

Similarly, we define the vectorized queries and keys as  $\mathbf{q}^{(t)}, \mathbf{k}^{(t)}$  with time-independent perturbations, and as  $\widetilde{\mathbf{q}}^{(t)}, \widetilde{\mathbf{k}}^{(t)}$  with time-dependent perturbations.

**Definition 5** (Vectorized Q & K). *Let  $\mathbf{W}_Q^{(t)}$  and  $\mathbf{W}_K^{(t)}$  be the QK matrices of the ViT at the  $t$ -th iteration of adversarial training. Then we define the vectorized Q and vectorized K as follows*

$$\begin{aligned} \mathbf{q}_+^{(t)} &= \widetilde{\mu}_+^\top \mathbf{W}_Q^{(t)}, & \mathbf{q}_-^{(t)} &= \widetilde{\mu}_-^\top \mathbf{W}_Q^{(t)}, & \mathbf{q}_{n,i}^{(t)} &= \widetilde{\xi}_{n,i}^\top \mathbf{W}_Q^{(t)}, \\ \mathbf{k}_+^{(t)} &= \widetilde{\mu}_+^\top \mathbf{W}_K^{(t)}, & \mathbf{k}_-^{(t)} &= \widetilde{\mu}_-^\top \mathbf{W}_K^{(t)}, & \mathbf{k}_{n,i}^{(t)} &= \widetilde{\xi}_{n,i}^\top \mathbf{W}_K^{(t)}, \\ \widetilde{\mathbf{q}}_+^{(t)} &= \widetilde{\mu}_+^{(t)\top} \mathbf{W}_Q^{(t)}, & \widetilde{\mathbf{q}}_-^{(t)} &= \widetilde{\mu}_-^{(t)\top} \mathbf{W}_Q^{(t)}, & \widetilde{\mathbf{q}}_{n,i}^{(t)} &= \widetilde{\xi}_{n,i}^{(t)\top} \mathbf{W}_Q^{(t)}, \\ \widetilde{\mathbf{k}}_+^{(t)} &= \widetilde{\mu}_+^{(t)\top} \mathbf{W}_K^{(t)}, & \widetilde{\mathbf{k}}_-^{(t)} &= \widetilde{\mu}_-^{(t)\top} \mathbf{W}_K^{(t)}, & \widetilde{\mathbf{k}}_{n,i}^{(t)} &= \widetilde{\xi}_{n,i}^{(t)\top} \mathbf{W}_K^{(t)} \end{aligned}$$

for  $i \in [M] \setminus \{1\}, n \in [N]$ .

Table 1: Notations

Symbols	Definitions
$x_{n,i}$	the $i$ -th token in the $n$ -th training sample if $i \in [M] \setminus \{1\}$ , $x_{n,i} = \xi_{n,i}$ .
$\varphi_{n,i}^{(t)}$	the $i$ -th row of attention for the $n$ -th sample, i.e., $\varphi_{n,i}^{(t)} := \varphi(x_{n,i}^\top \mathbf{W}_Q^{(t)} \mathbf{W}_K^{(t)\top} X_n^\top)$
$S_+, S_-$	the training samples with +1 labels and -1 labels, i.e., $S_+ := \{n \in [N] : y_n = 1\}$ , $S_- := \{n \in [N] : y_n = -1\}$
$\mathbf{q}_+^{(t)}, \mathbf{q}_-^{(t)}, \mathbf{q}_{n,i}^{(t)}$	vectorized Q, defined as $\mathbf{q}_+^{(t)} = \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q^{(t)}$ , $\mathbf{q}_-^{(t)} = \tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_Q^{(t)}$ , $\mathbf{q}_{n,i}^{(t)} = \tilde{\boldsymbol{\xi}}_{n,i}^\top \mathbf{W}_Q^{(t)}$
$\mathbf{k}_+^{(t)}, \mathbf{k}_-^{(t)}, \mathbf{k}_{n,i}^{(t)}$	vectorized K, defined as $\mathbf{k}_+^{(t)} = \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_K^{(t)}$ , $\mathbf{k}_-^{(t)} = \tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_K^{(t)}$ , $\mathbf{k}_{n,i}^{(t)} = \tilde{\boldsymbol{\xi}}_{n,i}^\top \mathbf{W}_K^{(t)}$
$V_+^{(t)}, V_-^{(t)}, V_{n,i}^{(t)}$	scalarized V, defined as $V_+^{(t)} := \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$ , $V_-^{(t)} := \tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$ , $V_{n,i}^{(t)} := \tilde{\boldsymbol{\xi}}_{n,i}^\top \mathbf{W}_V^{(t)} \mathbf{w}_O$
$\alpha_{\pm,\pm}^{(t)}, \alpha_{n,\pm,i}^{(t)}$	linear combinations coefficients for the dynamics of $\mathbf{q}_+^{(t)}$ and $\mathbf{q}_-^{(t)}$ , i.e., $\mathbf{q}_{\pm}^{(t+1)} - \mathbf{q}_{\pm}^{(t)} = \alpha_{\pm,\pm}^{(t)} \mathbf{k}_{\pm}^{(t)} + \sum_{n \in S_{\pm}} \sum_{i=2}^M \alpha_{n,\pm,i}^{(t)} \mathbf{k}_{n,i}^{(t)}$
$\alpha_{n,i,\pm}^{(t)}, \alpha_{n,i,n',i'}^{(t)}$	linear combinations coefficients for the dynamics of $\mathbf{q}_{n,i}^{(t)}$ , i.e., $\mathbf{q}_{n,i}^{(t+1)} - \mathbf{q}_{n,i}^{(t)} = \alpha_{n,i,+}^{(t)} \mathbf{k}_+^{(t)} + \alpha_{n,i,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n'=1}^N \sum_{i'=2}^M \alpha_{n,i,n',i'}^{(t)} \mathbf{k}_{n',i'}^{(t)}$
$\beta_{\pm,\pm}^{(t)}, \beta_{n,\pm,i}^{(t)}$	linear combinations coefficients for the dynamics of $\mathbf{k}_+^{(t)}$ and $\mathbf{k}_-^{(t)}$ , i.e., $\mathbf{k}_{\pm}^{(t+1)} - \mathbf{k}_{\pm}^{(t)} = \beta_{\pm,\pm}^{(t)} \mathbf{q}_{\pm}^{(t)} + \sum_{n \in S_{\pm}} \sum_{i=2}^M \beta_{n,\pm,i}^{(t)} \mathbf{q}_{n,i}^{(t)}$
$\beta_{n,i,\pm}^{(t)}, \beta_{n,i,n',i'}^{(t)}$	linear combinations coefficients for the dynamics of $\mathbf{k}_{n,i}^{(t)}$ , i.e., $\mathbf{k}_{n,i}^{(t+1)} - \mathbf{k}_{n,i}^{(t)} = \beta_{n,i,+}^{(t)} \mathbf{q}_+^{(t)} + \beta_{n,i,-}^{(t)} \mathbf{q}_-^{(t)} + \sum_{n'=1}^N \sum_{i'=2}^M \beta_{n,i,n',i'}^{(t)} \mathbf{q}_{n',i'}^{(t)}$
$\text{softmax}(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle)$	a general references to $\frac{\exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}$ for $n \in S_+$ , $i \in [M] \setminus \{1\}$ , and $\frac{\exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}$ for $n \in S_-$ , $i \in [M] \setminus \{1\}$
$\text{softmax}(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)$	a general references to $\frac{\exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}$ for $n \in S_+$ , $i, j \in [M] \setminus \{1\}$ , and $\frac{\exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}$ for $n \in S_-$ , $i, j \in [M] \setminus \{1\}$
$\text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)$	a general references to $\frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,+}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}$ for $n \in S_+$ , $i, j \in [M] \setminus \{1\}$ , and $\frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,-}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}$ for $n \in S_-$ , $i, j \in [M] \setminus \{1\}$
$\Lambda_{n,\pm,j}^{(t)}, \Lambda_{n,i,\pm,j}^{(t)}$	$\Lambda_{n,\pm,j}^{(t)} := \langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle - \langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle$ , $\Lambda_{n,i,\pm,j}^{(t)} := \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle - \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle$

**Definition 6 (Gradient Decomposition).** *There exist coefficients  $\alpha_{+,+}^{(t)}$ ,  $\alpha_{n,+i}^{(t)}$ ,  $\alpha_{-,-}^{(t)}$ ,  $\alpha_{n,-i}^{(t)}$ ,  $\alpha_{n,i,+}^{(t)}$ ,  $\alpha_{n,i,-}^{(t)}$ ,  $\alpha_{n,i,n',i'}^{(t)}$ ,  $\beta_{+,+}^{(t)}$ ,  $\beta_{n,+i}^{(t)}$ ,  $\beta_{-,-}^{(t)}$ ,  $\beta_{n,-i}^{(t)}$ ,  $\beta_{n,i,+}^{(t)}$ ,  $\beta_{n,i,-}^{(t)}$ ,  $\beta_{n,i,n',i'}^{(t)}$  such that*

$$\Delta \mathbf{q}_+^{(t)} := \mathbf{q}_+^{(t+1)} - \mathbf{q}_+^{(t)} = \alpha_{+,+}^{(t)} \mathbf{k}_+^{(t)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+i}^{(t)} \mathbf{k}_{n,i}^{(t)}$$

$$\Delta \mathbf{q}_-^{(t)} := \mathbf{q}_-^{(t+1)} - \mathbf{q}_-^{(t)} = \alpha_{-,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n \in S_-} \sum_{i=2}^M \alpha_{n,-i}^{(t)} \mathbf{k}_{n,i}^{(t)}$$

$$\Delta \mathbf{q}_{n,i}^{(t)} := \mathbf{q}_{n,i}^{(t+1)} - \mathbf{q}_{n,i}^{(t)} = \alpha_{n,i,+}^{(t)} \mathbf{k}_+^{(t)} + \alpha_{n,i,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n'=1}^N \sum_{i'=2}^M \alpha_{n,i,n',i'}^{(t)} \mathbf{k}_{n',i'}^{(t)}$$

$$\Delta \mathbf{k}_+^{(t)} := \mathbf{k}_+^{(t+1)} - \mathbf{k}_+^{(t)} = \beta_{+,+}^{(t)} \mathbf{q}_+^{(t)} + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+i}^{(t)} \mathbf{q}_{n,i}^{(t)}$$

$$\Delta \mathbf{k}_-^{(t)} := \mathbf{k}_-^{(t+1)} - \mathbf{k}_-^{(t)} = \beta_{-,-}^{(t)} \mathbf{q}_-^{(t)} + \sum_{n \in S_-} \sum_{i=2}^M \beta_{n,-i}^{(t)} \mathbf{q}_{n,i}^{(t)}$$

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

$$\Delta \mathbf{k}_{n,i}^{(t)} := \mathbf{k}_{n,i}^{(t+1)} - \mathbf{k}_{n,i}^{(t)} = \beta_{n,i,+}^{(t)} \mathbf{q}_+^{(t)} + \beta_{n,i,-}^{(t)} \mathbf{q}_-^{(t)} + \sum_{n'=1}^N \sum_{i'=2}^M \beta_{n,i,n',i'}^{(t)} \mathbf{q}_{n',i'}^{(t)}$$

for  $i, i' \in [M] \setminus \{1\}$  and  $n, n' \in [N]$ .

**Lemma 8** (Update Rule for V). *The coefficients  $\gamma_{V,+}^{(t)}, \gamma_{V,-}^{(t)}, \rho_{V,n,i}^{(t)}$  defined in Definition 1 satisfy the following iterative equations:*

$$\begin{aligned} \gamma_{V,+}^{(t+1)} &= \gamma_{V,+}^{(t)} - \frac{\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle}{NM} \sum_{n \in S_+} \tilde{\ell}_n^{(t)} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\ &\quad + \sum_{n \in S_+} \tilde{\ell}_n^{(t)} \sum_{i=2}^M \frac{-\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle}{NM} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\ \gamma_{V,-}^{(t+1)} &= \gamma_{V,-}^{(t)} - \frac{\eta \langle \tilde{\boldsymbol{\mu}}_-, \tilde{\boldsymbol{\mu}}_-^{(t)} \rangle}{NM} \sum_{n \in S_-} \tilde{\ell}_n^{(t)} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_-, \tilde{\mathbf{k}}_-^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_-, \tilde{\mathbf{k}}_-^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_-, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_-^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_-^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\ &\quad + \sum_{n \in S_+} \tilde{\ell}_n^{(t)} \sum_{i=2}^M \frac{-\eta \langle \tilde{\boldsymbol{\mu}}_-, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle}{NM} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_-, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_-, \tilde{\mathbf{k}}_-^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_-, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_-^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right), \\ \rho_{V,n,i}^{(t+1)} &= \rho_{V,n,i}^{(t)} - \frac{\eta}{NM} \sum_{n' \in S_+} \tilde{\ell}_{n'}^{(t)} \left( \left( \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \frac{\exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \right. \\ &\quad \left. \left. + \sum_{j=2}^M \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \right. \\ &\quad \left. \left. + \sum_{i'=2}^M \left( \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\xi}}_{n',i'}^{(t)} \rangle \frac{\exp(\langle \tilde{\mathbf{q}}_{n',i'}^{(t)}, \tilde{\mathbf{k}}_{n',i'}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n',i'}^{(t)}, \tilde{\mathbf{k}}_{n',i'}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n',i'}^{(t)}, \tilde{\mathbf{k}}_{n',k}^{(t)} \rangle)} \right. \right. \\ &\quad \left. \left. + \sum_{j=2}^M \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\xi}}_{n',i'}^{(t)} \rangle \frac{\exp(\langle \tilde{\mathbf{q}}_{n',j}^{(t)}, \tilde{\mathbf{k}}_{n',i'}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n',j}^{(t)}, \tilde{\mathbf{k}}_{n',i'}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n',j}^{(t)}, \tilde{\mathbf{k}}_{n',k}^{(t)} \rangle)} \right) \right) \\ &\quad - \frac{\eta}{NM} \sum_{n' \in S_-} (\cdot) \end{aligned}$$

for  $i \in [M] \setminus \{1\}, n \in [N]$ .

*Proof.* The gradient of  $\mathbf{W}_V$  can be obtained using the chain rule as follows

$$\nabla_{\mathbf{w}_V} L_S(\theta) = \frac{1}{N} \sum_{n=1}^N y_n \ell'(y_n f(\mathbf{X}_n, \theta)) \nabla_{\mathbf{w}_V} f(\mathbf{X}_n, \theta)$$

$$= \frac{1}{NM} \sum_{n=1}^N y_n \ell'_n(\theta) \left[ \mathbf{w}_O \sum_{l=1}^M \varphi(\mathbf{x}_{n,l} \mathbf{W}_Q \mathbf{W}_K^\top (\mathbf{X}_n)^\top) \mathbf{X}_n \right]^\top$$

Base on above, we have

$$\begin{aligned} \mathbf{x}^\top \nabla_{\mathbf{w}_V} \widetilde{L}_S(\theta) \mathbf{w}_O &= \frac{1}{NM} \sum_{n \in S_+} \tilde{\ell}'_n(\theta) \left( \langle \mathbf{x}, \tilde{\boldsymbol{\mu}}_+ \rangle \frac{\exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+)}{\exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+) + \sum_{k=2}^M \exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,k})} \right. \\ &+ \sum_{j=2}^M \langle \mathbf{x}, \tilde{\boldsymbol{\mu}}_+ \rangle \frac{\exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+)}{\exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+) + \sum_{k=2}^M \exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,k})} \\ &+ \sum_{i=2}^M \langle \mathbf{x}, \tilde{\boldsymbol{\xi}}_{n,i} \rangle \frac{\exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,i})}{\exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+) + \sum_{k=2}^M \exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,k})} \\ &+ \left. \sum_{j=2}^M \langle \mathbf{x}, \tilde{\boldsymbol{\xi}}_{n,i} \rangle \frac{\exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,i})}{\exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+) + \sum_{k=2}^M \exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,k})} \right) \|\mathbf{w}_O\|_2^2 \\ &+ \frac{1}{NM} \sum_{n \in S_-} (\cdot) \end{aligned}$$

where the second equality we expand  $X_n$  into vectors and make inner products with  $x$ , the third equality we materializing all the  $x_{n,i}$  (e.g.,  $x_{n,1} = \boldsymbol{\mu}_+$  for  $n \in S_+$ ). Note the orthogonality between  $\boldsymbol{\mu}$  and  $\boldsymbol{\xi}_{n,i}$ , we can remove many of the terms in this equation. For any  $x = \tilde{\boldsymbol{\mu}}'_+ \in B(\boldsymbol{\mu}_+, \tau)$ , we have

$$\begin{aligned} \tilde{\boldsymbol{\mu}}_+^{\prime\top} \nabla_{\mathbf{w}_V} \widetilde{L}_S(\theta) \mathbf{w}_O &= \frac{1}{NM} \sum_{n \in S_+} \tilde{\ell}'_n(\theta) \left( \langle \tilde{\boldsymbol{\mu}}'_+, \tilde{\boldsymbol{\mu}}_+ \rangle \frac{\exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+)}{\exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+) + \sum_{k=2}^M \exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,k})} \right. \\ &+ \sum_{j=2}^M \langle \tilde{\boldsymbol{\mu}}'_+, \tilde{\boldsymbol{\mu}}_+ \rangle \frac{\exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+)}{\exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+) + \sum_{k=2}^M \exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,k})} \\ &+ \sum_{i=2}^M \langle \tilde{\boldsymbol{\mu}}'_+, \tilde{\boldsymbol{\xi}}_{n,i} \rangle \frac{\exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,i})}{\exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+) + \sum_{k=2}^M \exp(\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,k})} \\ &+ \left. \sum_{j=2}^M \langle \tilde{\boldsymbol{\mu}}'_+, \tilde{\boldsymbol{\xi}}_{n,i} \rangle \frac{\exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,i})}{\exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\mu}}_+) + \sum_{k=2}^M \exp(\tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_Q \mathbf{W}_K^\top \tilde{\boldsymbol{\xi}}_{n,k})} \right) \|\mathbf{w}_O\|_2^2 \end{aligned}$$

Then we have

$$\begin{aligned} \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_V^{(t+1)} \mathbf{w}_O - \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_V^{(t)} \mathbf{w}_O &= \tilde{\boldsymbol{\mu}}_+^\top (-\eta \nabla_{\mathbf{w}_V} \widetilde{L}_S(\theta(t))) \mathbf{w}_O \\ &= -\frac{\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+ \rangle}{NM} \sum_{n \in S_+} \tilde{\ell}'_n(t) \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &+ \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \left. \right) \|\mathbf{w}_O\|_2^2 \\ &+ \sum_{n \in S_+} \tilde{\ell}'_n(t) \sum_{i=2}^M \frac{-\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i} \rangle}{NM} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &+ \left. \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \|\mathbf{w}_O\|_2^2 \end{aligned}$$

Dividing by  $\|\mathbf{w}_O\|_2^2$  we get

$$\begin{aligned} \gamma_{V,+}^{(t+1)} = & \gamma_{V,+}^{(t)} - \frac{\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle}{NM} \sum_{n \in S_+} \tilde{\rho}_n^{(t)} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ & \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\ & + \sum_{n \in S_+} \tilde{\rho}_n^{(t)} \sum_{i=2}^M \frac{-\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle}{NM} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ & \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \end{aligned}$$

This proves the update rule for  $\gamma_{V,+}^{(t)}$ . The proof for  $\gamma_{V,-}^{(t)}$  and  $\rho_{V,n,i}^{(t)}$  is similar to it.  $\square$

**Lemma 9** (Update Rule for QK, Lemma B.3 in Jiang et al. (2024)). *The dynamics of  $\mathbf{x}^\top \mathbf{W}_Q \mathbf{W}_K \mathbf{x}$  can be characterized as follows:*

$$\begin{aligned} & \langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \\ & = \alpha_+^{(t)} \|\mathbf{k}_+^{(t)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle \\ & \quad + \beta_+^{(t)} \|\mathbf{q}_+^{(t)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle \\ & \quad + \left( \alpha_+^{(t)} \mathbf{k}_+^{(t)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(t)} \mathbf{k}_{n,i}^{(t)} \right) \\ & \quad \cdot \left( \beta_+^{(t)} \mathbf{q}_+^{(t)\top} + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+}^{(t)} \mathbf{q}_{n,i}^{(t)\top} \right), \end{aligned} \tag{2}$$

$$\begin{aligned} & \langle \mathbf{q}_-^{(t+1)}, \mathbf{k}_-^{(t+1)} \rangle - \langle \mathbf{q}_-^{(t)}, \mathbf{k}_-^{(t)} \rangle \\ & = \alpha_{-,-}^{(t)} \|\mathbf{k}_-^{(t)}\|_2^2 + \sum_{n \in S_-} \sum_{i=2}^M \alpha_{n,-}^{(t)} \langle \mathbf{k}_-^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle \\ & \quad + \beta_{-,-}^{(t)} \|\mathbf{q}_-^{(t)}\|_2^2 + \sum_{n \in S_-} \sum_{i=2}^M \beta_{n,-}^{(t)} \langle \mathbf{q}_-^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle \\ & \quad + \left( \alpha_{-,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n \in S_-} \sum_{i=2}^M \alpha_{n,-}^{(t)} \mathbf{k}_{n,i}^{(t)} \right) \\ & \quad \cdot \left( \beta_{-,-}^{(t)} \mathbf{q}_-^{(t)\top} + \sum_{n \in S_-} \sum_{i=2}^M \beta_{n,-}^{(t)} \mathbf{q}_{n,i}^{(t)\top} \right), \end{aligned} \tag{3}$$

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

$$\begin{aligned}
& \langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle \\
&= \alpha_{n,i,+}^{(t)} \|\mathbf{k}_+^{(t)}\|_2^2 + \alpha_{n,i,-}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_-^{(t)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \alpha_{n,i,n',l}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_{n',l}^{(t)} \rangle \\
&+ \beta_{+,+}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle + \sum_{n' \in S_+} \sum_{l=2}^M \beta_{n',+,l}^{(t)} \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',l}^{(t)} \rangle \\
&+ \left( \alpha_{n,i,+}^{(t)} \mathbf{k}_+^{(t)} + \alpha_{n,i,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n'=1}^N \sum_{l=2}^M \alpha_{n,i,n',l}^{(t)} \mathbf{k}_{n',l}^{(t)} \right) \\
&\cdot \left( \beta_{+,+}^{(t)} \mathbf{q}_+^{(t)\top} + \sum_{n' \in S_+} \sum_{l=2}^M \beta_{n',+,l}^{(t)} \mathbf{q}_{n',l}^{(t)\top} \right), \tag{4}
\end{aligned}$$

$$\begin{aligned}
& \langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_-^{(t+1)} \rangle - \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_-^{(t)} \rangle \\
&= \alpha_{n,i,-}^{(t)} \|\mathbf{k}_-^{(t)}\|_2^2 + \alpha_{n,i,+}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_-^{(t)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \alpha_{n,i,n',l}^{(t)} \langle \mathbf{k}_-^{(t)}, \mathbf{k}_{n',l}^{(t)} \rangle \\
&+ \beta_{-,-}^{(t)} \langle \mathbf{q}_-^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle + \sum_{n' \in S_-} \sum_{l=2}^M \beta_{n',-,l}^{(t)} \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',l}^{(t)} \rangle \\
&+ \left( \alpha_{n,i,+}^{(t)} \mathbf{k}_+^{(t)} + \alpha_{n,i,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n'=1}^N \sum_{l=2}^M \alpha_{n,i,n',l}^{(t)} \mathbf{k}_{n',l}^{(t)} \right) \\
&\cdot \left( \beta_{-,-}^{(t)} \mathbf{q}_-^{(t)\top} + \sum_{n' \in S_-} \sum_{l=2}^M \beta_{n',-,l}^{(t)} \mathbf{q}_{n',l}^{(t)\top} \right), \tag{5}
\end{aligned}$$

$$\begin{aligned}
& \langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle \\
&= \alpha_{+,+}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle + \sum_{n' \in S_+} \sum_{l=2}^M \alpha_{n',+,l}^{(t)} \langle \mathbf{k}_{n,j}^{(t)}, \mathbf{k}_{n',l}^{(t)} \rangle \\
&+ \beta_{n,j,+}^{(t)} \|\mathbf{q}_+^{(t)}\|_2^2 + \beta_{n,j,-}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_-^{(t)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \beta_{n,j,n',l}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_{n',l}^{(t)} \rangle \\
&+ \left( \alpha_{+,+}^{(t)} \mathbf{k}_+^{(t)} + \sum_{n' \in S_+} \sum_{l=2}^M \alpha_{n',+,l}^{(t)} \mathbf{k}_{n',l}^{(t)} \right) \\
&\cdot \left( \beta_{n,j,+}^{(t)} \mathbf{q}_+^{(t)\top} + \beta_{n,j,-}^{(t)} \mathbf{q}_-^{(t)\top} + \sum_{n'=1}^N \sum_{l=2}^M \beta_{n,j,n',l}^{(t)} \mathbf{q}_{n',l}^{(t)\top} \right), \tag{6}
\end{aligned}$$

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

$$\begin{aligned}
& \langle \mathbf{q}_-^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle - \langle \mathbf{q}_-^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle \\
&= \alpha_{-,-}^{(t)} \langle \mathbf{k}_-^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle + \sum_{n' \in S_-} \sum_{l=2}^M \alpha_{n',-}^{(t)} \langle \mathbf{k}_{n',l}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle \\
&\quad + \beta_{n,j,-}^{(t)} \|\mathbf{q}_-^{(t)}\|_2^2 + \beta_{n,j,+}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_-^{(t)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \beta_{n,j,n',l}^{(t)} \langle \mathbf{q}_-^{(t)}, \mathbf{q}_{n',l}^{(t)} \rangle \\
&\quad + \left( \alpha_{-,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n' \in S_-} \sum_{l=2}^M \alpha_{n',-}^{(t)} \mathbf{k}_{n',l}^{(t)} \right) \\
&\quad \cdot \left( \beta_{n,j,+}^{(t)} \mathbf{q}_+^{(t)\top} + \beta_{n,j,-}^{(t)} \mathbf{q}_-^{(t)\top} + \sum_{n'=1}^N \sum_{l=2}^M \beta_{n,j,n',l}^{(t)} \mathbf{q}_{n',l}^{(t)\top} \right), \tag{7}
\end{aligned}$$

$$\begin{aligned}
& \langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle - \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle \\
&= \alpha_{n,i,+}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle + \alpha_{n,i,-}^{(t)} \langle \mathbf{k}_-^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \alpha_{n,i,n',l}^{(t)} \langle \mathbf{k}_{n',l}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle \\
&\quad + \beta_{n,j,+}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle + \beta_{n,j,-}^{(t)} \langle \mathbf{q}_-^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \beta_{n,j,n',l}^{(t)} \langle \mathbf{q}_{n',l}^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle \\
&\quad + \left( \alpha_{n,i,+}^{(t)} \mathbf{k}_+^{(t)} + \alpha_{n,i,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n'=1}^N \sum_{l=2}^M \alpha_{n,i,n',l}^{(t)} \mathbf{k}_{n',l}^{(t)} \right) \\
&\quad \cdot \left( \beta_{n,j,+}^{(t)} \mathbf{q}_+^{(t)\top} + \beta_{n,j,-}^{(t)} \mathbf{q}_-^{(t)\top} + \sum_{n'=1}^N \sum_{l=2}^M \beta_{n,j,n',l}^{(t)} \mathbf{q}_{n',l}^{(t)\top} \right), \tag{8}
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n \in [N]$ .

## E CONCENTRATION INEQUALITIES

In this section, we will give some concentration inequalities that show some important properties of the data and the ViT parameters at random initialization.

**Lemma 10** (Lemma B.1 in Cao et al. (2022)). *Suppose that  $\delta > 0$  and  $n \geq 8 \log(4/\delta)$ . Then with probability at least  $1 - \delta$ ,*

$$\frac{N}{4} \leq |\{n \in [N] : y_n = 1\}|, |\{n \in [N] : y_n = -1\}| \leq \frac{3N}{4}.$$

**Lemma 11** (Initialization of V, Lemma C.2 in Jiang et al. (2024)). *Suppose that  $\delta > 0$ . Then with probability at least  $1 - \delta$ ,*

$$|V_{\pm}^{(0)}| \leq d_h^{-\frac{1}{4}}, \quad |V_{n,i}^{(0)}| \leq d_h^{-\frac{1}{4}}$$

for  $i \in [M] \setminus \{1\}, n \in [N]$ .

**Lemma 12** (Initialization of QK, Lemma C.3 in Jiang et al. (2024)). *Suppose that  $\delta > 0$ . Then with probability at least  $1 - \delta$ ,*

$$\begin{aligned}
\frac{\|\boldsymbol{\mu}\|_2^2 \sigma_p^2 d_h}{2} &\leq \|\mathbf{q}_{\pm}^{(0)}\|_2^2 \leq \frac{3\|\boldsymbol{\mu}\|_2^2 \sigma_p^2 d_h}{2}, \\
\frac{\sigma_p^2 \sigma_h^2 d d_h}{2} &\leq \|\mathbf{q}_{n,i}^{(0)}\|_2^2 \leq \frac{3\sigma_p^2 \sigma_h^2 d d_h}{2}, \\
\frac{\|\boldsymbol{\mu}\|_2^2 \sigma_p^2 d_h}{2} &\leq \|\mathbf{k}_{\pm}^{(0)}\|_2^2 \leq \frac{3\|\boldsymbol{\mu}\|_2^2 \sigma_p^2 d_h}{2},
\end{aligned}$$

$$\begin{aligned}
1296 \quad & \frac{\sigma_p^2 \sigma_h^2 d d_h}{2} \leq \|\mathbf{k}_{n,i}^{(0)}\|_2^2 \leq \frac{3\sigma_p^2 \sigma_h^2 d d_h}{2}, \\
1297 \quad & |\langle \mathbf{q}_+^{(0)}, \mathbf{q}_-^{(0)} \rangle| \leq 2\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1298 \quad & |\langle \mathbf{q}_\pm^{(0)}, \mathbf{q}_{n,i}^{(0)} \rangle| \leq 2\|\boldsymbol{\mu}\|_2 \sigma_p \sigma_h^2 d^{\frac{3}{2}} \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1299 \quad & |\langle \mathbf{k}_\pm^{(0)}, \mathbf{k}_\pm^{(0)} \rangle| \leq 2\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1300 \quad & |\langle \mathbf{q}_\pm^{(0)}, \mathbf{k}_\pm^{(0)} \rangle| \leq 2\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1301 \quad & |\langle \mathbf{q}_\pm^{(0)}, \mathbf{k}_\mp^{(0)} \rangle| \leq 2\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1302 \quad & |\langle \mathbf{q}_{n,i}^{(0)}, \mathbf{k}_\pm^{(0)} \rangle| \leq 2\|\boldsymbol{\mu}\|_2 \sigma_p \sigma_h^2 d^{\frac{3}{2}} \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1303 \quad & |\langle \mathbf{q}_{n,i}^{(0)}, \mathbf{q}_{n',j}^{(0)} \rangle| \leq 2\sigma_p^2 \sigma_h^2 d \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1304 \quad & |\langle \mathbf{k}_{n,i}^{(0)}, \mathbf{k}_{n',j}^{(0)} \rangle| \leq 2\sigma_p^2 \sigma_h^2 d \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1305 \quad & |\langle \mathbf{k}_\pm^{(0)}, \mathbf{k}_{n,i}^{(0)} \rangle| \leq 2\|\boldsymbol{\mu}\|_2 \sigma_p \sigma_h^2 d^{\frac{3}{2}} \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1306 \quad & |\langle \mathbf{q}_\pm^{(0)}, \mathbf{k}_{n,i}^{(0)} \rangle| \leq 2\|\boldsymbol{\mu}\|_2 \sigma_p \sigma_h^2 d^{\frac{3}{2}} \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}, \\
1307 \quad & |\langle \mathbf{q}_{n,i}^{(0)}, \mathbf{k}_{n',j}^{(0)} \rangle| \leq 2\sigma_p^2 \sigma_h^2 d \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}$  and  $n, n' \in [N]$ .

**Lemma 13** (Lemma B.2 in Cao et al. (2022) and Lemma B.4 in Kou et al. (2023)). *Suppose that  $\delta > 0$  and  $d = \Omega(\log(4NM/\delta))$ . Then with probability at least  $1 - \delta$*

$$\begin{aligned}
1319 \quad & \frac{\sigma_p^2 d}{2} \leq \|\boldsymbol{\xi}_{n,i}\|_2^2 \leq \frac{3\sigma_p^2 d}{2}, \\
1320 \quad & |\langle \boldsymbol{\xi}_{n,i}, \boldsymbol{\xi}_{n',i'} \rangle| \leq 2\sigma_p^2 \cdot \sqrt{d \log(4N^2 M^2 / \delta)}, \\
1321 \quad & \frac{\sigma_p^2 d}{2} - 2\sigma_p \tau \sqrt{2 \log(4NM/\delta)} - \tau^2 \leq \|\tilde{\boldsymbol{\xi}}_{n,i}\|_2^2 \leq \frac{3\sigma_p^2 d}{2} + 2\sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2, \\
1322 \quad & (\|\boldsymbol{\mu}\|_2 - \tau)^2 \leq \langle \tilde{\boldsymbol{\mu}}_\pm, \tilde{\boldsymbol{\mu}}'_\pm \rangle \leq (\|\boldsymbol{\mu}\|_2 + \tau)^2, \\
1323 \quad & |\langle \tilde{\boldsymbol{\mu}}_\pm, \tilde{\boldsymbol{\xi}}_{n,i} \rangle| \leq \|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2 \\
1324 \quad & |\langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\xi}}_{n',i'} \rangle| \leq 2\sigma_p^2 \cdot \sqrt{d \log(4N^2 M^2 / \delta)} + 2\sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2
\end{aligned}$$

for  $i, i' \in [M] \setminus \{1\}$ ,  $n, n' \in [N]$ ,  $i \neq i'$  or  $n \neq n'$ .

## F BENIGN OVERFITTING IN CASE 1

In this section, we consider the benign overfitting regime under the condition that  $N \cdot \text{SNR}^2 = \Omega(1)$  and  $\tau \leq O(\frac{\|\boldsymbol{\mu}\|_2}{\log d_h})$ . We analyze the dynamics of  $V_\pm$ ,  $V_{n,i}$ , the inner product  $\mathbf{q}_\pm$ ,  $\mathbf{q}_{n,i}$ , and  $\mathbf{k}_\pm$ ,  $\mathbf{k}_{n,i}$  during adversarial training, and further give the upper bound for clean test error and robust test error. The proofs in this section are based on the results in Section E, which hold with high probability.

### F.1 STAGE I

In Stage I,  $V_\pm^{(t)}$ ,  $V_{n,i}^{(t)}$  begin to pull apart until  $|V_\pm^{(t)}|$  is sufficiently larger than  $|V_{n,i}^{(t)}|$ . At the same time, the inner products of  $q$  and  $k$  maintain their magnitude.

**Lemma 14** (Gradient of Loss). *As long as  $\max\{|V_+^{(t)}|, |V_-^{(t)}|, |V_{n,i}^{(t)}|\} = o(1)$ , we have  $-\ell'(y_n f(\tilde{\mathbf{X}}_n, \theta(t)))$  remains  $1/2 \pm o(1)$ .*

*Proof.* Note that  $\ell(z) = \log(1 + \exp(-z))$  and  $-\ell' = \exp(-z)/(1 + \exp(-z))$ , without loss of generality, we assume  $y_n = 1$ , we have

$$1348 \quad -\ell'(f(\tilde{\mathbf{X}}_n, \theta(t))) = \frac{1}{1349 \quad 1 + \exp\left(\frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^\top \mathbf{W}_Q^{(t)} \mathbf{W}_K^{(t)\top} \tilde{\mathbf{X}}_n) \tilde{\mathbf{X}}_n^\top \mathbf{W}_V^{(t)} \mathbf{w}_O\right)}.$$

Note that

$$-\max\{|V_+^{(t)}|, |V_-^{(t)}|, |V_{n,i}^{(t)}|\} \leq \frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^\top \mathbf{W}_Q^{(t)\top} \tilde{\mathbf{X}}_n^\top) \tilde{\mathbf{X}}_n \mathbf{W}_V^{(t)} \mathbf{w}_O \leq \max\{|V_+^{(t)}|, |V_-^{(t)}|, |V_{n,i}^{(t)}|\}.$$

Then we have

$$\begin{aligned} -\ell'(f(\tilde{\mathbf{X}}_n, \theta(t))) &\geq \frac{1}{1 + \exp(0 + o(1))} \geq \frac{1}{2 + o(1)} \geq \frac{1}{2} - o(1), \\ -\ell'(f(\tilde{\mathbf{X}}_n, \theta(t))) &\leq \frac{\exp(0 + o(1))}{1 + \exp(0 + o(1))} \leq \frac{1 + o(1)}{1 + 1 + o(1)} \leq \frac{1}{2} + o(1). \end{aligned}$$

□

**Lemma 15** (Bound of Attention). *As long as  $|\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle| = o(1)$ , we have*

$$\begin{aligned} \frac{1}{M} - o(1) &\leq \text{softmax}(\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle) \leq \frac{1}{M} + o(1), \\ \frac{1}{M} - o(1) &\leq \text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle) \leq \frac{1}{M} + o(1), \\ \frac{1}{M} - o(1) &\leq \text{softmax}(\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle) \leq \frac{1}{M} + o(1), \\ \frac{1}{M} - o(1) &\leq \text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle) \leq \frac{1}{M} + o(1). \end{aligned}$$

*Proof.* It is clear that  $\exp(o(1)) = 1 + o(1)$ . Therefore, as long as  $|\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle| = o(1)$ , we have

$$\begin{aligned} \frac{1}{M} - o(1) &= \frac{1}{1 + (M-1) + (M-1)o(1)} = \frac{1}{1 + (M-1)\exp(o(1))} = \\ &= \frac{\exp(-o(1))}{\exp(-o(1)) + (M-1)\exp(o(1))} \leq \text{softmax}(\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle) \leq \frac{\exp(o(1))}{\exp(o(1)) + (M-1)\exp(-o(1))} \\ &= \frac{\exp(o(1))}{\exp(o(1)) + (M-1)} = \frac{1 + o(1)}{1 + o(1) + (M-1)} = \frac{1}{M} + o(1) \end{aligned}$$

Similarly, we have

$$\begin{aligned} \frac{1}{M} - o(1) &\leq \text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle) \leq \frac{1}{M} + o(1), \\ \frac{1}{M} - o(1) &\leq \text{softmax}(\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle) \leq \frac{1}{M} + o(1), \\ \frac{1}{M} - o(1) &\leq \text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle) \leq \frac{1}{M} + o(1). \end{aligned}$$

□

**Lemma 16** (Upper bound of V). *Let  $T_0 = \mathcal{O}\left(\frac{1}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right)$ . Then under the same conditions as Theorem 2 we have*

$$|V_+^{(t)}|, |V_-^{(t)}|, |V_{n,i}^{(t)}| = \mathcal{O}(d_h^{-\frac{1}{4}})$$

for  $t \in [0, T_0]$ .

*Proof.* By Lemma 8, we have

$$\begin{aligned} |\gamma_{V,+}^{(t+1)} - \gamma_{V,+}^{(t)}| &\leq -\frac{\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle}{NM} \sum_{n \in S_+} \tilde{\ell}'_n(t) \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \|\mathbf{w}_O\|_2^2 \end{aligned}$$

$$\begin{aligned}
& + \frac{-\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle}{NM} \sum_{n \in S_+} \sum_{i=2}^M \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\
& + \left. \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\
& \leq \frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{NM} \cdot \frac{3N}{4} (1 + (M-1)) + \frac{\eta(\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2)}{NM} \cdot \frac{3N}{4} M(M-1) \\
& \leq O(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2),
\end{aligned}$$

where the second inequality is by Lemma 10 and Lemma 13 and  $-\tilde{\ell}_n^{(t)} \leq 1$ .

Similarly, we have

$$|\gamma_{V,-}^{(t+1)} - \gamma_{V,-}^{(t)}| \leq O(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2).$$

By Definition 4, we have

$$\begin{aligned}
|V_+^{(t)}| &= |V_+^{(0)} + \sum_{s=0}^{t-1} (\gamma_{V,+}^{(s+1)} - \gamma_{V,+}^{(s)}) \|\mathbf{w}_O\|_2^2| \\
&\leq |V_+^{(0)}| + \sum_{s=0}^{t-1} |\gamma_{V,+}^{(s+1)} - \gamma_{V,+}^{(s)}| \cdot \|\mathbf{w}_O\|_2^2 \\
&\leq d_h^{-\frac{1}{4}} + O(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2) \cdot \|\mathbf{w}_O\|_2^2 \cdot O\left(\frac{1}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right) \\
&= O(d_h^{-\frac{1}{4}}),
\end{aligned}$$

where the first inequality is by triangle inequality, the second inequality is by Lemma 11. Similarly, we have  $|V_-^{(t)}| = O(d_h^{-\frac{1}{4}})$ . By Lemma 8, we have

$$\begin{aligned}
|\rho_{V,n,i}^{(t+1)} - \rho_{V,n,i}^{(t)}| &\leq \left| -\frac{\eta}{NM} \sum_{n' \in S_+} \tilde{\ell}_{n'}^{(t)} \left( \left( \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \right. \\
&\quad + \left. \sum_{j=2}^M \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\
&\quad + \sum_{i=2}^M \left( \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\xi}}_{n',i'}^{(t)} \rangle \frac{\exp(\langle \tilde{\mathbf{q}}_{n',i'}^{(t)}, \tilde{\mathbf{k}}_{n',i'}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n',i'}^{(t)}, \tilde{\mathbf{k}}_{n',i'}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n',i'}^{(t)}, \tilde{\mathbf{k}}_{n',k}^{(t)} \rangle)} \right. \\
&\quad \left. \left. + \sum_{j=2}^M \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\xi}}_{n',i'}^{(t)} \rangle \frac{\exp(\langle \tilde{\mathbf{q}}_{n',j}^{(t)}, \tilde{\mathbf{k}}_{n',i'}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n',j}^{(t)}, \tilde{\mathbf{k}}_{n',i'}^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n',j}^{(t)}, \tilde{\mathbf{k}}_{n',k}^{(t)} \rangle)} \right) \right) \\
&\quad - \frac{\eta}{NM} \sum_{n' \in S_-} (\cdot) \Big| \\
&\leq \frac{3\eta\sigma_p^2 d}{2NM} \cdot M + \frac{\eta}{NM} \cdot MN \cdot 2\sigma_p^2 \cdot \sqrt{d \log(4N^2 M^2 / \delta)} \\
&\quad + \frac{\eta(\|\boldsymbol{\mu}\| \tau + (2M-1)\sigma_p \tau \sqrt{3d/2} + M\tau^2)}{NM} \cdot NM \\
&\leq \frac{2\eta\sigma_p^2 d}{N} + \frac{\eta(\|\boldsymbol{\mu}\| \tau + \sigma_p \tau \sqrt{2 \log(4MN/\delta)} + \tau^2)}{NM} \cdot NM^2 \\
&= O(\eta(\max\{\|\boldsymbol{\mu}\|_2 + \tau, \frac{\sigma_p^2 d}{N}\})) \\
&= O(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2)
\end{aligned}$$

where the second inequality is by Lemma 13 and  $-\ell_n^{(t)} \leq 1$ , the third inequality is by  $d = \tilde{\Omega}(\epsilon^{-2}N^2d_h) \geq 4N\sqrt{\log(4N^2M^2/\delta)}$ , the last inequality is by  $N \cdot \text{SNR}^2 = \Omega(1)$ . Then by Definition 4, we have

$$\begin{aligned} |V_{n,i}^{(t)}| &= |V_{n,i}^{(0)} + \sum_{s=0}^{t-1} (\rho_{V_{n,i}}^{(s+1)} - \rho_{V_{n,i}}^{(s)}) \|\mathbf{w}_O\|_2^2| \\ &\leq |V_{n,i}^{(0)}| + \sum_{s=0}^{t-1} |\rho_{V_{n,i}}^{(s+1)} - \rho_{V_{n,i}}^{(s)}| \cdot \|\mathbf{w}_O\|_2^2 \\ &\leq d_h^{-\frac{1}{4}} + O(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2) \cdot \|\mathbf{w}_O\|_2^2 \cdot O\left(\frac{1}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right) \\ &= O(d_h^{-\frac{1}{4}}), \end{aligned}$$

where the first inequality is by triangle inequality, the second inequality is by Lemma C.2, which completes the proof.  $\square$

**Lemma 17** (Inner Products Hold Magnitude). *Let  $T_0 = O\left(\frac{1}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right)$ . Then under the same conditions as Theorem 2, we have*

$$\begin{aligned} &|\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle|, |\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| \\ &= O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h} \log(6N^2M^2/\delta)\right), \\ &|\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{q}_{\mp}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{\mp}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',j}^{(t)} \rangle| \\ &= O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h} \log(6N^2M^2/\delta)\right), \\ &|\langle \mathbf{k}_{\pm}^{(t)}, \mathbf{k}_{\mp}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| \\ &= O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h} \log(6N^2M^2/\delta)\right), \\ &\|\mathbf{q}_{\pm}^{(t)}\|_2^2, \|\mathbf{k}_{\pm}^{(t)}\|_2^2 = \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h), \\ &\|\mathbf{q}_{n,i}^{(t)}\|_2^2, \|\mathbf{k}_{n,i}^{(t)}\|_2^2 = \Theta(\sigma_p^2 \sigma_h^2 d d_h) \end{aligned}$$

for  $i, j \in [M] \setminus \{1\}$ ,  $n, n' \in [N]$  and  $t \in [0, T_0]$ .

The proof for Lemma 17 is in Section I.4. Note that  $\sigma_h^2 \leq \min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}} \cdot (\log(6N^2M^2/\delta))^{-\frac{3}{2}}$ , thus  $O(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h} \log(6N^2M^2/\delta)) = o(1)$ .

**Lemma 18** (V's Beginning of Learning Signals). *Under the same conditions as Theorem 2, there exists*

$$T_1 = \frac{10M(3M+1)N}{\eta d_h^{\frac{1}{4}} \left( N(\|\boldsymbol{\mu}\|_2 - \tau)^2 - (N+30M^2)(\|\boldsymbol{\mu}\|_2\tau + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) - 60M^2\sigma_p^2d \right) \|\mathbf{w}_O\|_2}$$

such that the first element of the vector  $\tilde{X}_n \mathbf{W}_V^{(t)} \mathbf{w}_O$  dominates its other elements, that is,

$$\begin{aligned} V_+^{(t)} &\geq 3M \cdot |V_{n,i}^{(t)}|, \quad \text{for all } n \in S_+, i \in [M] \setminus \{1\}, \\ V_-^{(t)} &\leq -3M \cdot |V_{n,i}^{(t)}|, \quad \text{for all } n \in S_-, i \in [M] \setminus \{1\}. \end{aligned}$$

*Proof.* Let  $C$  be a constant larger than  $10M(3M+1)$ . As long as

$$N(\|\boldsymbol{\mu}\|_2 - \tau)^2 - (N+30M^2)(\|\boldsymbol{\mu}\|_2\tau + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) - 60M^2\sigma_p^2d \geq \frac{10M(3M+1)}{C} N(\|\boldsymbol{\mu}\|_2 + \tau)^2.$$

Thus, we further get

$$T_1 = \frac{10M(3M+1)N}{\eta d_h^{\frac{1}{4}} \left( N(\|\boldsymbol{\mu}\|_2 - \tau)^2 - (N+30M^2)(\|\boldsymbol{\mu}\|\tau + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) - 60M^2\sigma_p^2d \right) \|\mathbf{w}_O\|_2}$$

$$\leq \frac{C}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2} = \mathcal{O}\left(\frac{1}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2}\right),$$

which satisfies the time condition in Lemma 16 and Lemma 17. Then by Lemma 14 and Lemma 15 we have:

$$-\ell'_n(t) = \frac{1}{2} \pm o(1),$$

$$\frac{1}{M} - o(1) \leq \text{softmax}(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) \leq \frac{1}{M} + o(1),$$

$$\frac{1}{M} - o(1) \leq \text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) \leq \frac{1}{M} + o(1),$$

$$\frac{1}{M} - o(1) \leq \text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle) \leq \frac{1}{M} + o(1).$$

For  $i, j \in [M] \setminus \{1\}$ ,  $n \in [N]$ , and  $t \in [0, T_1]$ , plugging these into the update rule for  $\gamma_{V_+}^{(t)}$  shown in Lemma 8, we have:

$$\begin{aligned} \gamma_{V_+}^{(t+1)} - \gamma_{V_+}^{(t)} &= -\frac{\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle}{NM} \sum_{n \in S_+} \tilde{\ell}'_n(t) \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\ &\quad + \sum_{n \in S_+} \tilde{\ell}'_n(t) \sum_{i=2}^M \frac{-\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle}{NM} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\ &\geq \frac{\eta(\|\boldsymbol{\mu}\|_2 - \tau)_2^2}{NM} \cdot \frac{N}{4} \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot M \left(\frac{1}{M} \pm o(1)\right) \\ &\quad - \frac{\eta(\|\boldsymbol{\mu}\|\tau + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2)}{NM} \cdot \frac{N}{4} \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot M(M-1) \left(\frac{1}{M} \pm o(1)\right) \\ &\geq \frac{\eta \left( (\|\boldsymbol{\mu}\|_2 - \tau)_2^2 - (\|\boldsymbol{\mu}\|\tau + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) \right)}{10M} \end{aligned}$$

Then by Definition 4 and summing over  $T_1$  steps, we have:

$$\begin{aligned} V_+^{(T_1)} &\geq -|V_+^{(0)}| + T_1 \cdot \frac{\eta \left( (\|\boldsymbol{\mu}\|_2 - \tau)_2^2 - (\|\boldsymbol{\mu}\|\tau + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) \right)}{10M} \cdot \|\mathbf{w}_O\|_2^2 \\ &= -d_h^{-\frac{1}{4}} + T_1 \cdot \frac{\eta \left( (\|\boldsymbol{\mu}\|_2 - \tau)_2^2 - (\|\boldsymbol{\mu}\|\tau + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) \right)}{10M} \cdot \|\mathbf{w}_O\|_2^2 \end{aligned} \tag{9}$$

1566 Similarly, we have:

$$1567 V_-^{(T_1)} \leq d_h^{-\frac{1}{4}} - T_1 \cdot \frac{\eta \left( (\|\boldsymbol{\mu}\|_2 - \tau)_2^2 - (\|\boldsymbol{\mu}\|_2 + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) \right)}{10M} \cdot \|\mathbf{w}_O\|_2^2 \quad (10)$$

1570 Similarly, by the bound

$$1571 |\rho_{V_{n,i}^{(t+1)}} - \rho_{V_{n,i}^{(t)}}| \leq \frac{2\eta\sigma_p^2 d}{N} + \eta(\|\boldsymbol{\mu}\|_2 + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2)$$

1572 given in equation (9), we have

$$1573 |V_{n,i}^{(T_1)}| \leq |V_{n,i}^{(0)}| + T_1 \cdot \left( \frac{2\eta\sigma_p^2 d}{N} + \eta(\|\boldsymbol{\mu}\|_2 + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) \right) \cdot \|\mathbf{w}_O\|_2^2 \quad (11)$$

$$1574 = d_h^{-\frac{1}{4}} + T_1 \cdot \left( \frac{2\eta\sigma_p^2 d}{N} + \eta(\|\boldsymbol{\mu}\|_2 + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2) \right) \cdot \|\mathbf{w}_O\|_2^2$$

1575 According to equations (9), (10), and (11), it is easy to verify that

$$1576 V_+^{(T_1)} - 3M \cdot |V_{n,i}^{(T_1)}| \geq 0 \quad \text{and} \quad V_-^{(T_1)} + 3M \cdot |V_{n,i}^{(T_1)}| \leq 0,$$

1577 which completes the proof.

□

## 1588 F.2 STAGE II

1589 In stage II,  $\langle \mathbf{q}_+, \mathbf{k}_+ \rangle$ ,  $\langle \mathbf{q}_{n,i}, \mathbf{k}_+ \rangle$  grows while  $\langle \mathbf{q}_+, \mathbf{k}_{n,j} \rangle$ ,  $\langle \mathbf{q}_{n,i}, \mathbf{k}_{n,j} \rangle$  decreases, resulting in attention focusing more and more on the signals and less on the noises. By the results of stage I, we have the following conditions at the beginning of stage II

$$1590 V_+^{(T_1)} \geq 3M \cdot |V_{n,i}^{(T_1)}|,$$

$$1591 V_-^{(T_1)} \leq -3M \cdot |V_{n,i}^{(T_1)}|,$$

$$1592 |V_+^{(T_1)}|, |V_-^{(T_1)}|, |V_{n,i}^{(T_1)}| = O(d_h^{-\frac{1}{4}}),$$

$$1593 |\langle \mathbf{q}_\pm^{(T_1)}, \mathbf{k}_\pm^{(T_1)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_\pm^{(T_1)} \rangle|, |\langle \mathbf{q}_\pm^{(T_1)}, \mathbf{k}_{n,j}^{(T_1)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{n',j}^{(T_1)} \rangle|$$

$$1594 = O \left( \max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)} \right),$$

$$1595 |\langle \mathbf{q}_\pm^{(T_1)}, \mathbf{q}_\mp^{(T_1)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{q}_\pm^{(T_1)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{q}_{n',j}^{(T_1)} \rangle|$$

$$1596 = O \left( \max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)} \right),$$

$$1597 |\langle \mathbf{k}_\pm^{(T_1)}, \mathbf{k}_\pm^{(T_1)} \rangle|, |\langle \mathbf{k}_{n,i}^{(T_1)}, \mathbf{k}_\pm^{(T_1)} \rangle|, |\langle \mathbf{k}_{n,i}^{(T_1)}, \mathbf{k}_{n',j}^{(T_1)} \rangle|$$

$$1598 = O \left( \max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)} \right),$$

$$1599 \|\mathbf{q}_\pm^{(T_1)}\|_2, \|\mathbf{k}_\pm^{(T_1)}\|_2 = \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h),$$

$$1600 \|\mathbf{q}_{n,i}^{(T_1)}\|_2, \|\mathbf{k}_{n,i}^{(T_1)}\|_2 = \Theta(\sigma_p^2 \sigma_h^2 d d_h)$$

1601 for  $i, j \in [M] \setminus \{1\}$ ,  $n, n' \in [N]$ .

1602 *Notations.* To better characterize the gap between different inner products, we define the following notations:

- denote  $\Lambda_{n,+j}^{(t)} = \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle$ ,  $n \in S_+$ .
- denote  $\Lambda_{n,-j}^{(t)} = \langle \mathbf{q}_-^{(t)}, \mathbf{k}_-^{(t)} \rangle - \langle \mathbf{q}_-^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle$ ,  $n \in S_-$ .
- denote  $\Lambda_{n,i,+j}^{(t)} = \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle - \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle$ ,  $n \in S_+$ .
- denote  $\Lambda_{n,i,-j}^{(t)} = \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_-^{(t)} \rangle - \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle$ ,  $n \in S_-$ .

**Lemma 19** (Upper bound of V). *Let  $T_0 = O\left(\frac{1}{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 \log(6N^2M^2/\delta)}\right)$ . Then under the same conditions as Theorem 2, we have*

$$|V_+^{(t)}|, |V_-^{(t)}|, |V_{n,i}^{(t)}| = o(1)$$

for  $t \in [0, T_0]$ .

The proof of Lemma 19 is similar to that of Lemma 16, except that the time  $T_0$  is changed. Let  $T_2 = \Omega\left(\frac{1}{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 \log(6N^2M^2/\delta)}\right)$ , then by Lemma 19 and Lemma 14 we have  $\frac{1}{2} - o(1) \leq -\tilde{\ell}_n^{(t)} \leq \frac{1}{2} + o(1)$  for  $n \in [N], t \in [T_1, T_2]$ , which can simplify the calculations of  $\alpha$  and  $\beta$  defined in Definition 6 by their bounds. Next we prove the following four propositions  $\mathcal{B}(t), \mathcal{C}(t), \mathcal{D}(t), \mathcal{E}(t)$  by induction on  $t$  for  $t \in [T_1, T_2]$ :

- $\mathcal{B}(t)$ :

$$V_+^{(t)} \geq \eta C_3 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_1)$$

$$V_-^{(t)} \leq \eta C_3 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_1)$$

$$V_+^{(t)} \geq 3M \cdot |V_{n,i}^{(t)}|,$$

$$V_-^{(t)} \leq -3M \cdot |V_{n,i}^{(t)}|,$$

$$|V_+^{(t)}| \leq O(d_h^{-\frac{1}{4}}) + \eta C_4 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_1)$$

$$|V_-^{(t)}| \leq O(d_h^{-\frac{1}{4}}) + \eta C_4 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_1)$$

for  $i \in [M] \setminus \{1\}, n \in [N]$ .

- $\mathcal{C}(t)$ :

$$\|\mathbf{q}_\pm^{(t)}\|_2, \|\mathbf{k}_\pm^{(t)}\|_2 = \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h),$$

$$\|\mathbf{q}_{n,i}^{(t)}\|_2, \|\mathbf{k}_{n,i}^{(t)}\|_2 = \Theta(\sigma_p^2 \sigma_h^2 d d d_h),$$

$$|\langle \mathbf{q}_+^{(t)}, \mathbf{q}_-^{(t)} \rangle|, |\langle \mathbf{q}_+^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',j}^{(t)} \rangle| = o(1),$$

$$|\langle \mathbf{k}_+^{(t)}, \mathbf{k}_-^{(t)} \rangle|, |\langle \mathbf{k}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| = o(1),$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], i \neq j$  or  $n \neq n'$ .

- $\mathcal{D}(t)$ :

$$\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle \geq \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle$$

$$\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle \geq \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle$$

$$\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle \leq \langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle$$

$$\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle \leq \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle$$

$$\Lambda_{n,\pm,j}^{(t+1)} \geq \log \left( \exp(\Lambda_{n,\pm,j}^{(T_1)}) + \frac{\eta^2 C_8 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2M^2/\delta))^2} \cdot (t - T_1)(t - T_1 + 1) \right)$$

$$\Lambda_{n,i,\pm,j}^{(t+1)} \geq \log \left( \exp(\Lambda_{n,i,\pm,j}^{(T_1)}) + \frac{\eta^2 C_8 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h}{N (\log(6N^2M^2/\delta))^2} \cdot (t - T_1)(t - T_1 + 1) \right).$$

for  $i, j \in [M] \setminus \{1\}, n \in [N]$ .

•  $\mathcal{E}(t)$ :

$$|\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_{n,j}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| \leq \log(d_h^{1/4})$$

$$|\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle| = o(1)$$

for  $i, j \in [M] \setminus \{1\}$ ,  $n, \bar{n} \in [N]$ ,  $n \neq \bar{n}$ .

By the results of Stage I, we know that  $\mathcal{B}(T_1), \mathcal{C}(T_1), \mathcal{E}(T_1)$  are true. To prove that  $\mathcal{B}(t), \mathcal{C}(t), \mathcal{D}(t)$  and  $\mathcal{E}(t)$  are true in Stage II, we will prove the following claims holds for  $t \in [T_1, T_2]$ :

**Claim 1.**  $\mathcal{D}(T_1), \dots, \mathcal{D}(t-1), \mathcal{E}(T_1), \dots, \mathcal{E}(t) \implies \mathcal{B}(t+1)$

**Claim 2.**  $\mathcal{B}(T_1), \dots, \mathcal{B}(t), \mathcal{C}(T_1), \dots, \mathcal{C}(t), \mathcal{D}(T_1), \dots, \mathcal{D}(t-1), \mathcal{E}(T_1), \dots, \mathcal{E}(t) \implies \mathcal{D}(t)$

**Claim 3.**  $\mathcal{B}(T_1), \dots, \mathcal{B}(t), \mathcal{D}(T_1), \dots, \mathcal{D}(t-1), \mathcal{E}(T_1), \dots, \mathcal{E}(t) \implies \mathcal{C}(t+1)$

**Claim 4.**  $\mathcal{B}(T_1), \dots, \mathcal{B}(t), \mathcal{C}(T_1), \dots, \mathcal{C}(t), \mathcal{D}(T_1), \dots, \mathcal{D}(t-1), \mathcal{E}(T_1), \dots, \mathcal{E}(t) \implies \mathcal{E}(t+1)$

First, we emphasize that when the perturbation is sufficiently small, the softmax of the inner product  $\langle \mathbf{q}, \mathbf{k} \rangle$  remains robust with respect to such perturbations. This is formally established in Lemma 4. Importantly, this lemma enables us to uniformly control the behavior of the softmax function, rather than being restricted to the specific case of  $\tilde{\mathbf{q}}^{(t)}$  and  $\tilde{\mathbf{k}}^{(t)}$  at iteration  $t$ .

**Lemma 20.** *Suppose the perturbation satisfies  $\tau \leq O(\frac{\|\boldsymbol{\mu}\|_2}{\log d_h})$  and  $t \geq \Omega\left(\frac{1}{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_o\|_2^2 \log(6N^2M^2/\delta)}\right)$  ( $\mathcal{E}(t)$  holds). Then there exists a universal constant  $C \leq e/2$  such that*

$$\max_{\tilde{\mathbf{X}} \in B(X, \tau)} \text{softmax}\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right) / \min_{\tilde{\mathbf{X}} \in B(X, \tau)} \text{softmax}\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right) \leq C.$$

*Proof.* For any perturbed query-key pair, the softmax weight can be written as

$$\begin{aligned} \text{softmax}\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right) &= \frac{\exp\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right)}{\exp\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle\right) + \sum_{j'=2}^M \exp\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j'}^{(t)} \rangle\right)} \\ &= \frac{1}{\exp\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right) + \sum_{j'=2}^M \exp\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j'}^{(t)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right)}. \end{aligned} \quad (12)$$

Next, we analyze the difference in the logits. Expanding the perturbation terms yields

$$\begin{aligned} &\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle \\ &= \langle \mathbf{q}_+, \mathbf{k}_+ - \mathbf{k}_{n,j} \rangle + \langle \tau_+ \mathbf{W}_Q^{(t)}, \mathbf{k}_+ - \mathbf{k}_{n,j} \rangle \\ &\quad + \langle \mathbf{q}_+, \tau_+ \mathbf{W}_K^{(t)} - \tau_{n,j} \mathbf{W}_K^{(t)} \rangle + \langle \tau_+ \mathbf{W}_Q^{(t)}, \tau_+ \mathbf{W}_K^{(t)} - \tau_{n,j} \mathbf{W}_K^{(t)} \rangle \\ &= \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} - \mathbf{k}_{n,j}^{(t)} \rangle \pm o(1). \end{aligned} \quad (13)$$

The first equality follows from decomposing the perturbed terms, while the second uses the perturbation bound  $\tau \leq O(\frac{\|\boldsymbol{\mu}\|_2}{\log d_h})$ , together with the assumption that the logit magnitudes satisfy

$$|\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle| \leq \log(d_h^{1/2}).$$

An analogous relation holds for differences involving  $\tilde{\mathbf{k}}_{n,j'}^{(t)}$  as well.

Substituting equation 13 into equation 12, we obtain

$$\frac{1}{C} \text{softmax}\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right) \leq \text{softmax}\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right) \leq C \text{softmax}\left(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle\right),$$

for some absolute constant  $1 \leq C \leq e/2$ . This establishes the claim.

Thus, in the following discussion, we show that the updates related to  $\mathbf{q}^{(t)}$  and  $\mathbf{k}^{(t)}$  during adversarial training can be bounded, while the effect of the perturbation does not accumulate over time. Since  $C$  is a very small function, when computing the single-step update, we approximate the perturbed softmax at the previous time step by its clean state.  $\square$

## F.2.1 PROOF OF CLAIM 1

By the results of Stage I, we have

$$|\langle \mathbf{q}_{\pm}^{(T_1)}, \mathbf{k}_{\pm}^{(T_1)} \rangle|, |\langle \mathbf{q}_{\pm}^{(T_1)}, \mathbf{k}_{n,j}^{(T_1)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{\pm}^{(T_1)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{n,j}^{(T_1)} \rangle| = o(1)$$

Assume that  $\mathcal{D}(T_1), \dots, \mathcal{D}(t-1)$  ( $t \in [T_1, T_2]$ ) are true, then  $\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle$ ,  $\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle$  are monotonically non-decreasing and  $\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$ ,  $\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$  are monotonically non-increasing for  $s \in [T_1, t-1]$ , so we have

$$\begin{aligned} \langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle, \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle &\geq -o(1), \\ \langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle, \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle &\leq o(1), \end{aligned}$$

for  $s \in [T_1, t]$ . Further we have the lower bounds for the attention on signal  $\mu_+$  as follows for  $s \in [T_1, t]$ :

$$\begin{aligned} \text{softmax}(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) &\geq \frac{\exp(-o(1))}{\exp(-o(1)) + (M-1)\exp(o(1))} \\ &= \frac{1}{1 + (M-1)\exp(o(1))} \\ &= \frac{1}{1 + (M-1) + (M-1)o(1)} \\ &= \frac{1}{M} - o(1), \end{aligned} \tag{14}$$

where the second equality is by  $\exp(o(1)) = 1 + o(1)$ . Similarly, we have

$$\text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) \geq \frac{1}{M} - o(1). \tag{15}$$

Plugging them in the update rule for  $\gamma_{V_+}^{(s)}$  shown in Lemma 8 and we have

$$\begin{aligned} &\gamma_{V_+}^{(s+1)} - \gamma_{V_+}^{(s)} \\ &= -\frac{\eta \langle \tilde{\mu}_+, \tilde{\mu}_+^{(t)} \rangle}{NM} \sum_{n \in S_+} \tilde{\rho}_n^{(t)} \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_+, \mathbf{k}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \mathbf{q}_{n,j}^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,j}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{n,j}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)} \right) \\ &+ \sum_{n \in S_+} \tilde{\rho}_n^{(t)} \sum_{i=2}^M \frac{-\eta \langle \tilde{\mu}_+, \tilde{\xi}_{n,i}^{(t)} \rangle}{NM} \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_+, \mathbf{k}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \mathbf{q}_{n,j}^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,j}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{n,j}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)} \right) \\ &\geq \frac{\eta (\|\mu\|_2 - \tau)_2^2}{NMC} \cdot \frac{N}{4} \cdot \left( \frac{1}{2} \pm o(1) \right) \cdot M \left( \frac{1}{M} \pm o(1) \right) \\ &- \frac{\eta (\|\mu\| \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2)}{NMC} \cdot \frac{N}{4} \cdot \left( \frac{1}{2} \pm o(1) \right) \cdot M(M-1) \left( \frac{1}{M} \pm o(1) \right) \\ &\geq \frac{\eta \left( (\|\mu\|_2 - \tau)^2 - (\|\mu\| \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \right)}{10MC} = O \left( \frac{\eta (\|\mu\|_2 - \tau)^2}{10MC} \right), \end{aligned}$$

for  $s \in [T_1, t]$ . The first inequality is by Lemma 4. Then by Definition 4 and taking a summation, we have

$$\begin{aligned} V_+^{(t+1)} &\geq V_+^{(T_1)} + (t - T_1 + 1) \frac{\eta (\|\mu\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2}{10MC} \\ &\geq \eta C_3 (\|\mu\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_1 + 1), \end{aligned} \tag{16}$$

where the last inequality is by  $V_+^{(T_1)} \geq 0$  and  $M = \Theta(1)$ . Similarly, we have

$$V_-^{(t+1)} \leq -\eta C_3 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_1 + 1). \quad (17)$$

By  $|\rho_{V_{n,i}^{(t+1)}} - \rho_{V_{n,i}^{(t)}}| \leq \eta (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2 + \frac{2\sigma_p^2 d}{N})$  in (9) and taking a summation, we have

$$|V_{n,i}^{(t+1)}| \leq |V_{n,i}^{(t)}| + (t - T_1 + 1) \eta (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2 + \frac{2\sigma_p^2 d}{N}) \|\mathbf{w}_O\|_2^2. \quad (18)$$

Combining (16) and (18) we have

$$\begin{aligned} & V_+^{(t+1)} - 3M \cdot |V_{n,i}^{(t+1)}| \\ & \geq V_+^{(T_1)} + (t - T_1 + 1) \frac{\eta (\|\boldsymbol{\mu}\|_2 - \tau)^2}{10M} \|\mathbf{w}_O\|_2^2 \\ & \quad - 3M \cdot (|V_{n,i}^{(T_1)}| + (t - T_1 + 1) \eta (\max\{\|\boldsymbol{\mu}\|_2 \tau + \tau^2, \frac{2C_p^2 \sigma_p^2 d}{N}\}) \|\mathbf{w}_O\|_2^2) \\ & \geq V_+^{(T_1)} - 3M \cdot |V_{n,i}^{(T_1)}| + (t - T_1 + 1) \frac{\eta (\|\boldsymbol{\mu}\|_2 - \tau)^2}{10M} \|\mathbf{w}_O\|_2^2 \\ & \quad - 3M \cdot (|V_{n,i}^{(T_1)}| + (t - T_1 + 1) \eta (\max\{\|\boldsymbol{\mu}\|_2 \tau + \tau^2, \frac{2C_p^2 \sigma_p^2 d}{N}\}) \|\mathbf{w}_O\|_2^2) \\ & \geq 0, \end{aligned} \quad (19)$$

where the last inequality is by  $V_+^{(T_1)} \geq 3M \cdot |V_{n,i}^{(T_1)}|$  and requires  $N \cdot \text{SNR}^2 \geq 60M^2 C_p^2$ . The proof for  $V_-^{(t+1)} \leq -3M \cdot |V_{n,i}^{(t+1)}|$  is the same.

Next, we prove the upper bound for  $V_+$  and  $V_{n,i}$ . Based on the upper bound of attention ( $< 1$ ) and  $-\tilde{\ell}'_n \leq 1$  we have

$$\begin{aligned} \gamma_{V_+}^{(s+1)} & \leq \gamma_{V_+}^{(s)} - \frac{\eta}{NM} \sum_{n \in S_+} \tilde{\ell}'_n^{(s)} ((\|\boldsymbol{\mu}\|_2 + \tau)^2 + \sum_{j=2}^M (\|\boldsymbol{\mu}\|_2 + \tau)^2) \\ & \quad - \frac{\eta}{NM} \sum_{n \in S_+} \tilde{\ell}'_n^{(s)} (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \cdot M \\ & \leq \gamma_{V_+}^{(s)} + \frac{3\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2}{4} \\ & \leq \gamma_{V_+}^{(s)} + \eta C_4 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \end{aligned} \quad (20)$$

Then we can get that

$$\begin{aligned} |V_+^{(t+1)}| & \leq V_+^{(T_1)} + (\gamma_{V_+}^{(t+1)} - \gamma_{V_+}^{(T_1)}) \|\mathbf{w}_O\|_2 \\ & \leq V_+^{(T_1)} + \sum_{s=T_1}^t \eta C_4 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 \\ & \leq O(d_h^{-\frac{1}{4}}) + \eta C_4 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_1 + 1) \end{aligned} \quad (21)$$

where the first inequality is by the monotonicity of  $\gamma_{V_+}$  and the definition of  $V_+$ , the last inequality is by the result of stage 1 where  $V_+^{(T_1)} = O(d^{-1})$ . Similarly, we have

$$|V_-^{(t+1)}| \leq O(d_h^{-\frac{1}{4}}) + \eta C_4 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_1 + 1) \quad (22)$$

which completes the proof for the upper bound of  $V_+$ .

1836 Expanding (18) yields  
1837

$$1838 \quad |V_{n,i}^{(t+1)}| \leq |V_{n,i}^{(T_1)}| + \eta(\max\{\|\boldsymbol{\mu}\|_2 + \tau^2, \frac{2C_p^2\sigma_p^2d}{N}\}) \cdot \|\mathbf{w}_O\|_2^2(t - T_1 + 1) \quad (23)$$

$$1840 \quad \leq O(d_h^{-\frac{1}{4}}) + \eta C_4(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2(t - T_1 + 1)$$

1842 where the last inequality is by the result of phase 1 where  $|V_{n,i}^{(T_1)}| = O(d_h^{-\frac{1}{4}})$  and the condition that  
1843  $N \cdot \text{SNR}^2 \geq \Omega(1)$ .  
1844  
1845

## 1846 F.2.2 PROOF OF CLAIM 2

1847 By the results of I.6, we have the dynamic of  $\langle \mathbf{q}, \mathbf{k} \rangle$  as follows

$$1848 \quad \langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle$$

$$1849 \quad \geq \frac{\eta^2 C_6 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+}^{(s)})}, \quad (24)$$

$$1850 \quad \langle \mathbf{q}_-^{(s+1)}, \mathbf{k}_-^{(s+1)} \rangle - \langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle$$

$$1851 \quad \geq \frac{\eta^2 C_6 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,-}^{(s)})}, \quad (25)$$

$$1852 \quad \langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$$

$$1853 \quad \leq - \frac{\eta^2 C_6 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+}^{(s)})}, \quad (26)$$

$$1854 \quad \langle \mathbf{q}_-^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)} \rangle - \langle \mathbf{q}_-^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$$

$$1855 \quad \leq - \frac{\eta^2 C_6 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,-}^{(s)})}, \quad (27)$$

$$1856 \quad \langle \mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle$$

$$1857 \quad \geq \frac{\eta^2 C_6 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,i,+}^{(s)})}, \quad (28)$$

$$1858 \quad \langle \mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_-^{(s+1)} \rangle - \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_-^{(s)} \rangle$$

$$1859 \quad \geq \frac{\eta^2 C_6 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,i,-}^{(s)})}, \quad (29)$$

$$1860 \quad \langle \mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)} \rangle - \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$$

$$1861 \quad \leq - \frac{\eta^2 C_6 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,i,\pm,j}^{(s)})} \quad (30)$$

1862 for  $s \in [T_1, t]$ . The seven equations above show that  $\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle$ ,  $\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_\pm^{(s)} \rangle$  are monotonically  
1863 increasing and  $\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$ ,  $\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$  are monotonically decreasing. Next, we provide the loga-  
1864 rithmic increasing lower bounds of  $\Lambda_{n,\pm,j}^{(s+1)}$  and  $\Lambda_{n,\pm,j}^{(s+1)}$ .  
1865  
1866

We have

$$\begin{aligned}
\Lambda_{n,+j}^{(s+1)} - \Lambda_{n,+j}^{(s)} &= (\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) - (\langle \mathbf{q}_{n,j}^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)} \rangle - \langle \mathbf{q}_{n,j}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle) \\
&\geq \frac{\eta^2 C_6 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+j}^{(s)})} \\
&\quad + \frac{\eta^2 C_6 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)(s - T_1 + 1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+j}^{(s)})} \\
&\geq \frac{\eta^2 C_7 \max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+j}^{(s)})} \\
&\geq \frac{\eta^2 C_7 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}} (s - T_1)^2}{N (\log(6N^2 M^2 / \delta))^2} \cdot \frac{1}{\exp(\Lambda_{n,+j}^{(s)})}
\end{aligned} \tag{31}$$

where the last inequality is by  $\sigma_h^2 \geq \min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} d_h^{-\frac{1}{2}} (\log(6N^2 M^2 / \delta))^{-2}$ . Multiply both sides simultaneously by  $\exp(\Lambda_{n,+j}^{(s)})$  and get

$$\exp(\Lambda_{n,+j}^{(s)}) (\Lambda_{n,+j}^{(s+1)} - \Lambda_{n,+j}^{(s)}) \geq \frac{\eta^2 C_7 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}} (s - T_1)}{N (\log(6N^2 M^2 / \delta))^2} \cdot \frac{1}{\exp(\Lambda_{n,+j}^{(s)})}. \tag{32}$$

Taking a summation from  $T_1$  to  $t$  and get

$$\sum_{s=T_1}^t \exp(\Lambda_{n,+j}^{(s)}) (\Lambda_{n,+j}^{(s+1)} - \Lambda_{n,+j}^{(s)}) \geq \sum_{s=T_1}^t \frac{\eta^2 C_7 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}} (s - T_1)}{N (\log(6N^2 M^2 / \delta))^2} \tag{33}$$

$$\geq \frac{\eta^2 C_8 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (t - T_1)(t - T_1 + 1). \tag{34}$$

By the property that  $\Lambda_{n,+j}^{(t)}$  is monotonically increasing, we have

$$\int_{\Lambda_{n,+j}^{(T_1)}}^{\Lambda_{n,+j}^{(t)}} \exp(x) dx \geq \sum_{s=T_1}^t \exp(\Lambda_{n,+j}^{(s)}) (\Lambda_{n,+j}^{(s+1)} - \Lambda_{n,+j}^{(s)}) \geq \frac{\eta^2 C_8 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (t - T_1)(t - T_1 + 1). \tag{35}$$

By  $\int_{\Lambda_{n,+j}^{(T_1)}}^{\Lambda_{n,+j}^{(t+1)}} \exp(x) dx = \exp(\Lambda_{n,+j}^{(t+1)}) - \exp(\Lambda_{n,+j}^{(T_1)})$  we get

$$\Lambda_{n,+j}^{(t+1)} \geq \log \left( \exp(\Lambda_{n,+j}^{(T_1)}) + \frac{\eta^2 C_8 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (t - T_1)(t - T_1 + 1) \right).$$

Similarly, we have

$$\Lambda_{n,-j}^{(t+1)} \geq \log \left( \exp(\Lambda_{n,-j}^{(T_1)}) + \frac{\eta^2 C_8 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (t - T_1)(t - T_1 + 1) \right).$$

Thus, we have

$$\begin{aligned}
\Lambda_{n,i,+j}^{(s+1)} - \Lambda_{n,i,+j}^{(s)} &= \langle (\mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_+^{(s+1)}) - (\mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)}) \rangle - \langle (\mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)}) - (\mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)}) \rangle \\
&\geq \frac{\eta^2 C_6 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,i,+j}^{(s)})} \\
&\quad + \frac{\eta^2 C_6 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,i,+j}^{(s)})} \\
&\geq \frac{\eta^2 C_7 \max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,i,+j}^{(s)})} \\
&\geq \frac{\eta^2 C_7 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}} (s - T_1)}{N (\log(6N^2 M^2 / \delta))^2} \cdot \frac{1}{\exp(\Lambda_{n,i,+j}^{(s)})}. \tag{36}
\end{aligned}$$

Then using the similar method as for  $\Lambda_{n,i,+j}^{(t)}$ , we get

$$\Lambda_{n,i,+j}^{(t+1)} \geq \log \left( \exp(\Lambda_{n,i,+j}^{(T_1)}) + \frac{\eta^2 C_8 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (t - T_1)(t - T_1 + 1) \right),$$

which complete the proof. The proof for Claim 3 is in Section I.9.

### F.2.3 PROOF OF CLAIM 4

By the results of I.7, we have

$$\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle \leq \frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}$$

for  $s \in [T_1, t]$ . Further we have

$$\begin{aligned}
\exp(\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle) &\leq \exp \left( \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle + \frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \right) \\
&= \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) \cdot \exp \left( \frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \right) \\
&\leq C_{11} \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle). \tag{37}
\end{aligned}$$

For the last inequality, by  $\eta \leq \tilde{O}(\min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}})$ ,  $\sigma_h^2 \leq \min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} d_h^{-\frac{1}{2}} (\log(6N^2 M^2 / \delta))^{-\frac{3}{2}}$ ,  $\langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_+^{(T_1)} \rangle = o(1)$  and the monotonicity of  $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle$  for  $s \in [T_1, t]$ , we have  $\exp \left( \frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \right) \leq \exp(o(1)) \leq C_{11}$ .

Multiplying both sides by  $(\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)$  simultaneously gives

$$\begin{aligned}
&\exp(\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle) (\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) \\
&\leq C_{11} \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) \cdot (\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) \\
&\leq \eta C_{12} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h, \tag{38}
\end{aligned}$$

where the last inequality is by plugging (37). Taking a summation, we obtain

$$\begin{aligned}
& \int_{\langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_+^{(T_1)} \rangle}^{\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle} \exp(x) dx \\
& \leq \sum_{s=T_1}^t \exp(\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle) \left( \langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle \right) \\
& \leq \sum_{s=T_1}^t \eta C_{12} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h \\
& \leq T_2 \cdot \eta C_{12} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h \\
& \leq \frac{d_h^{1/2}}{\log^2(6N^2 M^2 / \delta)}.
\end{aligned} \tag{39}$$

where the first inequality is due to  $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle$  is monotone increasing, the last inequality is by  $T_2 = \Theta(\eta^{-1} \|\boldsymbol{\mu}\|_2^{-2} \|\mathbf{w}_O\|_2^{-2} \log(6N^2 M^2 / \delta)^{-1})$  and  $\sigma_h^2 \leq \min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} d_h^{-\frac{1}{2}} (\log(6N^2 M^2 / \delta))^{-\frac{3}{2}}$ . By

$$\int_{\langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_+^{(T_1)} \rangle}^{\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle} \exp(x) dx = \exp(\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle) - \exp(\langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_+^{(T_1)} \rangle),$$

we have

$$\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle \leq \log \left( \langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_+^{(T_1)} \rangle + \frac{d_h^{\frac{1}{2}}}{\log(6N^2 M^2 / \delta)} \right) \leq \log \left( d_h^{\frac{1}{2}} \right), \tag{40}$$

By the results of I.7, we also have

$$\langle \mathbf{q}_-^{(s+1)}, \mathbf{k}_-^{(s+1)} \rangle - \langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle \leq \frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h}{\exp(\langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle)} \tag{41}$$

$$\langle \mathbf{q}_\pm^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)} \rangle - \langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle \geq -\frac{\eta C_{10} \sigma_p^2 d (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h}{N} \cdot \exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle). \tag{42}$$

$$\langle \mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_\pm^{(s+1)} \rangle - \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_\pm^{(s)} \rangle \leq \frac{\eta C_{10} \sigma_p^2 d (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h}{N \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_\pm^{(s)} \rangle)}. \tag{43}$$

$$\langle \mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)} \rangle - \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle \geq -\frac{\eta C_{10} \sigma_p^2 d (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \sigma_h^2 d_h}{N} \cdot \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle). \tag{44}$$

Then using the similar method as for  $\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle$ , we get

$$\langle \mathbf{q}_-^{(t+1)}, \mathbf{k}_-^{(t+1)} \rangle \leq \log \left( d_h^{\frac{1}{2}} \right), \tag{45}$$

$$\langle \mathbf{q}_\pm^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle \geq -\log \left( d_h^{\frac{1}{2}} \right), \tag{46}$$

$$\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_\pm^{(t+1)} \rangle \leq \log \left( d_h^{\frac{1}{2}} \right), \tag{47}$$

$$\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle \geq -\log \left( d_h^{\frac{1}{2}} \right). \tag{48}$$

Next we provide the upper bound for  $|\langle \mathbf{q}_\pm^{(t+1)}, \mathbf{k}_\pm^{(t+1)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{n',j}^{(t+1)} \rangle|$ . By the results of I.8, we have

$$\sum_{s=T_1}^t |\alpha_{+,+}^{(s)}|, \sum_{s=T_1}^t |\alpha_{-,-}^{(s)}|, \sum_{s=T_1}^t |\beta_{+,+}^{(s)}|, \sum_{s=T_1}^t |\beta_{-,-}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,i,+}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,i,-}^{(s)}| = O \left( N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \right), \tag{49}$$

2052 for  $i \in [M] \setminus \{1\}, n \in S_{\pm}$ .

2053

2054

2055

2056

2057

$$\sum_{s=T_1}^t |\alpha_{n,+}^{(s)}|, \sum_{s=T_1}^t |\alpha_{n,-}^{(s)}| = O\left(N^{-\frac{1}{2}} d_h^{-\frac{1}{4}}\right), \quad (50)$$

2058 for  $i \in [M] \setminus \{1\}, n \in S_{\pm}$ .

2059

2060

2061

2062

$$\sum_{s=T_1}^t |\beta_{n,+}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,-}^{(s)}| = O\left(\text{SNR} \cdot N^{-\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \quad (51)$$

2063 for  $i \in [M] \setminus \{1\}, n \in S_{\pm}$ .

2064

2065

2066

2067

2068

$$\sum_{s=T_1}^t |\alpha_{n,i,n',j}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,j,n',i}^{(s)}| = O\left(d^{-\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2 / \delta)\right) \quad (52)$$

2069 for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], n \neq n'$ . Plugging these into the update rule of  
2070  $\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle, \langle \mathbf{q}_{n,i}, \mathbf{k}_{n,j}^{(t)} \rangle$  and assume that propositions  $\mathcal{C}(T_1), \dots, \mathcal{C}(t)$  hold, we have

2071

2072

2073

2074

2075

2076

2077

2078

2079

2080

2081

2082

2083

2084

2085

2086

2087

2088

2089

$$\begin{aligned} |\langle \mathbf{q}_{+}^{(t+1)}, \mathbf{k}_{-}^{(t+1)} \rangle| &\leq |\langle \mathbf{q}_{+}^{(T_1)}, \mathbf{k}_{-}^{(T_1)} \rangle| + \sum_{s=T_1}^t |\langle \mathbf{q}_{+}^{(s+1)}, \mathbf{k}_{-}^{(s+1)} \rangle - \langle \mathbf{q}_{+}^{(s)}, \mathbf{k}_{-}^{(s)} \rangle| \\ &\leq |\langle \mathbf{q}_{+}^{(T_1)}, \mathbf{k}_{-}^{(T_1)} \rangle| \\ &\quad + \sum_{s=T_1}^t \left| \alpha_{+,+}^{(s)} \langle \mathbf{k}_{+}^{(s)}, \mathbf{k}_{-}^{(s)} \rangle + \sum_{n \in S_{+}} \sum_{i=2}^M \alpha_{n,+}^{(s)} \langle \mathbf{k}_{n,i}^{(s)}, \mathbf{k}_{-}^{(s)} \rangle \right. \\ &\quad \left. + \beta_{-,-}^{(s)} \langle \mathbf{q}_{+}^{(s)}, \mathbf{q}_{-}^{(s)} \rangle + \sum_{n \in S_{-}} \sum_{i=2}^M \beta_{n,-}^{(s)} \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{q}_{+}^{(s)} \rangle \right. \\ &\quad \left. + \left( \alpha_{+,+}^{(s)} \mathbf{k}_{+}^{(s)} + \sum_{n \in S_{+}} \sum_{i=2}^M \alpha_{n,+}^{(s)} \mathbf{k}_{n,i}^{(s)} \right) \right. \\ &\quad \left. \cdot \left( \beta_{-,-}^{(s)\top} \mathbf{q}_{-}^{(s)\top} + \sum_{n \in S_{-}} \sum_{i=2}^M \beta_{n,-}^{(s)\top} \mathbf{q}_{n,i}^{(s)\top} \right) \right| \\ &\leq |\langle \mathbf{q}_{+}^{(T_1)}, \mathbf{k}_{-}^{(T_1)} \rangle| \\ &\quad + \sum_{s=T_1}^t |\alpha_{+,+}^{(s)}| |\langle \mathbf{k}_{+}^{(s)}, \mathbf{k}_{-}^{(s)} \rangle| + \sum_{n \in S_{+}} \sum_{i=2}^M \sum_{s=T_1}^t |\alpha_{n,+}^{(s)}| |\langle \mathbf{k}_{n,i}^{(s)}, \mathbf{k}_{-}^{(s)} \rangle| \\ &\quad + \sum_{s=T_1}^t |\beta_{-,-}^{(s)}| |\langle \mathbf{q}_{+}^{(s)}, \mathbf{q}_{-}^{(s)} \rangle| + \sum_{n \in S_{-}} \sum_{i=2}^M \sum_{s=T_1}^t |\beta_{n,-}^{(s)}| |\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{q}_{+}^{(s)} \rangle| \\ &\quad + \{\text{lower order term}\} \\ &= |\langle \mathbf{q}_{+}^{(T_1)}, \mathbf{k}_{-}^{(T_1)} \rangle| \\ &\quad + O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \cdot o(1) + N \cdot M \cdot O\left(N^{-\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \cdot o(1) \\ &\quad + O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \cdot o(1) + N \cdot M \cdot O\left(\text{SNR} \cdot N^{-\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \cdot o(1) \\ &= |\langle \mathbf{q}_{+}^{(T_1)}, \mathbf{k}_{-}^{(T_1)} \rangle| + o\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}}\right) + o\left(\text{SNR} \cdot N^{\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \\ &= o(1), \end{aligned} \quad (53)$$

2090

2091

2092

2093

2094

2095

2096

2097

2098

2099

2100

2101

2102

2103

2104

2105

where the first inequality is by triangle inequality, the last equality is by  $|\langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_-^{(T_1)} \rangle| = o(1)$  and  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$ . Similarly we have  $|\langle \mathbf{q}_-^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle| = o(1)$ .

$$\begin{aligned}
|\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{\bar{n},j}^{(t+1)} \rangle| &\leq |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{\bar{n},j}^{(T_1)} \rangle| + \sum_{s=T_1}^t |\langle \mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_{\bar{n},j}^{(s+1)} \rangle - \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle| \\
&\leq |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{\bar{n},j}^{(T_1)} \rangle| \\
&\quad + \sum_{s=T_1}^t \left| \alpha_{n,i,+}^{(s)} \langle \mathbf{k}_+^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle + \alpha_{n,i,-}^{(s)} \langle \mathbf{k}_-^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \alpha_{n,i,n',l}^{(s)} \langle \mathbf{k}_{n',l}^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle \right. \\
&\quad + \beta_{\bar{n},j,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle + \beta_{\bar{n},j,-}^{(s)} \langle \mathbf{q}_-^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \beta_{\bar{n},j,n',l}^{(s)} \langle \mathbf{q}_{n',l}^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle \\
&\quad + \left( \alpha_{n,i,+}^{(s)} \langle \mathbf{k}_+^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle + \alpha_{n,i,-}^{(s)} \langle \mathbf{k}_-^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \alpha_{n,i,n',l}^{(s)} \langle \mathbf{k}_{n',l}^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle \right) \\
&\quad \cdot \left( \beta_{\bar{n},j,+}^{(s)\top} \langle \mathbf{q}_+^{(s)\top}, \mathbf{q}_{n,i}^{(s)\top} \rangle + \beta_{\bar{n},j,-}^{(s)\top} \langle \mathbf{q}_-^{(s)\top}, \mathbf{q}_{n,i}^{(s)\top} \rangle + \sum_{n'=1}^N \sum_{l=2}^M \beta_{\bar{n},j,n',l}^{(s)\top} \langle \mathbf{q}_{n',l}^{(s)\top}, \mathbf{q}_{n,i}^{(s)\top} \rangle \right) \Big| \\
&\leq |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{\bar{n},j}^{(T_1)} \rangle| \\
&\quad + \sum_{s=T_1}^t |\alpha_{n,i,+}^{(s)}| |\langle \mathbf{k}_+^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle| + \sum_{s=T_1}^t |\alpha_{n,i,-}^{(s)}| |\langle \mathbf{k}_-^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle| \\
&\quad + \sum_{s=T_1}^t |\alpha_{n,i,\bar{n},j}^{(s)}| |\langle \mathbf{k}_{\bar{n},j}^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle| + \sum_{l=2}^M \sum_{s=T_1}^t |\alpha_{n,i,n,l}^{(s)}| |\langle \mathbf{k}_{n,l}^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle| \\
&\quad + \sum_{n' \neq \bar{n} \wedge (l \neq j \vee n' \neq \bar{n})} \sum_{s=T_1}^t \sum_{n'=1}^N \sum_{l=2}^M |\alpha_{n,i,n',l}^{(s)}| |\langle \mathbf{k}_{n',l}^{(s)}, \mathbf{k}_{\bar{n},j}^{(s)} \rangle| \\
&\quad + \sum_{s=T_1}^t |\beta_{\bar{n},j,+}^{(s)}| |\langle \mathbf{q}_+^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle| + \sum_{s=T_1}^t |\beta_{\bar{n},j,-}^{(s)}| |\langle \mathbf{q}_-^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle| \\
&\quad + \sum_{s=T_1}^t |\beta_{\bar{n},j,n,i}^{(s)}| |\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle| + \sum_{l=2}^M \sum_{s=T_1}^t |\beta_{\bar{n},j,n,l}^{(s)}| |\langle \mathbf{q}_{n,l}^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle| \\
&\quad + \sum_{n' \neq \bar{n} \wedge (l \neq i \vee n' \neq n)} \sum_{s=T_1}^t \sum_{n'=1}^N \sum_{l=2}^M |\beta_{\bar{n},j,n',l}^{(s)}| |\langle \mathbf{q}_{n',l}^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle| \\
&\quad + \{\text{lower order term}\} \\
&= |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{\bar{n},j}^{(T_1)} \rangle| \\
&\quad + O(d_h^{-\frac{1}{4}}) \cdot o(1) + O(d_h^{-\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2 / \delta)) \cdot \Theta(\sigma_p^2 \sigma_h^2 d d_h) \\
&\quad + M \cdot O(d_h^{-\frac{1}{4}}) \cdot o(1) + N \cdot M \cdot O(d_h^{-\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2 / \delta)) \cdot o(1) \\
&\quad + O(N^{\frac{1}{2}} d_h^{-\frac{1}{4}}) \cdot o(1) \\
&= |\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{\bar{n},j}^{(T_1)} \rangle| + o(d_h^{-\frac{1}{4}}) \\
&\quad + O(d_h^{-\frac{1}{2}} d_h^{\frac{1}{4}}) + o(N d_h^{-\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2 / \delta)) \\
&= o(1),
\end{aligned} \tag{54}$$

where the first inequality is by triangle inequality, the second equality is by  $\sigma_h^2 \leq \min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}} \cdot (\log(6N^2 M^2 / \delta))^{-\frac{3}{2}}$ , the last equality is by  $|\langle \mathbf{q}_{n,i}^{(T_1)}, \mathbf{k}_{n',j}^{(T_1)} \rangle| = o(1)$ ,  $d = \tilde{\Omega}(\epsilon^{-2} N^2 d_h)$  and  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$ .

### F.3 STAGE III

In Stage III, the outputs of ViT grow up and the loss derivatives are no longer at  $o(1)$ . We will carefully compute the growth rate of  $V_{\pm}$  and  $V_{n,i}$  while keeping monitoring the monotonicity of  $\langle q, k \rangle$ . By substituting  $t = T_2 = \Theta\left(\frac{1}{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right)$  into propositions  $\mathcal{B}(t)$ ,  $\mathcal{C}(t)$ ,  $\mathcal{D}(t)$ ,  $\mathcal{E}(t)$  in Stage II, we have the following conditions at the beginning of stage III

$$\begin{aligned} |V_+^{(T_2)}|, |V_-^{(T_2)}|, |V_{n,i}^{(T_2)}| &= o(1), \\ V_+^{(T_2)} &\geq 3M \cdot |V_{n,i}^{(T_2)}|, \\ V_-^{(T_2)} &\leq -3M \cdot |V_{n,i}^{(T_2)}|, \\ \|\mathbf{q}_+^{(T_2)}\|_2^2, \|\mathbf{k}_+^{(T_2)}\|_2^2 &= \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h), \\ \|\mathbf{q}_{n,i}^{(T_2)}\|_2^2, \|\mathbf{k}_{n,i}^{(T_2)}\|_2^2 &= \Theta(\sigma_p^2 \sigma_h^2 d d_h), \\ |\langle \mathbf{q}_+^{(T_2)}, \mathbf{q}_-^{(T_2)} \rangle|, |\langle \mathbf{q}_+^{(T_2)}, \mathbf{q}_{n,i}^{(T_2)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_2)}, \mathbf{q}_{n',j}^{(T_2)} \rangle| &= o(1), \\ |\langle \mathbf{k}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle|, |\langle \mathbf{k}_+^{(T_2)}, \mathbf{k}_{n,i}^{(T_2)} \rangle|, |\langle \mathbf{k}_{n,i}^{(T_2)}, \mathbf{k}_{n',j}^{(T_2)} \rangle| &= o(1), \end{aligned}$$

for  $i, j \in [M] \setminus \{1\}$ ,  $n, n' \in [N]$ ,  $i \neq j$  or  $n \neq n'$ .

$$\begin{aligned} \Lambda_{n,\pm,j}^{(T_2)} &\geq \log\left(\exp(\Lambda_{n,\pm,j}^{(T_1)}) + \Theta\left(\frac{d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2 / \delta))^3}\right)\right) \\ \Lambda_{n,i,\pm,j}^{(T_2)} &\geq \log\left(\exp(\Lambda_{n,i,\pm,j}^{(T_1)}) + \Theta\left(\frac{\sigma_p^2 d d_h^{\frac{1}{2}}}{N\|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}\right)\right) \\ |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle|, |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_{n,j}^{(T_2)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_2)}, \mathbf{k}_{n',j}^{(T_2)} \rangle| &\leq \log(d_h^{\frac{1}{2}}) \\ |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_2)}, \mathbf{k}_{n,j}^{(T_2)} \rangle| &= o(1) \end{aligned}$$

for  $i, j \in [M] \setminus \{1\}$ ,  $n, \bar{n} \in [N]$ ,  $n \neq \bar{n}$ .

Let  $T_3 = \Theta\left(\frac{1}{\eta\epsilon(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right)$ . Next we prove the following four propositions  $\mathcal{F}(t)$ ,  $\mathcal{G}(t)$ ,  $\mathcal{H}(t)$ ,  $\mathcal{I}(t)$  by induction on  $t$  for  $t \in [T_2, T_3]$ :

- $\mathcal{F}(t)$ :

$$\begin{aligned} V_+^{(t)} &\geq 3M \cdot |V_{n,i}^{(t)}|, \\ V_-^{(t)} &\leq -3M \cdot |V_{n,i}^{(t)}|, \\ |V_{n,i}^{(t)}| &= o(1), \end{aligned}$$

$$\begin{aligned} \log\left(\exp(V_+^{(T_2)}) + \eta C_{17}(\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2)\right) &\leq V_+^{(t)} \leq 2 \log\left(O\left(\frac{1}{\epsilon}\right)\right), \\ -2 \log\left(O\left(\frac{1}{\epsilon}\right)\right) &\leq V_-^{(t)} \leq -\log\left(\exp(-V_-^{(T_2)}) + \eta C_{17}(\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2)\right) \end{aligned}$$

for  $i \in [M] \setminus \{1\}$ ,  $n \in [N]$ .

- $\mathcal{G}(t)$ :

$$\|\mathbf{q}_{\pm}^{(t)}\|_2^2, \|\mathbf{k}_{\pm}^{(t)}\|_2^2 = \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h),$$

$$\begin{aligned}
& \|\mathbf{q}_{n,i}^{(t)}\|_2^2, \|\mathbf{k}_{n,i}^{(t)}\|_2^2 = \Theta(\sigma_p^2 \sigma_h^2 dd_h), \\
& |\langle \mathbf{q}_+^{(t)}, \mathbf{q}_-^{(t)} \rangle|, |\langle \mathbf{q}_\pm^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',j}^{(t)} \rangle| = o(1), \\
& |\langle \mathbf{k}_+^{(t)}, \mathbf{k}_-^{(t)} \rangle|, |\langle \mathbf{k}_\pm^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| = o(1)
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], i \neq j$  or  $n \neq n'$ .

•  $\mathcal{H}(t)$ :

$$\begin{aligned}
& \langle \mathbf{q}_\pm^{(t+1)}, \mathbf{k}_\pm^{(t+1)} \rangle \geq \langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle, \\
& \langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_\pm^{(t+1)} \rangle \geq \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle, \\
& \langle \mathbf{q}_\pm^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle \leq \langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle, \\
& \langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle \leq \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n \in [N]$ .

•  $\mathcal{I}(t)$ :

$$\begin{aligned}
& |\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle| \leq \log(\epsilon^{-1} d_h^{\frac{1}{2}}), \\
& |\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\mp^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{\bar{n},j}^{(t)} \rangle| = o(1)
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], n \neq \bar{n}$ .

By the results of Stage II, we know that  $\mathcal{F}(T_2), \mathcal{G}(T_2), \mathcal{I}(T_2)$  are true. To prove that  $\mathcal{F}(t), \mathcal{G}(t), \mathcal{H}(t)$  and  $\mathcal{I}(t)$  are true in stage 3, we will prove the following claims holds for  $t \in [T_2, T_3]$ :

- Claim 5.  $\mathcal{H}(T_2), \dots, \mathcal{H}(t-1), \mathcal{I}(T_2), \dots, \mathcal{I}(t) \implies \mathcal{F}(t+1)$
- Claim 6.  $\mathcal{F}(t), \mathcal{G}(t), \mathcal{H}(T_2), \dots, \mathcal{H}(t-1), \mathcal{I}(T_2), \dots, \mathcal{I}(t-1) \implies \mathcal{H}(t)$
- Claim 7.  $\mathcal{F}(T_2), \dots, \mathcal{F}(t), \mathcal{G}(t), \mathcal{H}(T_2), \dots, \mathcal{H}(t-1), \mathcal{I}(T_2), \dots, \mathcal{I}(t) \implies \mathcal{G}(t+1)$
- Claim 8.  $\mathcal{F}(T_2), \dots, \mathcal{F}(t), \mathcal{G}(T_2), \dots, \mathcal{G}(t), \mathcal{H}(T_2), \dots, \mathcal{H}(t-1), \mathcal{I}(T_2), \dots, \mathcal{I}(t) \implies \mathcal{I}(t+1)$

### F.3.1 PROOF OF CLAIM 5

The proofs for  $V_+^{(t)} \geq 3M \cdot |V_{n,i}^{(t)}|$  and  $V_-^{(t)} \leq -3M \cdot |V_{n,i}^{(t)}|$  are the same as for F.2.1. Based on  $\mathcal{H}(T_2), \dots, \mathcal{H}(t)$  where  $\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle$  and  $\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_\pm^{(s)} \rangle$  are monotonically non-decreasing and  $\max_j \langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle, \max_j \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$  are monotonically non-increasing for  $s \in [T_2, t-1]$ , we have

$$\Lambda_{n,\pm,j}^{(s)} \geq \Lambda_{n,\pm,j}^{(T_2)} \geq \log \left( \exp(\Lambda_{n,\pm,j}^{(T_1)}) + \Theta \left( \frac{d_h^{\frac{1}{2}}}{N(\log(6N^2M^2/\delta))^3} \right) \right), \quad (55)$$

$$\Lambda_{n,i,\pm,j}^{(s)} \geq \Lambda_{n,i,\pm,j}^{(T_2)} \geq \log \left( \exp(\Lambda_{n,i,\pm,j}^{(T_1)}) + \Theta \left( \frac{\sigma_p^2 dd_h^{\frac{1}{2}}}{N \|\boldsymbol{\mu}\|_2^2 (\log(6N^2M^2/\delta))^3} \right) \right) \quad (56)$$

for  $i, j \in [M] \setminus \{1\}, n \in [N], s \in [T_2, t]$ . We further get

$$\begin{aligned}
& \frac{\exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle)}{\exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \leq \frac{\exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{C \exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle)} = \frac{1}{C \exp(\Lambda_{n,\pm,j}^{(s)})} \\
& \leq \frac{1}{C \exp(\Lambda_{n,\pm,j}^{(T_1)}) + \Theta \left( \frac{d_h^{\frac{1}{2}}}{N(\log(6N^2M^2/\delta))^3} \right)} = O \left( \frac{N(\log(6N^2M^2/\delta))^3}{d_h^{\frac{1}{2}}} \right). \quad (57)
\end{aligned}$$

For the first inequality, by the monotonicity of  $\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle$  ( $\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle$  is increasing and  $\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$  is decreasing), there exist a constant  $C$  such that  $C \exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle) \geq \exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle) +$

2268  $\sum_{j'=2}^M \exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)$ . The second inequality is by plugging 55. Similarly, we have

$$2269 \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \leq \frac{1}{C \exp(\Lambda_{n,i,\pm,j}^{(s)})}$$

$$2270 = O\left(\frac{N \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}{\sigma_p^2 d d_h^{\frac{1}{2}}}\right). \quad (58)$$

2276 Plugging (57) and (58) into the update rule of  $\rho_{V,n,i}$  in Lemma 8 and get

$$2277 |\rho_{V,n,i}^{(s+1)} - \rho_{V,n,i}^{(s)}| \leq \frac{\eta}{NM} |\tilde{\ell}_n^{(s)}| \cdot \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle \cdot \left( O\left(\frac{N (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{3}{2}}}\right) + O\left(\frac{N \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}{\sigma_p^2 d d_h^{\frac{1}{2}}}\right) \right)$$

$$2278 + \frac{\eta}{NM} \sum_{n' \neq n, i' \neq i} |\tilde{\ell}_{n'}^{(t)}| \cdot \langle \tilde{\boldsymbol{\xi}}_{n,i}, \tilde{\boldsymbol{\xi}}_{n',i'}^{(t)} \rangle \cdot \left( O\left(\frac{N (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{3}{2}}}\right) + O\left(\frac{N \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}{\sigma_p^2 d d_h^{\frac{1}{2}}}\right) \right) + o(1)$$

$$2279 \leq \frac{3\eta(\sigma_p^2 d + o(1))}{2NM} \left( O\left(\frac{N (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{1}{2}}}\right) + O\left(\frac{N \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}{\sigma_p^2 d d_h^{\frac{1}{2}}}\right) \right)$$

$$2280 + \frac{\eta}{NM} NM (2\sigma_p^2 \sqrt{d \log(4N^2 M^2 / \delta)} \sqrt{d} + o(1)) \cdot \left( O\left(\frac{N (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{1}{2}}}\right) + O\left(\frac{N \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}{\sigma_p^2 d d_h^{\frac{1}{2}}}\right) \right)$$

$$2281 \leq \frac{2\eta}{NM} \left( O\left(\frac{N (\sigma_p^2 d + o(1)) (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{3}{2}}}\right) + O\left(\frac{N \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{3}{2}}}\right) \right)$$

$$2282 = O\left(\frac{\eta(\sigma_p^2 d + o(1)) (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{1}{2}}} + \frac{\eta \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{1}{2}}}\right) \quad (59)$$

2296 where the second inequality is by Lemma 13 and  $|\tilde{\ell}_n^{(t)}| \leq 1$ . For the last inequality, since  $d =$

2297  $\tilde{\Omega}(\epsilon^{-2} N^2 d_h)$ , we have  $N \cdot M \cdot 2\sigma_p^2 \sqrt{d \log(4N^2 M^2 / \delta)} \leq \frac{1}{2} \sigma_p^2 d$ . By Definition 4 and taking a

2298 summation we have

$$2299 |V_{n,i}^{(t+1)}| \leq |V_{n,i}^{(T_2)}| + \sum_{s=T_2}^t |\rho_{V,n,i}^{(s+1)} - \rho_{V,n,i}^{(s)}| \cdot \|\mathbf{w}_O\|_2^2$$

$$2300 \leq |V_{n,i}^{(T_2)}| + T_3 \cdot |\rho_{V,n,i}^{(t+1)} - \rho_{V,n,i}^{(t)}| \cdot \|\mathbf{w}_O\|_2^2$$

$$2301 \leq o(1) + \Theta\left(\frac{1}{\eta \epsilon (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right) \cdot O\left(\frac{\eta(\sigma_p^2 d + o(1)) (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{1}{2}}} + \frac{\eta \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{1}{2}}}\right) \cdot \|\mathbf{w}_O\|_2^2$$

$$2302 = o(1) + O\left(\frac{(\sigma_p^2 d + o(1)) (\log(6N^2 M^2 / \delta))^3}{\epsilon (\|\boldsymbol{\mu}\|_2 + \tau)^2 d_h^{\frac{1}{2}}} + \frac{(\log(6N^2 M^2 / \delta))^3}{\epsilon d_h^{\frac{1}{2}}}\right)$$

$$2303 = o(1) + o(1)$$

$$2304 = o(1), \quad (60)$$

2312 where the first equality is by  $N \cdot \text{SNR}^2 \geq \Omega(1)$ , the second equality is by  $d_h =$

2313  $\tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$ . Then we have a constant upper bound for the sum of  $V_{n,i}$  as

2314 follows:

$$2315 \sum_{i \in [M] \setminus \{1\}} |V_{n,i}^{(s)}| = (M-1) \cdot o(1) \leq C_{15},$$

2317 for  $n \in [N], s \in [T_2, t]$ .

2319 Expanding (57) and (58), we have

$$2320 \frac{\exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} = O\left(\frac{N (\log(6N^2 M^2 / \delta))^3}{d_h^{3/2}}\right) = o(1), \quad (61)$$

where the equality is by  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\}N^2\epsilon^{-2})$ .

Similarly,

$$\frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} = O\left(\frac{N\|\boldsymbol{\mu}\|_2^2(\log(6N^2M^2/\delta))^3}{\sigma_p^2 d d_h^{1/2}}\right) = o(1), \quad (62)$$

where the equality is by the same choice of  $d_h$ .

Then we have

$$\text{softmax}(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) \geq 1 - (M-1) \cdot o(1) \geq 1 - o(1), \quad (63)$$

$$\text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) \geq 1 - (M-1) \cdot o(1) \geq 1 - o(1), \quad (64)$$

for  $i \in [M] \setminus \{1\}$ ,  $n \in [N]$ ,  $s \in [T_2, t]$ .

Thus, in the following discussion, we omit the noise-related terms as  $o(1)$ .

Next we provide the bounds for  $-\tilde{\ell}_n^{(s)}$ . Note that  $\ell(z) = \log(1 + \exp(-z))$  and  $-\ell'(z) = \exp(-z)/(1 + \exp(-z))$ . Without loss of generality, assume  $y_n = 1$ . We have

$$\begin{aligned} -\ell'(f(\tilde{\mathbf{X}}_n, \theta(s))) &= \frac{1}{1 + \exp\left(\frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^{(s)\top} \mathbf{W}_Q^{(s)} \mathbf{W}_K^{(s)\top} (\tilde{\mathbf{X}}_n^{(s)})^\top \tilde{\mathbf{X}}_n^{(s)} \mathbf{W}_V^{(s)} \mathbf{w}_O)\right)} \\ &= \frac{1}{M} \left( \text{softmax}(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) + \sum_{l=2}^M \text{softmax}(\langle \mathbf{q}_{n,l}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) \right) \cdot \tilde{\boldsymbol{\mu}}_+^{(s)\top} \mathbf{W}_V^{(s)} \mathbf{w}_O \\ &\quad + \sum_{j \in [M] \setminus \{1\}} \left( \text{softmax}(\langle \mathbf{q}_{n,j}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle) + \sum_{l=2}^M \text{softmax}(\langle \mathbf{q}_{n,l}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle) \right) \cdot \tilde{\boldsymbol{\xi}}_{n,j}^{(s)\top} \mathbf{W}_V^{(s)} \mathbf{w}_O \\ &= \frac{1}{M} \left( M \cdot (1 - o(1)) \cdot V_+^{(s)} + M \cdot o(1) \cdot \sum_{j \in [M] \setminus \{1\}} V_{n,i}^{(s)} \right) \\ &\geq \frac{1}{2} V_+^{(s)}, \end{aligned} \quad (65)$$

for  $s \in [T_2, t]$ , where the second equality is by plugging equations above, and the inequality follows from  $V_+^{(s)} \geq 3M \cdot |V_{n,i}^{(s)}|$ .

Similarly, we have

$$\frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^{(s)\top} \mathbf{W}_Q^{(s)} \mathbf{W}_K^{(s)\top} (\tilde{\mathbf{X}}_n^{(s)})^\top \tilde{\mathbf{X}}_n^{(s)} \mathbf{W}_V^{(s)} \mathbf{w}_O) \leq \max_{i \in [M] \setminus \{1\}} \{V_+^{(s)}, V_{n,i}^{(s)}\} = V_+^{(s)}. \quad (66)$$

Then we have

$$\begin{aligned} -\ell'(f(\tilde{\mathbf{X}}_n, \theta(s))) &= \frac{1}{1 + \exp\left(\frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^{(s)\top} \mathbf{W}_Q^{(s)} \mathbf{W}_K^{(s)\top} (\tilde{\mathbf{X}}_n^{(s)})^\top \tilde{\mathbf{X}}_n^{(s)} \mathbf{W}_V^{(s)} \mathbf{w}_O)\right)} \\ &\geq \frac{1}{1 + \exp(V_+^{(s)})} \\ &\geq \frac{C_{16}}{\exp(V_+^{(s)})} \end{aligned} \quad (67)$$

where the first inequality is by plugging (66). For the last inequality, note that  $V_+^{(T_2)} \geq 0$  and  $V_+^{(s)}$  is monotonically increasing, so there exist a constant  $C_{16}$  such that  $\frac{1}{1 + \exp(V_+^{(s)})} \geq \frac{C_{16}}{\exp(V_+^{(s)})}$ . We also

2376 have the upper bound  
2377

$$\begin{aligned}
2378 \quad -\ell'(f(\mathbf{X}_n, \theta(s))) &= \frac{1}{1 + \exp\left(\frac{1}{M} \sum_{l=1}^M \varphi(\mathbf{x}_{n,l}^\top \mathbf{W}_Q^{(s)} \mathbf{W}_K^{(s)\top} (\mathbf{X}_n)^\top) \mathbf{X}_n \mathbf{W}_V^{(s)} \mathbf{w}_O\right)} \\
2379 &\leq \frac{1}{1 + \exp(V_+^{(s)}/2)} \\
2380 &\leq \frac{1}{\exp(V_+^{(s)}/2)}
\end{aligned} \tag{68}$$

2386 Then by the update rule of  $\gamma_{V_+,+}^{(t)}$  and in Lemma 8 and get  
2387

$$\begin{aligned}
2388 \quad \gamma_{V_+,+}^{(s+1)} - \gamma_{V_+,+}^{(s)} &= -\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(s)} \rangle \sum_{n \in S_+} \tilde{\ell}_n^{(s)} \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)} \right. \\
2389 &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \mathbf{q}_{n,j}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,j}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{n,j}^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)} \right) + o(1) \\
2390 &\geq -\eta (\|\boldsymbol{\mu}\|_2 - \tau)^2 \sum_{n \in S_+} \tilde{\ell}_n^{(s)} (M \cdot (1 - o(1))) \\
2391 &\geq \eta (\|\boldsymbol{\mu}\|_2 - \tau)^2 \cdot \frac{N}{4} \cdot (1 - o(1)) \cdot \frac{C_{16}}{\exp(V_+^{(s)})} \\
2392 &\geq \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \frac{1}{\exp(V_+^{(s)})}
\end{aligned} \tag{69}$$

2402 where the second inequality is by (67). Then by definition 4, we get

$$2403 \quad V_+^{(s+1)} - V_+^{(s)} = (\gamma_{V_+,+}^{(s+1)} - \gamma_{V_+,+}^{(s)}) \|\mathbf{w}_O\|_2^2 \geq \frac{\eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2}{\exp(V_+^{(s)})} \tag{70}$$

2406 Multiply both sides simultaneously by  $\exp(V_+^{(s)})$  and get

$$2407 \quad \exp(V_+^{(s)}) (V_+^{(s+1)} - V_+^{(s)}) \geq \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 \tag{71}$$

2409 Taking a summation from  $T_2$  to  $t$  and get

$$\begin{aligned}
2410 \quad \sum_{s=T_2}^t \exp(V_+^{(s)}) (V_+^{(s+1)} - V_+^{(s)}) &\geq \sum_{s=T_2}^t \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 \\
2411 &\geq \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1)
\end{aligned} \tag{72}$$

2415 By the property that  $V_+^{(s)}$  is monotonically increasing, we have

$$\begin{aligned}
2416 \quad \int_{V_+^{(T_2)}}^{V_+^{(t+1)}} \exp(x) dx &\geq \sum_{s=T_2}^t \exp(V_+^{(s)}) (V_+^{(s+1)} - V_+^{(s)}) \\
2417 &\geq \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1)
\end{aligned} \tag{73}$$

2422 By  $\int_{V_+^{(T_2)}/2}^{V_+^{(t+1)}/2} \exp(x) dx = \exp(V_+^{(t+1)}/2) - \exp(V_+^{(T_2)}/2)$  we get

$$\begin{aligned}
2423 \quad V_+^{(t+1)} &\geq \log\left(\exp(V_+^{(T_2)}/2) + \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1)\right) \\
2424 &\geq \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1)
\end{aligned} \tag{74}$$

2428 Similarly, we have

$$2429 \quad V_-^{(t+1)} \leq -\log\left(\exp(V_-^{(T_2)}) + \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1)\right) \tag{75}$$

Next we provide upper bounds for  $V_+^{(t+1)}$  and  $V_-^{(t+1)}$ . By the update rule of  $\gamma_{V_+}^{(t)}$  and in Lemma 8 we have

$$\begin{aligned}
\gamma_{V_+}^{(s+1)} - \gamma_{V_+}^{(s)} &= -\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(s)} \rangle \sum_{n \in S_+} \tilde{\ell}_n^{(s)} \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)} \right. \\
&\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \mathbf{q}_{n,j}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,j}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{k=2}^M \exp(\langle \mathbf{q}_{n,j}^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)} \right) + o(1) \\
&\leq \eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sum_{n \in S_+} (-\tilde{\ell}_n^{(s)}) \cdot M \\
&\leq \eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 \cdot \frac{3N}{4} \cdot \exp(-V_+^{(s)}/2) \\
&= \frac{3\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2}{4} \exp(-V_+^{(s)}/2).
\end{aligned} \tag{76}$$

where the second inequality is by (68). Then by definition 4, we get

$$\begin{aligned}
V_+^{(s+1)} - V_+^{(s)} &= (\gamma_{V_+}^{(s+1)} - \gamma_{V_+}^{(s)}) \|\mathbf{w}_O\|_2^2 \\
&\leq \frac{3\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}{4 \exp(V_+^{(s)}/2)}
\end{aligned} \tag{77}$$

Further we have

$$\begin{aligned}
\exp(V_+^{(s+1)}/2) &\leq \exp(V_+^{(s)}/2) + \frac{3\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}{8 \exp(V_+^{(s)}/2)} \\
&= \exp(V_+^{(s)}/2) \cdot \exp\left(\frac{3\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}{8 \exp(V_+^{(s)}/2)}\right) \\
&\leq C_{18} \exp(V_+^{(s)}/2)
\end{aligned} \tag{78}$$

For the last inequality, by  $\eta \leq \tilde{O}(\min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}})$ ,  $V_+^{(T_2)} = \Theta(1)$  and the monotonicity of  $V_+^{(s)}$ , we have  $\exp(\frac{3\eta \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2}{8 \exp(V_+^{(s)}/2)}) \leq C_{18}$ . Multiplying both sides by  $(V_+^{(s+1)}/2 - V_+^{(s)}/2)$  simultaneously gives

$$\begin{aligned}
\exp(V_+^{(s)}) (V_+^{(s+1)}/2 - V_+^{(s)}/2) &\leq C_{18} \exp(V_+^{(s)}/2) (V_+^{(s+1)}/2 - V_+^{(s)}/2) \\
&\leq \frac{3\eta C_{18} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}{8}
\end{aligned} \tag{79}$$

where the last inequality is by plugging (78). Taking a summation we have

$$\begin{aligned}
\int_{V_+^{(T_2)}/2}^{V_+^{(t+1)}/2} \exp(x) dx &\leq \sum_{s=T_2}^{T_3} \exp(V_+^{(s+1)}/2) (V_+^{(s+1)}/2 - V_+^{(s)}/2) \\
&\leq \sum_{s=T_2}^{T_3} \frac{3\eta C_{18} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}{8} \\
&\leq \Theta\left(\frac{1}{\eta \epsilon (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right) \cdot \frac{3\eta C_{18} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}{8} = O\left(\frac{1}{\epsilon}\right)
\end{aligned} \tag{80}$$

By  $\int_{V_+^{(T_2)}/2}^{V_+^{(t+1)}/2} \exp(x) dx = \exp(V_+^{(t+1)}/2) - \exp(V_+^{(T_2)}/2)$  we have

$$V_+^{(t+1)} \leq 2 \log \left( \exp(V_+^{(T_2)}/2) + O\left(\frac{1}{\epsilon}\right) \right) = 2 \log \left( O\left(\frac{1}{\epsilon}\right) \right)$$

2484 Similarly, we have

$$2485 V_-^{(t+1)} \geq -2 \log \left( O \left( \frac{1}{\epsilon} \right) \right)$$

### 2488 F.3.2 PROOF OF CLAIM 6

2489 By  $\mathcal{H}(T_2), \dots, \mathcal{H}(t-1)$ , we have  $\text{softmax}(\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle), \text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle) = 1 - o(1)$  and  
 2491  $\text{softmax}(\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle), \text{softmax}(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle) = o(1)$ , which have been proved in F.3.1. By the  
 2492 results of I.5, we have the signs of  $\alpha$  and  $\beta$  as follows:

$$2493 \alpha_{+,+}^{(t)}, \alpha_{-,-}^{(t)}, \beta_{+,+}^{(t)}, \beta_{-,-}^{(t)}, \alpha_{n,i,+}^{(t)}, \alpha_{n,i,-}^{(t)}, \beta_{n,+i}^{(t)}, \beta_{n,-i}^{(t)} \geq 0,$$

$$2494 \alpha_{n,+i}^{(t)}, \alpha_{n,-i}^{(t)}, \alpha_{n,i,n,j}^{(t)}, \beta_{n,i,+}^{(t)}, \beta_{n,i,-}^{(t)}, \beta_{n,j,n,i}^{(t)} \leq 0.$$

2497 Then combined with  $\mathcal{G}(T)$  and we have the dynamics of  $\langle \mathbf{q}, \mathbf{k} \rangle$  as follows:

$$2498 \begin{aligned} 2499 \langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle &= \alpha_{+,+}^{(t)} \|\mathbf{k}_+^{(t)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+i}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle \\ 2500 &+ \beta_{+,+}^{(t)} \|\mathbf{q}_+^{(t)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+i}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle \\ 2501 &+ \left( \alpha_{+,+}^{(t)} \mathbf{k}_+^{(t)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+i}^{(t)} \mathbf{k}_{n,i}^{(t)} \right) \\ 2502 &\cdot \left( \beta_{+,+}^{(t)} \mathbf{q}_+^{(t)} + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+i}^{(t)} \mathbf{q}_{n,i}^{(t)} \right)^\top \\ 2503 &= \alpha_{+,+}^{(t)} \|\mathbf{k}_+^{(t)}\|_2^2 + \beta_{+,+}^{(t)} \|\mathbf{q}_+^{(t)}\|_2^2 + \{\text{lower order term}\} \\ 2504 &\geq 0 \end{aligned} \quad (81)$$

2515 Similarly, we have

$$2516 \langle \mathbf{q}_\pm^{(t+1)}, \mathbf{k}_\pm^{(t+1)} \rangle - \langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle \geq 0,$$

$$2517 \langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_\pm^{(t+1)} \rangle - \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle \geq 0,$$

$$2518 \langle \mathbf{q}_\pm^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle - \langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle \leq 0,$$

$$2519 \langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle - \langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle \leq 0$$

2522 which completes the proof. The proof for Claim 7 is in Section I.12

### 2524 F.3.3 PROOF OF CLAIM 8

2525 By the results of I.10, we have

$$2526 \langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \leq \frac{\eta C_{10} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h \log(O(\frac{1}{\epsilon}))}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}, \quad (82)$$

2530 Further we have

$$2531 \exp(\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle) \leq \exp \left( \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle + \frac{\eta C_{10} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h \log(O(\frac{1}{\epsilon}))}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)} \right)$$

$$2532 = \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) \cdot \exp \left( \frac{\eta C_{10} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h \log(O(\frac{1}{\epsilon}))}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)} \right)$$

$$2533 \leq C_{11} \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle).$$

For the last inequality, by  $\eta \leq \tilde{O}(\min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}})$ ,  $\sigma_h^2 \leq \min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}} \cdot (\log(6N^2 M^2 / \delta))^{-\frac{3}{2}}$ ,  $\langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_+^{(T_1)} \rangle = o(1)$  and the monotonicity of  $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle$  for  $s \in [T_1, t]$ , we have  $\exp\left(\frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^4 \sigma_h^2 d_h \log(O(\frac{1}{\epsilon}))}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}\right) \leq \exp(o(1)) \leq C_{11}$ . Multiplying both sides by  $(\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)$  simultaneously gives

$$\begin{aligned} & \exp(\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle) \left( \langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \right) \\ & \leq C_{11} \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) \cdot \left( \langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \right) \\ & \leq \eta C_{12} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h \log\left(O\left(\frac{1}{\epsilon}\right)\right), \end{aligned} \quad (83)$$

where the last inequality is by plugging (82). Taking a summation we have

$$\begin{aligned} \int_{\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle}^{\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle} \exp(x) dx & \leq \sum_{s=T_2}^t \exp(\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle) \left( \langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle \right) \\ & \leq \sum_{s=T_2}^t \eta C_{12} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h \log\left(O\left(\frac{1}{\epsilon}\right)\right) \\ & \leq T_3 \cdot \eta C_{12} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h \log\left(O\left(\frac{1}{\epsilon}\right)\right) \\ & = O\left(\frac{d_h^{\frac{1}{2}} \log\left(O\left(\frac{1}{\epsilon}\right)\right)}{\epsilon (\log(6N^2 M^2 / \delta))^{\frac{3}{2}}}\right). \end{aligned} \quad (84)$$

where the first inequality is due to  $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle$  is monotone increasing, the last equality is by  $T_3 = \Theta(\eta^{-1} \epsilon^{-1} (\|\boldsymbol{\mu}\|_2 + \tau)^{-2} \|\mathbf{w}_O\|_2^{-2})$ ,  $\|\mathbf{w}_O\|_2^2 = \Theta(1)$  and  $\sigma_h^2 \leq \min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}} \cdot (\log(6N^2 M^2 / \delta))^{-\frac{3}{2}}$ . By  $\int_{\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle}^{\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle} \exp(x) dx = \exp(\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle) - \exp(\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle)$ , we have

$$\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle \leq \log\left(\exp(\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle) + O\left(\frac{d_h^{\frac{1}{2}} \log\left(O\left(\frac{1}{\epsilon}\right)\right)}{\epsilon (\log(6N^2 M^2 / \delta))^{\frac{3}{2}}}\right)\right) \leq \log(\epsilon^{-1} d_h^{\frac{1}{2}})$$

where the last inequality is by  $\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle \leq \log(d_h^{\frac{1}{2}})$ . By the results of I.10, we also have

$$\langle \mathbf{q}_-^{(t+1)}, \mathbf{k}_-^{(t+1)} \rangle - \langle \mathbf{q}_-^{(t)}, \mathbf{k}_-^{(t)} \rangle \leq \frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h \log\left(O\left(\frac{1}{\epsilon}\right)\right)}{\exp(\langle \mathbf{q}_-^{(t)}, \mathbf{k}_-^{(t)} \rangle)}. \quad (85)$$

$$\langle \mathbf{q}_\pm^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)} \rangle - \langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle \geq -\frac{\eta C_{10} \sigma_p^2 d (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h \log\left(O\left(\frac{1}{\epsilon}\right)\right)}{N} \cdot \exp(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle). \quad (86)$$

$$\langle \mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_\pm^{(s+1)} \rangle - \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_\pm^{(s)} \rangle \leq \frac{\eta C_{10} \sigma_p^2 d (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h \log\left(O\left(\frac{1}{\epsilon}\right)\right)}{N \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_\pm^{(s)} \rangle)}. \quad (87)$$

$$\begin{aligned} \langle \mathbf{q}_{n,i}^{(s+1)}, \mathbf{k}_{n,j}^{(s+1)} \rangle - \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle & \geq -\frac{\eta C_{10} \sigma_p^2 d (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \sigma_h^2 d_h \log\left(O\left(\frac{1}{\epsilon}\right)\right)}{N} \\ & \quad \cdot \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle). \end{aligned} \quad (88)$$

Then using the similar method as for  $\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle$ , we get

$$\begin{aligned}
\langle \mathbf{q}_-^{(t+1)}, \mathbf{k}_-^{(t+1)} \rangle &\leq \log(\epsilon^{-1} d_h^{\frac{1}{2}}), \\
\langle \mathbf{q}_\pm^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle &\geq -\log(\epsilon^{-1} d_h^{\frac{1}{2}}), \\
\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_\pm^{(t+1)} \rangle &\leq \log(\epsilon^{-1} d_h^{\frac{1}{2}}), \\
\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{n,j}^{(t+1)} \rangle &\geq -\log(\epsilon^{-1} d_h^{\frac{1}{2}}),
\end{aligned} \tag{89}$$

Next we provide the upper bound for  $|\langle \mathbf{q}_\pm^{(t+1)}, \mathbf{k}_\pm^{(t+1)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{n',j}^{(t+1)} \rangle|$ . By the results of I.11, we have

$$\sum_{s=T_2}^t |\beta_{n,+,i}^{(t)}|, \sum_{s=T_2}^t |\beta_{n,-,i}^{(t)}| = O\left(\frac{\text{SNR}^2(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}}\right), \tag{90}$$

for  $i \in [M] \setminus \{1\}, n \in S_\pm$ .

$$\begin{aligned}
&\sum_{s=T_2}^t |\alpha_{+,+}^{(t)}|, \sum_{s=T_2}^t |\alpha_{-,-}^{(t)}|, \sum_{s=T_2}^t |\beta_{+,+}^{(t)}|, \sum_{s=T_2}^t |\beta_{-,-}^{(t)}|, \sum_{s=T_2}^t |\alpha_{n,i,+}^{(t)}|, \sum_{s=T_2}^t |\beta_{n,i,-}^{(t)}| \\
&= O\left(\frac{N(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}}\right),
\end{aligned} \tag{91}$$

for  $i \in [M] \setminus \{1\}, n \in S_\pm$ .

$$\begin{aligned}
&\sum_{s=T_2}^t |\alpha_{n,+,i}^{(t)}|, \sum_{s=T_2}^t |\alpha_{n,-,i}^{(t)}|, \sum_{s=T_2}^t |\alpha_{n,i,+}^{(t)}|, \sum_{s=T_2}^t |\alpha_{n,i,-}^{(t)}|, \sum_{s=T_2}^t |\alpha_{n,i,n,j}^{(t)}|, \sum_{s=T_2}^t |\beta_{n,j,n,i}^{(t)}| \\
&= O\left(\frac{(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}}\right),
\end{aligned} \tag{92}$$

for  $i, j \in [M] \setminus \{1\}, n \in S_\pm$ .

$$\sum_{s=T_2}^t |\alpha_{n,i,n',j}^{(t)}|, \sum_{s=T_2}^t |\beta_{n,j,n',i}^{(t)}| = O\left(\frac{(\log(6N^2M^2/\delta))^4 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}}\right) \tag{93}$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], n \neq n'$ . Plugging these and proposition  $\mathcal{G}(t)$  into the update rule of  $|\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\mp^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{\bar{n},j}^{(t)} \rangle|$  and get

$$\begin{aligned}
|\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_-^{(t+1)} \rangle| &\leq |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle| + \sum_{s=T_2}^t |\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_-^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_-^{(s)} \rangle| \\
&\leq |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle| \\
&\quad + \sum_{s=T_2}^t \left| \alpha_{+,+}^{(s)} \langle \mathbf{k}_+^{(s)}, \mathbf{k}_-^{(s)} \rangle + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \langle \mathbf{k}_{n,i}^{(s)}, \mathbf{k}_-^{(s)} \rangle \right. \\
&\quad \left. + \beta_{-,-}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{q}_-^{(s)} \rangle + \sum_{n \in S_-} \sum_{i=2}^M \beta_{n,-}^{(s)} \langle \mathbf{q}_{n,i}^{(s)}, \mathbf{q}_-^{(s)} \rangle \right| \\
&\quad + \left( \alpha_{+,+}^{(s)} \mathbf{k}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \mathbf{k}_{n,i}^{(s)} \right) \\
&\quad \cdot \left( \beta_{-,-}^{(s)} \mathbf{q}_-^{(s)\top} + \sum_{n \in S_-} \sum_{i=2}^M \beta_{n,-}^{(s)} \mathbf{q}_{n,i}^{(s)\top} \right) \Big| \\
&\leq |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle| \\
&\quad + \sum_{s=T_2}^t |\alpha_{+,+}^{(s)}| |\langle \mathbf{k}_+^{(s)}, \mathbf{k}_-^{(s)} \rangle| + \sum_{n \in S_+} \sum_{i=2}^M \sum_{s=T_2}^t |\alpha_{n,+}^{(s)}| |\langle \mathbf{k}_{n,i}^{(s)}, \mathbf{k}_-^{(s)} \rangle| \\
&\quad + \sum_{s=T_2}^t |\beta_{-,-}^{(s)}| |\langle \mathbf{q}_+^{(s)}, \mathbf{q}_-^{(s)} \rangle| + \sum_{n \in S_-} \sum_{i=2}^M \sum_{s=T_2}^t |\beta_{n,-}^{(s)}| |\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{q}_-^{(s)} \rangle| \\
&\quad + \{\text{lower order term}\} \\
&= |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle| \\
&\quad + O \left( \frac{N(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}} \right) \cdot o(1) + N \cdot M \cdot O \left( \frac{\text{SNR}^2(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}} \right) \cdot o(1) \\
&\quad + O \left( \frac{N(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}} \right) \cdot o(1) + N \cdot M \cdot O \left( \frac{\text{SNR}^2(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}} \right) \cdot o(1) \\
&= |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle| + o \left( \frac{N(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}} \right) + o \left( \frac{N \cdot \text{SNR}^2(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{1}{2}}} \right) \\
&= o(1), \tag{94}
\end{aligned}$$

where the first inequality is by triangle inequality, the second inequality is by results in D.2, the last equality is by  $|\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle| = o(1)$  and  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\}N^2\epsilon^{-2})$ . Similarly we have  $|\langle \mathbf{q}_-^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle| = o(1)$  and  $|\langle \mathbf{q}_{n,i}^{(t+1)}, \mathbf{k}_{\bar{n},j}^{(t+1)} \rangle| = o(1)$ .

**Lemma 21** (Convergence of Training Loss, Lemma D.7 in (Jiang et al., 2024)). *There exist  $T = \frac{C_{19}}{\eta\epsilon(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}$  such that*

$$L_S(\theta(T)) \leq \epsilon \tag{95}$$

*Proof.* As we have the same conditions at the end of stage III as (Jiang et al., 2024), thus we have: Substituting  $t = T = \frac{C_{19}}{\eta\epsilon(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}$  into propositions  $\mathcal{F}(t)$  and get

$$V_+^{(t)} \geq \log \left( \exp(V_+^{(T_2)}) + \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2) \right)$$

$$\begin{aligned}
&\geq \log \left( \exp(V_+^{(T_2)}) + \frac{C_{20}}{\epsilon} \right) \\
&\geq \log \left( \frac{C_{20}}{\epsilon} \right),
\end{aligned}$$

$$|V_{n,i}^{(t)}| = O(1).$$

Thus, for  $n \in S_+$ , we bound  $f(\tilde{\mathbf{X}}_n, \theta(t))$  as follows :

$$f(\tilde{\mathbf{X}}_n, \theta) = \frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^\top \mathbf{W}_Q \mathbf{W}_K^\top (\tilde{\mathbf{X}}_n)^\top) \tilde{\mathbf{X}}_n \mathbf{W}_V \mathbf{w}_O \geq \log \left( \frac{1}{\epsilon} \right) \quad (96)$$

And

$$\begin{aligned}
\tilde{\ell}_n^{(t)} &= \log \left( 1 + \exp(-f(\tilde{\mathbf{X}}_n, \theta(t))) \right) \\
&\leq \exp(-f(\tilde{\mathbf{X}}_n, \theta(t))) \\
&\leq \exp \left( -\log \left( \frac{1}{\epsilon} \right) \right) \\
&\leq \epsilon.
\end{aligned}$$

Similarly, we have  $\tilde{\ell}_n^{(t)} \leq \epsilon$  for  $n \in S_-$ . Therefore, we have

$$L_S(\theta(T)) = \frac{1}{N} \sum_{n=1}^N \tilde{\ell}_n^{(t)} \leq \epsilon.$$

□

#### F.4 TEST ERROR

In this section, we denote clean  $V_+$ ,  $V_-$  and  $V_\xi$  as  $\boldsymbol{\mu}_+^\top \mathbf{W}_V \mathbf{w}_O$ ,  $\boldsymbol{\mu}_-^\top \mathbf{W}_V \mathbf{w}_O$  and  $\boldsymbol{\xi}^\top \mathbf{W}_V \mathbf{w}_O$ , and perturbed  $\tilde{V}_+$ ,  $\tilde{V}_-$  and  $\tilde{V}_\xi$  as  $\tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_V \mathbf{w}_O$ ,  $\tilde{\boldsymbol{\mu}}_-^\top \mathbf{W}_V \mathbf{w}_O$  and  $\tilde{\boldsymbol{\xi}}^\top \mathbf{W}_V \mathbf{w}_O$ .

##### F.4.1 CLEAN TEST ERROR

**Theorem 22.** *Under Assumption 1, in the theoretical analysis of test error in the second and third stages of benign overfitting, we define  $g(\boldsymbol{\xi})$  as  $V_\xi^{(t)} = \langle \boldsymbol{\xi}, \mathbf{W}_V^{(t)} \mathbf{w}_O \rangle$ . Then, we know that for any  $x \geq 0$ , if  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  is a Lipschitz function and  $c$  is a constant, the following inequality holds for the test loss.*

$$\mathbb{P} \left( \sum_{j=2}^M \alpha_j (g(\boldsymbol{\xi}_j) - \mathbb{E}g(\boldsymbol{\xi}_j)) \geq x \right) \leq \exp \left( - \frac{cx^2}{\sigma_p^2 (\sum_{j=2}^M \alpha_j^2) \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2^2} \right).$$

*Proof.* According to Theorem 5.2.2 in Vershynin (2018), we know that for any  $x \geq 0$ , if  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  is a Lipschitz function, it holds that

$$\mathbb{P} \left( \sum_{i=1}^N \alpha_i (g(\boldsymbol{\xi}_i) - \mathbb{E}g(\boldsymbol{\xi}_i)) \geq x \right) \leq \exp \left( - \frac{cx^2}{\sigma_p^2 (\sum_{i=1}^N \alpha_i^2) \|g\|_{\text{Lip}}^2} \right) \quad (97)$$

where  $g(\boldsymbol{\xi})$  is defined as  $|V_\xi^{(t)}| = |\langle \boldsymbol{\xi}, \mathbf{W}_V^{(t)} \mathbf{w}_O \rangle|$ , we have

$$\begin{aligned}
|g(\boldsymbol{\xi}) - g(\boldsymbol{\xi}')| &= \left| \left| \langle \boldsymbol{\xi}, \mathbf{W}_V^{(t)} \mathbf{w}_O \rangle \right| - \left| \langle \boldsymbol{\xi}', \mathbf{W}_V^{(t)} \mathbf{w}_O \rangle \right| \right| \\
&\leq \left| \langle \boldsymbol{\xi} - \boldsymbol{\xi}', \mathbf{W}_V^{(t)} \mathbf{w}_O \rangle \right|
\end{aligned}$$

$$\leq \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2 \|\boldsymbol{\xi} - \boldsymbol{\xi}'\|_2$$

So, we can get

$$\|g\|_{\text{Lip}} \leq \left\| \mathbf{W}_V^{(t)} \mathbf{w}_O \right\|_2. \quad (98)$$

By plugging (98) into (97), we can get the result.  $\square$

The following inequality holds according to the update rules of the  $V$  vector in Lemma 8, the first equality is derived from the update of triangle inequality, and the second equality is due to the initialization of the  $V$  vector.

$$\begin{aligned} \left\| \mathbf{W}_V^{(t)} \mathbf{w}_O \right\|_2 &\leq \left\| \mathbf{W}_V^{(0)} \mathbf{w}_O \right\|_2 + \sum_{t'=0}^{t-1} \left\| \mathbf{W}_V^{(t'+1)} \mathbf{w}_O - \mathbf{W}_V^{(t')} \mathbf{w}_O \right\|_2 \\ &= \left\| \mathbf{W}_V^{(0)} \mathbf{w}_O \right\|_2 + tO \left( \eta \cdot \max \left\{ \|\boldsymbol{\mu}\|_2, \sigma_p \sqrt{d} \right\} \cdot \|\mathbf{w}_O\|^2 \right) \\ &= O \left( \sigma_V \|\mathbf{w}_O\|_2 \sqrt{d} + t\eta \|\mathbf{w}_O\|^2 \max \left\{ \|\boldsymbol{\mu}\|_2, \sigma_p \sqrt{d} \right\} \right) \\ &\leq O \left( t\eta \|\mathbf{w}_O\|^2 \max \left\{ \|\boldsymbol{\mu}\|_2, \sigma_p \sqrt{d} \right\} \right) \end{aligned} \quad (99)$$

Since  $g(\boldsymbol{\xi})$  as  $|V_{\boldsymbol{\xi}}^{(t)}| = |\langle \boldsymbol{\xi}, \mathbf{W}_V^{(t)} \mathbf{w}_O \rangle|$ , and since  $\langle \boldsymbol{\xi}, \mathbf{W}_V^{(t)} \mathbf{w}_O \rangle \sim \mathcal{N}(0, \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2^2 \sigma_p^2)$ , so we can get:

$$\mathbb{E}g(\boldsymbol{\xi}) = \mathbb{E}|\langle \boldsymbol{\xi}, \mathbf{W}_V^{(t)} \mathbf{w}_O \rangle| = \sqrt{\frac{2}{\pi}} \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2 \sigma_p$$

The test error can be interpreted as the probability that the noise term dominates the signal term. Formally, this corresponds to the event that the cumulative contribution of the random perturbation exceeds the deterministic signal margin. After centralization, we can apply Theorem 22 to obtain a high-probability upper bound on the test error.

$$\begin{aligned} P(y(f(\theta, \mathbf{X})) \leq 0) &= P \left( \left[ \left( \sum_{i=1}^M S_{i,1} (V_+^{(t)} - V_-^{(t)}) + \sum_{j=2}^M \left( \sum_{i=1}^M S_{i,j} \right) V_{\boldsymbol{\xi}_j}^{(t)} \right) \leq 0 \right] \right) \\ &\leq P \left( \sum_{j=2}^M \left( \sum_{i=1}^M S_{i,j} \right) |V_{\boldsymbol{\xi}_j}^{(t)}| \geq \left( \sum_{i=1}^M S_{i,1} \right) (V_+^{(t)} - V_-^{(t)}) \right) \\ &= P \left( \sum_{j=2}^M \alpha_j (g(\boldsymbol{\xi}_j) - E\{g(\boldsymbol{\xi}_j)\}) \geq \alpha_1 (V_+^{(t)} - V_-^{(t)}) - \sigma_p \sqrt{\frac{2}{\pi}} \left( \sum_{j=2}^M \alpha_j \right) \left\| \mathbf{W}_V^{(t)} \mathbf{w}_O \right\|_2 \right) \\ &\leq \exp \left[ - \frac{c_2 \left( \sum_r \alpha_1 (V_+^{(t)} - V_-^{(t)}) - \sigma_p \sqrt{\frac{2}{\pi}} \left( \sum_{j=2}^M \alpha_j \right) \left\| \mathbf{W}_V^{(t)} \mathbf{w}_O \right\|_2 \right)^2}{\sigma_p^2 \left( \sum_{j=2}^M \alpha_j^2 \right) \left( \left\| \mathbf{W}_V^{(t)} \mathbf{w}_O \right\|_2 \right)^2} \right] \\ &\leq \exp \left( \frac{c_4}{\pi} \right) \cdot \exp \left[ - \frac{c_5}{2} \left( \frac{\sum_r \alpha_1 (V_+^{(t)} - V_-^{(t)})}{\sigma_p \sqrt{\sum_{j=2}^M \alpha_j^2} \left\| \mathbf{W}_V^{(t)} \mathbf{w}_O \right\|_2} \right)^2 \right] \end{aligned}$$

We denote  $\alpha_j$  as  $\sum_i S_{i,j}$ , where  $S_{i,j}$  is  $\text{softmax}(\langle \mathbf{q}_i^{(t)}, \mathbf{k}_j^{(t)} \rangle)$ . When the subscript is 1, it represents the signal, and when the subscript is from 2 to  $M$ , it represents noise.

Then, by the lower bound of  $V_{\pm}^{(t)}$  and upper bound of  $\|\mathbf{W}_V^{(t)}\mathbf{w}_O\|_2$ , we can further bound the test error with following inequality:

$$\begin{aligned} P(y(f(\theta, \mathbf{X}) \leq 0) &\leq \exp\left(\frac{c_4}{\pi}\right) \cdot \exp\left[-\frac{c_5}{2} \left(\frac{\sum_r \alpha_1 (V_+^{(t)} - V_-^{(t)})}{\sigma_p \sqrt{\sum_{j=2}^M \alpha_j^2} \|\mathbf{W}_V^{(t)}\mathbf{w}_O\|_2}\right)^2\right] \\ &\leq \exp\left(\frac{c_{12}}{\pi}\right) \exp\left[-\frac{c_{13}}{2} O\left(\frac{V_+^{(t)} - V_-^{(t)}}{\frac{\sigma_p (V_+^{(t)} - V_-^{(t)})}{\|\boldsymbol{\mu}\|_2}}\right)^2\right] \\ &= \exp\left(\frac{c_{12}}{\pi}\right) \exp\left[-\frac{c_{13}}{2} O(dSNR^2)\right] \end{aligned}$$

where the second inequality is by  $\mathbf{W}_V^{(t)}\mathbf{w}_O$  is almost aligned with  $V_+^{(t)}$  and  $V_-^{(t)}$ , thus  $\|\mathbf{W}_V^{(t)}\mathbf{w}_O\|_2 \sim \frac{V_+^{(t)} - V_-^{(t)}}{\|\boldsymbol{\mu}\|_2}$ .

#### F.4.2 ROBUST TEST ERROR

We start by writing the prediction score under perturbation  $\tilde{X}$  as

$$yf(\theta, \tilde{X}) = \left(\sum_{i=1}^M \tilde{S}_{i1}\right)(\tilde{V}_+^{(t)} + \tilde{V}_-^{(t)}) + \sum_{j=2}^M \left(\sum_{i=1}^M \tilde{S}_{ij}\right)\tilde{V}_{\xi_j}^{(t)},$$

We first take the first term as an example and derive an upper bound on the maximum discrepancy between the perturbed input  $\tilde{X} \in B(X, \tau)$  and the clean input.

$$\begin{aligned} S_{11}V_+^{(t)} - \tilde{S}_{11}\tilde{V}_+^{(t)} &= (S_{11} - \tilde{S}_{11})V_+^{(t)} + \tilde{S}_{11}(V_+^{(t)} - \tilde{V}_+^{(t)}) \\ &\leq (1 - 1/C)S_{11}V_+^{(t)} + \tilde{S}_{11}|\langle \tilde{\boldsymbol{\mu}}_+ - \boldsymbol{\mu}_+, \mathbf{W}_V^{(t)}\mathbf{w}_O \rangle| \\ &\leq (1 - 1/C)S_{11}V_+^{(t)} + \tilde{S}_{11}\|\mathbf{W}_V^{(t)}\mathbf{w}_O\|\tau, \end{aligned}$$

where the inequality comes from Lemma 20, and the definition of  $V_+^{(t)}$ . This shows that the deviation can be bounded linearly in the perturbation magnitude, with only a small residual term since  $(C - 1) = o(1)$ . Similarly, an analogous upper bound holds for  $\tilde{S}_{i1}\tilde{V}_{\pm}^{(t)} - S_{i1}V_{\pm}^{(t)}$  and  $\tilde{S}_{ij}\tilde{V}_{\xi_j}^{(t)} - S_{ij}V_{\xi_j}^{(t)}$  for  $i \in [M]$ ,  $j \in [M] \setminus \{1\}$ .

Aggregating the deviations across all components, the worst-case perturbation satisfies

$$\begin{aligned} yf(\theta, X) - \min_{\tilde{X} \in B(X, \tau)} yf(\theta, \tilde{X}) &\leq \left(\sum_i S_{i1}\right)\|\mathbf{W}_V^{(t)}\mathbf{w}_O\|\tau + \sum_{j=2}^M \left(\sum_i S_{ij}\right)\|\mathbf{W}_V^{(t)}\mathbf{w}_O\|\tau + (1 - 1/C)S_{11}(\tilde{V}_+^{(t)} + \tilde{V}_-^{(t)}) + o(1) \\ &\lesssim \left(\sum_{j=1}^M \sum_i S_{ij}\right)\|\mathbf{W}_V^{(t)}\mathbf{w}_O\|\tau + (1 - 1/C)\left(\sum_i S_{i1}\right)(\tilde{V}_+^{(t)} + \tilde{V}_-^{(t)}) \\ &= M\|\mathbf{W}_V^{(t)}\mathbf{w}_O\|\tau + (1 - 1/C)\left(\sum_i S_{i1}\right)(\tilde{V}_+^{(t)} + \tilde{V}_-^{(t)}) \end{aligned} \tag{100}$$

where the first inequality comes from  $\tilde{V}_{\xi_j}^{(t)} = o(1)$  for  $j \in [M] \setminus \{1\}$ . Thus the adversarial effect scales with both the cumulative magnitude of the perturbed coefficients and the operator norm of the weight matrices. Then we can bound the robust test error:

$$P\left(\min_{\tilde{X} \in B(X, \tau)} yf(\theta, \tilde{X}) \leq 0\right) = P\left(yf(\theta, X) + \left(yf(\theta, X) - \min_{\tilde{X} \in B(X, \tau)} yf(\theta, \tilde{X})\right) \leq 0\right)$$

2862  
2863  
2864  
2865  
2866  
2867  
2868  
2869  
2870  
2871  
2872  
2873  
2874  
2875  
2876  
2877  
2878  
2879  
2880  
2881  
2882  
2883  
2884  
2885  
2886  
2887  
2888  
2889  
2890  
2891  
2892  
2893  
2894  
2895  
2896  
2897  
2898  
2899  
2900  
2901  
2902  
2903  
2904  
2905  
2906  
2907  
2908  
2909  
2910  
2911  
2912  
2913  
2914  
2915

$$\begin{aligned}
&\leq P\left(\sum_{j=2}^M \alpha_j (g(\xi_j) - E\{g(\xi_j)\}) \geq \alpha_1 (V_+^{(t)} - V_-^{(t)}) - \sigma_p \sqrt{\frac{2}{\pi}} \left(\sum_{j=2}^M \alpha_j\right) \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2 \right. \\
&\quad \left. - M \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2 \tau - (1 - 1/C) \alpha_1 (V_+^{(t)} - V_-^{(t)})\right) \\
&\leq \exp\left(\frac{c_{12}}{\pi}\right) \exp\left[-\frac{c_{13}}{2} \left(\frac{\alpha_1 (V_+^{(t)} - V_-^{(t)})}{\sigma_p \sqrt{\sum_{j=2}^M \alpha_j^2} \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2} - \frac{M \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2 \tau}{\sigma_p \sqrt{\sum_{j=2}^M \alpha_j^2} \|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2}\right)^2\right] \\
&\leq \exp\left(\frac{c_{12}}{\pi}\right) \exp\left[-\frac{c_{13}}{2} O\left(\sqrt{d} SNR \left(1 - \frac{\tau}{\|\boldsymbol{\mu}\|_2}\right)\right)^2\right]
\end{aligned}$$

The first inequality follows from (100). The third inequality uses the fact that  $C \leq e/2$ , which we absorb into the constant term. The last inequality follows from the bound on  $V_{\pm}^{(t)}$  and  $\|\mathbf{W}_V^{(t)} \mathbf{w}_O\|_2$ . This completes the proof.

## G BENIGN OVERFITTING IN CASE 2

Stage I stay same with F.1. Next, we aim to prove that under condition  $\tau = (1 - o(1))\|\boldsymbol{\mu}\|$ , the attention component in a ViT will remain in its initialization state and fail to learn meaningful signal-to-signal or noise-to-signal interactions. This is because, under such perturbations, any newly emerging margin can be immediately neutralized, preventing the signal-to-signal attention from accumulating advantages. In this regime, the ViT effectively degenerates into a linear model.

At the end of Stage I, since  $V_+^{T_1} \geq 3M|V_{n,i}^{T_1}|$  and  $q = o(1)$  at this point, the perturbation has no significant effect. Consequently,  $\langle \mathbf{q}_+, \mathbf{k}_+ \rangle$  and  $\langle \mathbf{q}_{n,i}, \mathbf{k}_+ \rangle$  experiences a temporary increase, but it does not exceed  $\Theta(\log C)$  (as we will demonstrate later). Therefore, we assume that after Stage II, when  $\langle q, k \rangle$  has stabilized and the loss derivatives are no longer at the  $o(1)$  scale. However,  $|V_{n,i}^{T_3}|$  is not  $o(1)$ ; it is of the same order as  $|V_+^{T_3}|$ , which implies that a larger SNR is required. The following conditions hold at the beginning of Stage III.

$$\begin{aligned}
&|V_+^{(T_2)}|, |V_-^{(T_2)}|, |V_{n,i}^{(T_2)}| = o(1), \\
&V_+^{(T_2)} \geq 3M \cdot |V_{n,i}^{(T_2)}|, \\
&V_-^{(T_2)} \leq -3M \cdot |V_{n,i}^{(T_2)}|, \\
&\|\mathbf{q}_+^{(T_2)}\|_2^2, \|\mathbf{k}_+^{(T_2)}\|_2^2 = \Theta(\log C), \\
&\|\mathbf{q}_{n,i}^{(T_2)}\|_2^2, \|\mathbf{k}_{n,i}^{(T_2)}\|_2^2 = \Theta(\log C), \\
&|\langle \mathbf{q}_+^{(T_2)}, \mathbf{q}_-^{(T_2)} \rangle|, |\langle \mathbf{q}_+^{(T_2)}, \mathbf{q}_{n,i}^{(T_2)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_2)}, \mathbf{q}_{n',j}^{(T_2)} \rangle| = o(1), \\
&|\langle \mathbf{k}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle|, |\langle \mathbf{k}_+^{(T_2)}, \mathbf{k}_{n,i}^{(T_2)} \rangle|, |\langle \mathbf{k}_{n,i}^{(T_2)}, \mathbf{k}_{n',j}^{(T_2)} \rangle| = o(1),
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], i \neq j$  or  $n \neq n'$ .

$$\begin{aligned}
&|\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle|, |\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_{n,j}^{(T_2)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_2)}, \mathbf{k}_+^{(T_2)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_2)}, \mathbf{k}_{n',j}^{(T_2)} \rangle| = \Theta(\log C) \\
&|\langle \mathbf{q}_+^{(T_2)}, \mathbf{k}_-^{(T_2)} \rangle|, |\langle \mathbf{q}_{n,i}^{(T_2)}, \mathbf{k}_{\bar{n},j}^{(T_2)} \rangle| = o(1)
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n, \bar{n} \in [N], n \neq \bar{n}$ .

Let  $T_3 = \Theta\left(\frac{M}{\eta \epsilon (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right)$ . Next we prove the following four propositions  $\mathcal{J}(t), \mathcal{K}(t), \mathcal{L}(t)$  by induction on  $t$  for  $t \in [T_2, T_3]$ :

2916  
2917  
2918  
2919  
2920  
2921  
2922  
2923  
2924  
2925  
2926  
2927  
2928  
2929  
2930  
2931  
2932  
2933  
2934  
2935  
2936  
2937  
2938  
2939  
2940  
2941  
2942  
2943  
2944  
2945  
2946  
2947  
2948  
2949  
2950  
2951  
2952  
2953  
2954  
2955  
2956  
2957  
2958  
2959  
2960  
2961  
2962  
2963  
2964  
2965  
2966  
2967  
2968  
2969

- $\mathcal{J}(t)$ :

$$\begin{aligned} V_+^{(t)} &\geq 3M \cdot |V_{n,i}^{(t)}|, \\ V_-^{(t)} &\leq -3M \cdot |V_{n,i}^{(t)}|, \\ |V_{n,i}^{(t)}| &= o(1), \end{aligned}$$

$$\begin{aligned} \log \left( \exp(V_+^{(T_2)}) + \frac{\eta}{M} C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2) \right) &\leq V_+^{(t)} \leq 2 \log \left( O \left( \frac{1}{\epsilon} \right) \right), \\ -2 \log \left( O \left( \frac{1}{\epsilon} \right) \right) &\leq V_-^{(t)} \leq -\log \left( \exp(-V_-^{(T_2)}) + \frac{\eta}{M} C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2) \right) \end{aligned}$$

for  $i \in [M] \setminus \{1\}, n \in [N]$ .

- $\mathcal{K}(t)$ :

$$\begin{aligned} \|\mathbf{q}_\pm^{(t)}\|_2^2, \|\mathbf{k}_\pm^{(t)}\|_2^2 &= \Theta(\log C), \\ \|\mathbf{q}_{n,i}^{(t)}\|_2^2, \|\mathbf{k}_{n,i}^{(t)}\|_2^2 &= \Theta(\log C), \\ |\langle \mathbf{q}_+^{(t)}, \mathbf{q}_-^{(t)} \rangle|, |\langle \mathbf{q}_\pm^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',j}^{(t)} \rangle| &= o(1), \\ |\langle \mathbf{k}_+^{(t)}, \mathbf{k}_-^{(t)} \rangle|, |\langle \mathbf{k}_\pm^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| &= o(1) \end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], i \neq j$  or  $n \neq n', C = O(1)$ .

- $\mathcal{L}(t)$ :

$$\begin{aligned} |\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_\pm^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle| &= \Theta(\log C), \\ |\langle \mathbf{q}_\pm^{(t)}, \mathbf{k}_\mp^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{\bar{n},j}^{(t)} \rangle| &= o(1) \end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], n \neq \bar{n}, C = O(1)$ .

By the results of Stage II, we know that  $\mathcal{F}(T_1), \mathcal{G}(T_2), \mathcal{I}(T_2)$  are true. To prove that  $\mathcal{F}(t), \mathcal{G}(t), \mathcal{H}(t)$  and  $\mathcal{I}(t)$  are true in stage III, we will prove the following claims holds for  $t \in [T_2, T_3]$ :

- Claim 9.  $\mathcal{L}(T_2), \dots, \mathcal{L}(t) \implies \mathcal{J}(t+1)$
- Claim 10.  $\mathcal{J}(t), \mathcal{L}(t), \mathcal{K}(t) \implies \mathcal{K}(t+1)$
- Claim 11.  $\mathcal{J}(t), \mathcal{K}(t), \mathcal{L}(t) \implies \mathcal{L}(t+1)$

### G.1 PROOF OF CLAIM 9

The proofs for  $V_+^{(t)} \geq 3M \cdot |V_{n,i}^{(t)}|$  and  $V_-^{(t)} \leq -3M \cdot |V_{n,i}^{(t)}|$  are the same as for F.2.1.

we provide the bounds for  $-\tilde{\ell}_n^{(s)}$ . Note that  $\ell(z) = \log(1 + \exp(-z))$  and  $-\ell'(z) = \exp(-z)/(1 + \exp(-z))$ . Without loss of generality, assume  $y_n = 1$ . We have

$$\begin{aligned} -\tilde{\ell}'(f(\tilde{\mathbf{X}}_n, \theta(s))) &= \frac{1}{1 + \exp\left(\frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^\top \mathbf{W}_Q^{(s)} \mathbf{W}_K^{(s)\top} (\tilde{\mathbf{X}}_n)^\top) \tilde{\mathbf{X}}_n \mathbf{W}_V^{(s)} \mathbf{w}_O\right)} \\ &= \frac{1}{M} \left( \text{softmax}(\langle \tilde{\mathbf{q}}_\pm^{(s)}, \tilde{\mathbf{k}}_\pm^{(s)} \rangle) + \sum_{l=2}^M \text{softmax}(\langle \tilde{\mathbf{q}}_{n,l}^{(s)}, \tilde{\mathbf{k}}_\pm^{(s)} \rangle) \right) \cdot \tilde{\boldsymbol{\mu}}_+^\top \mathbf{W}_V^{(s)} \mathbf{w}_O \\ &\quad + \sum_{j \in [M] \setminus \{1\}} \left( \text{softmax}(\langle \tilde{\mathbf{q}}_{n,j}^{(s)}, \tilde{\mathbf{k}}_{n,j}^{(s)} \rangle) + \sum_{l=2}^M \text{softmax}(\langle \tilde{\mathbf{q}}_{n,l}^{(s)}, \tilde{\mathbf{k}}_{n,j}^{(s)} \rangle) \right) \cdot \tilde{\boldsymbol{\xi}}_{n,j}^\top \mathbf{W}_V^{(s)} \mathbf{w}_O \\ &= \frac{1}{M} \left( M \cdot \frac{C}{C + M - 1} \cdot V_+^{(s)} + M \cdot \frac{1}{C + M - 1} \cdot \sum_{j \in [M] \setminus \{1\}} V_{n,i}^{(s)} \right) \\ &\geq \frac{1}{2M} V_+^{(s)}, \end{aligned} \tag{101}$$

for  $s \in [T_2, t]$ , where the second equality is by  $\mathcal{L}(t)$  that  $\Lambda_{n,\pm,j}^{(s)} = \Theta(\log C)$  and  $\Lambda_{n,i,\pm,j}^{(s)} = \Theta(\log C)$ , and the inequality follows from  $V_+^{(s)} \geq 3M \cdot |V_{n,i}^{(s)}|$ .

Similarly, we have

$$\begin{aligned} \frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^\top \mathbf{W}_Q^{(s)} \mathbf{W}_K^{(s)\top} (\tilde{\mathbf{X}}_n)^\top) \tilde{\mathbf{X}}_n \mathbf{W}_V^{(s)} \mathbf{w}_O &\leq \max_{i \in [M] \setminus \{1\}} \{V_+^{(s)}, V_{n,i}^{(s)}\} \\ &= V_+^{(s)}. \end{aligned} \quad (102)$$

Then, we have

$$\begin{aligned} -\tilde{\ell}'(f(\tilde{\mathbf{X}}_n, \theta(s))) &= \frac{1}{1 + \exp\left(\frac{1}{M} \sum_{l=1}^M \varphi(\tilde{\mathbf{x}}_{n,l}^\top \mathbf{W}_Q^{(s)} \mathbf{W}_K^{(s)\top} (\tilde{\mathbf{X}}_n)^\top) \tilde{\mathbf{X}}_n \mathbf{W}_V^{(s)} \mathbf{w}_O\right)} \\ &\geq \frac{1}{1 + \exp(V_+^{(s)})} \\ &\geq \frac{C_{16}}{\exp(V_+^{(s)})} \end{aligned} \quad (103)$$

For the last inequality, note that  $V_+^{(T_2)} \geq 0$  and  $V_+^{(s)}$  is monotonically increasing, so there exist a constant  $C_{16}$  such that  $\frac{1}{1 + \exp(V_+^{(s)})} \geq \frac{C_{16}}{\exp(V_+^{(s)})}$ . We also have the upper bound

$$\begin{aligned} -\ell'(f(X_n, \theta(s))) &= \frac{1}{1 + \exp\left(\frac{1}{M} \sum_{l=1}^M \varphi(x_{n,l}^\top \mathbf{W}_Q^{(s)} \mathbf{W}_K^{(s)\top} (X_n)^\top) X_n \mathbf{W}_V^{(s)} \mathbf{w}_O\right)} \\ &\leq \frac{1}{1 + \exp(V_+^{(s)}/2M)} \\ &\leq \frac{1}{\exp(V_+^{(s)}/2M)} \end{aligned} \quad (104)$$

By the update rule of  $\gamma_{V,+}^{(t)}$  and in Lemma 8 and get

$$\begin{aligned} \gamma_{V,+}^{(s+1)} - \gamma_{V,+}^{(s)} &= -\frac{\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle}{NM} \sum_{n \in S_+} \tilde{\ell}_n^{(t)} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\ &\quad + \sum_{n \in S_+} \tilde{\ell}_n^{(t)} \sum_{i=2}^M \frac{-\eta \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle}{NM} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right. \\ &\quad \left. + \sum_{j=2}^M \frac{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{k=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,j}^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)} \right) \\ &\geq -\frac{\eta(\|\boldsymbol{\mu}\|_2 - \tau)^2}{NM} \sum_{n \in S_+} \ell_n^{(s)} \left( M \cdot \frac{C}{C + M - 1} \right) + \frac{\eta\tau \|\boldsymbol{\mu}\|_2}{NM} \sum_{n \in S_+} \ell_n^{(s)} \left( M \cdot \frac{1}{C + M - 1} \right) \\ &\geq \frac{\eta(\|\boldsymbol{\mu}\|_2 - \tau)^2}{NM} \cdot \frac{N}{4} \cdot \frac{C_{16}}{\exp(V_+^{(s)})} - \frac{\eta\tau \|\boldsymbol{\mu}\|_2}{NM} \cdot \frac{N}{4} \cdot \frac{C_{16}}{\exp(V_+^{(s)}/2M)} \\ &\geq \frac{\eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2}{M} \frac{1}{\exp(V_+^{(s)})} \end{aligned}$$

where the second inequality is by (103)(104), the last inequality is by  $N \cdot SNR^2 = \Omega(\frac{1}{\epsilon})$ . Then by definition 4, we get

$$V_+^{(s+1)} - V_+^{(s)} = (\gamma_{V,+}^{(s+1)} - \gamma_{V,+}^{(s)}) \|\mathbf{w}_O\|_2^2 \geq \frac{\eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2}{M \exp(V_+^{(s)})} \quad (105)$$

Multiply both sides simultaneously by  $\exp(V_+^{(s)})$  and get

$$\exp(V_+^{(s)})(V_+^{(s+1)} - V_+^{(s)}) \geq \frac{\eta}{M} C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 \quad (106)$$

Taking a summation from  $T_2$  to  $t$  and get

$$\begin{aligned} \sum_{s=T_2}^t \exp(V_+^{(s)})(V_+^{(s+1)} - V_+^{(s)}) &\geq \sum_{s=T_2}^t \frac{\eta}{M} C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 \\ &\geq \frac{\eta}{M} C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1) \end{aligned} \quad (107)$$

By the property that  $V_+^{(s)}$  is monotonically increasing, we have

$$\begin{aligned} \int_{V_+^{(T_2)}}^{V_+^{(t+1)}} \exp(x) dx &\geq \sum_{s=T_2}^t \exp(V_+^{(s)})(V_+^{(s+1)} - V_+^{(s)}) \\ &\geq \frac{\eta}{M} C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1) \end{aligned} \quad (108)$$

By  $\int_{V_+^{(T_2)}}^{V_+^{(t+1)}} \exp(x) dx = \exp(V_+^{(t+1)}) - \exp(V_+^{(T_2)})$  we get

$$V_+^{(t+1)} \geq \log \left( \exp(V_+^{(T_2)}) + \frac{\eta}{M} C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1) \right) \quad (109)$$

Similarly, we have

$$V_-^{(t+1)} \leq -\log \left( \exp(V_-^{(T_2)}) + \eta C_{17} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\mathbf{w}_O\|_2^2 (t - T_2 + 1) \right) \quad (110)$$

Similarly, we have the lower bound, the proofs are the same as for F.3.1

## G.2 PROOF OF CLAIM 10

We first consider the increment of  $\langle \mathbf{q}, \mathbf{k} \rangle$  at the  $t$ -th update when using the original clean data. We then show that, at this step, the effect introduced by the perturbation under adversarial samples dominates the increment learned from the clean data; hence  $\langle \mathbf{q}, \mathbf{k} \rangle$  remains stable and bounded.

By the update rule of  $\langle \mathbf{q}, \mathbf{k} \rangle$  we have

$$\begin{aligned} \langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \\ = \alpha_{+,+}^{(t)} \|\mathbf{k}_+^{(t)}\|_2^2 + \beta_{+,+}^{(t)} \|\mathbf{q}_+^{(t)}\|_2^2 + \{\text{lower order term}\} \end{aligned} \quad (111)$$

Subsequently, we establish upper bounds for  $\alpha$  and  $\beta$  for clean data.

$$\begin{aligned}
\alpha_{+,+}^{(t)} &= \frac{\eta}{NM} \sum_{n \in S_+} -\ell'_n{}^{(t)} \|\boldsymbol{\mu}\|_2^2 \\
&\cdot \left( V_+^{(t)} \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \\
&\quad \left. - \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right)^2 \right. \\
&\quad \left. - \sum_{i=2}^M (V_{n,i}^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \\
&\quad \left. \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right) \\
&\leq \frac{\eta}{NM} \sum_{n \in S_+} \|\boldsymbol{\mu}\|_2^2 (V_+^{(t)}) \\
&\leq \frac{3\eta}{2M} \|\boldsymbol{\mu}\|_2^2 V_+^{(t)}
\end{aligned} \tag{112}$$

Then, we can then compute the incremental growth of  $\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle$  after one update step on clean data.

$$\begin{aligned}
&\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \\
&= \alpha_{+,+}^{(t)} \|\mathbf{k}_+^{(t)}\|_2^2 + \beta_{+,+}^{(t)} \|\mathbf{q}_+^{(t)}\|_2^2 + \{\text{lower order term}\} \\
&\leq \frac{3\eta}{2M} \|\boldsymbol{\mu}\|_2^2 V_+^{(t)} \Theta(\log C) + \frac{3\eta}{2M} \|\boldsymbol{\mu}\|_2^2 V_+^{(t)} \Theta(\log C)
\end{aligned} \tag{113}$$

Subsequently, we compute at the  $t^{\text{th}}$  iteration the magnitude of the effect that the perturbation imposes on  $\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle$ .

$$\begin{aligned}
&\max_{\tilde{\mathbf{X}}^{(t)} \in B(\mathbf{X}^{(t)}, \tau)} \langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \\
&= \left( \left( 1 + \frac{\tau}{\|\boldsymbol{\mu}\|_2} \right)^2 - 1 \right) \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \\
&= \left( \left( 1 + \frac{\tau}{\|\boldsymbol{\mu}\|_2} \right)^2 - 1 \right) \Theta(\log C)
\end{aligned} \tag{114}$$

As  $\tau$  and  $\|\boldsymbol{\mu}\|$  are same order, and  $\frac{3\eta}{2M} \|\boldsymbol{\mu}\|_2^2 V_+^{(t)} = o(\frac{1}{N})$  by  $V_+^{(t)} \leq 2 \log(O(\frac{1}{\epsilon}))$ ,  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$  and  $\eta \leq \tilde{O}(\min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}})$ . Thus, we have

$$\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \leq \max_{\tilde{\mathbf{X}}^{(t)} \in B(\mathbf{X}^{(t)}, \tau)} \langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle$$

which indicates that the perturbation's effect exceeds the one-step update on clean data. Therefore,  $\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle$  stay  $\Theta(\log C)$ . By similar methods, we can bound the other  $\langle \mathbf{q}^{(t)}, \mathbf{k}^{(t)} \rangle$ , which complete the proof.

### G.3 PROOF OF CLAIM 11

We use similar methods in G.2. We first consider the increment of  $\|\mathbf{q}\|_2^2$  and  $\|\mathbf{k}\|_2^2$  at the  $t$ -th update when using the original clean data. We then show that, at this step, the effect introduced by the perturbation under adversarial samples dominates the increment learned from the clean data; hence  $\|\mathbf{q}\|_2^2$  and  $\|\mathbf{k}\|_2^2$  remains stable and bounded.

By the update rule of  $\|\mathbf{q}\|_2^2$  and  $\|\mathbf{k}\|_2^2$  we have

$$\begin{aligned}
\|\mathbf{q}_+^{(t+1)}\|_2^2 - \|\mathbf{q}_+^{(t)}\|_2^2 &= 2\langle \Delta \mathbf{q}_+^{(t)}, \mathbf{q}_+^{(t)} \rangle + \langle \Delta \mathbf{q}_+^{(t)}, \Delta \mathbf{q}_+^{(t)} \rangle \\
&= 2\alpha_{+,+}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle + 2 \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle \\
&\quad + \left( \alpha_{+,+}^{(t)} \mathbf{k}_+^{(t)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(t)} \mathbf{k}_{n,i}^{(t)} \right) \\
&\quad \cdot \left( \alpha_{+,+}^{(t)} \mathbf{k}_+^{(t)\top} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(t)} \mathbf{k}_{n,i}^{(t)\top} \right) \\
&\leq 2|\alpha_{+,+}^{(t)}| |\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle| + 2 \sum_{n \in S_+} \sum_{i=2}^M |\alpha_{n,+}^{(t)}| |\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle| + \{\text{lower order term}\} \\
&\leq 2 \frac{3\eta}{2M} \|\boldsymbol{\mu}\|_2^2 \log \left( O\left(\frac{1}{\epsilon}\right) \right) \Theta(\log C) + 2 \sum_{n \in S_+} \sum_{i=2}^M \frac{3\eta}{2NM} \|\boldsymbol{\mu}\|_2^2 \log \left( O\left(\frac{1}{\epsilon}\right) \right) \Theta(\log C) \\
&\leq 12\eta \|\boldsymbol{\mu}\|_2^2 \log \left( O\left(\frac{1}{\epsilon}\right) \right) \Theta(\log C)
\end{aligned} \tag{115}$$

where the second inequality comes from the upper bounds for  $\alpha$  and  $\beta$  on clean data similar to (112).

Subsequently, we compute at the  $t^{\text{th}}$  iteration the magnitude of the effect that the perturbation imposes on  $\|\mathbf{q}_+^{(t)}\|_2^2$ .

$$\begin{aligned}
&\max_{\tilde{\mathbf{X}}^{(t)} \in B(\mathbf{X}^{(t)}, \tau)} \|\tilde{\mathbf{q}}_+^{(t)}\|_2^2 - \|\mathbf{q}_+^{(t)}\|_2^2 \\
&\leq \left( \left(1 + \frac{\tau}{\|\boldsymbol{\mu}\|_2}\right)^2 - 1 \right) \|\mathbf{q}_+^{(t)}\|_2^2 \\
&\leq \left( \left(1 + \frac{\tau}{\|\boldsymbol{\mu}\|_2}\right)^2 - 1 \right) \Theta(\log C)
\end{aligned} \tag{116}$$

As  $\tau$  and  $\|\boldsymbol{\mu}\|_2$  are same order, and  $12\eta \|\boldsymbol{\mu}\|_2^2 \log \left( O\left(\frac{1}{\epsilon}\right) \right) = o\left(\frac{1}{NM}\right)$  by  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$  and  $\eta \leq \tilde{O}(\min\{\|\boldsymbol{\mu}\|_2^{-2}, (\sigma_p^2 d)^{-1}\} \cdot d_h^{-\frac{1}{2}})$ . Thus, we have

$$\|\mathbf{q}_+^{(t+1)}\|_2^2 - \|\mathbf{q}_+^{(t)}\|_2^2 \leq \max_{\tilde{\mathbf{X}}^{(t)} \in B(\mathbf{X}^{(t)}, \tau)} \|\tilde{\mathbf{q}}_+^{(t)}\|_2^2 - \|\mathbf{q}_+^{(t)}\|_2^2$$

which indicates that the perturbation's effect exceeds the one-step update on clean data. Therefore,  $\|\mathbf{q}_+^{(t+1)}\|_2^2$  stay  $\Theta(\log C)$ . By similar methods, we can bound the other  $\|\mathbf{q}\|_2^2$  and  $\|\mathbf{k}\|_2^2$ , which complete the proof.

The proof of convergence and test error is similar with Lemma 21 and Section F.4

## H PROOF OF THEOREM 3

*Proof.* Fix an arbitrary  $\theta$ . Consider the two classes separately. For the positive class ( $y = +$ ):

$$\delta_1^{(+)} = -\boldsymbol{\mu}_+, \quad \delta_2^{(+)} = 0$$

is a valid perturbation since  $\|\delta_1^{(+)}\|_2 = \|\boldsymbol{\mu}_+\|_2 \leq \tau$ ,  $\|\delta_2^{(+)}\|_2 = 0 \leq \tau$ . Then the adversarially perturbed point

$$\tilde{\mathbf{x}}^{(+)} = [\boldsymbol{\mu}_+ - \boldsymbol{\mu}_+, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_M] = [0, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_M]$$

lies in  $\mathcal{B}([\boldsymbol{\mu}, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_M], \tau)$ , so there exists a perturbation that can potentially flip the classifier's output for the positive class.



3240  
3241  
3242  
3243  
3244  
3245  
3246  
3247  
3248  
3249  
3250  
3251  
3252  
3253  
3254  
3255  
3256  
3257  
3258  
3259  
3260  
3261  
3262  
3263  
3264  
3265  
3266  
3267  
3268  
3269  
3270  
3271  
3272  
3273  
3274  
3275  
3276  
3277  
3278  
3279  
3280  
3281  
3282  
3283  
3284  
3285  
3286  
3287  
3288  
3289  
3290  
3291  
3292  
3293

$$\begin{aligned}
\alpha_{n,+i}^{(t)} = & -\frac{\eta}{NM} \ell_n^{(t)} \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \\
& \cdot \left( -V_+^{(t)} \cdot \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right. \\
& + \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \\
& + V_{n,i}^{(t)} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right. \\
& \left. \left. - \left( \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right)^2 \right) \right) \\
& - \sum_{k \neq i} \left( V_{n,k}^{(t)} \cdot \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right. \\
& \left. \cdot \frac{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,k}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_+^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right) \\
& + \sum_{k=2}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,k}^{(t)} \rangle \cdot \left( -V_+^{(t)} \cdot \frac{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right. \\
& \left. \cdot \frac{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right) \\
& + V_{n,i}^{(t)} \left( \frac{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right. \\
& \left. - \left( \frac{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right)^2 \right) \\
& - \sum_{l \neq i} \left( V_{n,l}^{(t)} \cdot \frac{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,i}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right. \\
& \left. \cdot \frac{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,l}^{(t)} \rangle)}{\exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \tilde{\mathbf{q}}_{n,k}^{(t)}, \tilde{\mathbf{k}}_{n,j}^{(t)} \rangle)} \right)
\end{aligned}$$

We can also derive the calculations for other  $\alpha$  and  $\beta$ , since they follow the same procedure as in Section F.1 of Jiang et al. (2024).

## I.2 PROOF OF LEMMA 17

### I.3 UPDATE RULES FOR INNER PRODUCTS

In this subsection, we give the update rules for the inner products of  $\mathbf{q}$  and  $\mathbf{k}$ .

3294  
3295  
3296  
3297  
3298  
3299  
3300  
3301  
3302  
3303  
3304  
3305  
3306  
3307  
3308  
3309  
3310  
3311  
3312  
3313  
3314  
3315  
3316  
3317  
3318  
3319  
3320  
3321  
3322  
3323  
3324  
3325  
3326  
3327  
3328  
3329  
3330  
3331  
3332  
3333  
3334  
3335  
3336  
3337  
3338  
3339  
3340  
3341  
3342  
3343  
3344  
3345  
3346  
3347

$$\begin{aligned}
& \langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle \\
&= \alpha_{+,+}^{(t)} \|\mathbf{k}_+^{(t)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle \\
&\quad + \beta_{+,+}^{(t)} \|\mathbf{q}_+^{(t)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle \\
&\quad + \left( \alpha_{+,+}^{(t)} \mathbf{k}_+^{(t)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(t)} \mathbf{k}_{n,i}^{(t)} \right) \\
&\quad \cdot \left( \beta_{+,+}^{(t)} \mathbf{q}_+^{(t)\top} + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+}^{(t)} \mathbf{q}_{n,i}^{(t)\top} \right),
\end{aligned}$$

$$\begin{aligned}
& \langle \mathbf{q}_-^{(t+1)}, \mathbf{k}_-^{(t+1)} \rangle - \langle \mathbf{q}_-^{(t)}, \mathbf{k}_-^{(t)} \rangle \\
&= \alpha_{-,-}^{(t)} \|\mathbf{k}_-^{(t)}\|_2^2 + \sum_{n \in S_-} \sum_{i=2}^M \alpha_{n,-}^{(t)} \langle \mathbf{k}_-^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle \\
&\quad + \beta_{-,-}^{(t)} \|\mathbf{q}_-^{(t)}\|_2^2 + \sum_{n \in S_-} \sum_{i=2}^M \beta_{n,-}^{(t)} \langle \mathbf{q}_-^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle \\
&\quad + \left( \alpha_{-,-}^{(t)} \mathbf{k}_-^{(t)} + \sum_{n \in S_-} \sum_{i=2}^M \alpha_{n,-}^{(t)} \mathbf{k}_{n,i}^{(t)} \right) \\
&\quad \cdot \left( \beta_{-,-}^{(t)} \mathbf{q}_-^{(t)\top} + \sum_{n \in S_-} \sum_{i=2}^M \beta_{n,-}^{(t)} \mathbf{q}_{n,i}^{(t)\top} \right),
\end{aligned}$$

We can also derive the update rules for other  $\mathbf{q}$  and  $\mathbf{k}$ , since they follow the same procedure as in Section F.2 of Jiang et al. (2024).

## I.4 PROOF OF LEMMA 17

Let  $T_0 = O\left(\frac{1}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_0\|_2^2}\right)$ . By Lemma 16, we have  $|V_+^{(t)}|, |V_-^{(t)}|, |V_{n,i}^{(t)}| = O(d_h^{-\frac{1}{4}})$  for  $t \in [0, T_0]$  by Lemma 16. Plugging this into the expression for  $\alpha$  and  $\beta$  gives

$$\begin{aligned}
|\alpha_{+,+}^{(t)}| &= \left| \frac{\eta}{NM} \sum_{n \in S_+} -\tilde{\ell}'_n(\theta) \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \right. \\
&\quad \cdot \left( V_+^{(t)} \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right)^2 \right) \right. \\
&\quad \left. - \sum_{i=2}^M \left( V_{n,i}^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \left. \left. \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right) \right) \\
&\quad + \sum_{i=2}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle \cdot \left( V_+^{(t)} \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right)^2 \right) \right. \\
&\quad \left. - \sum_{k=2}^M \left( V_{n,i}^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \left. \left. \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right) \right) \Big| \\
&\leq \frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{NM} \cdot \frac{3NM}{4} \cdot O(d_h^{-\frac{1}{4}}) + \frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau + 2\sigma_p\tau\sqrt{2\log(4NM/\delta)} + \tau^2)}{NM} \cdot \frac{3NM(M-1)}{4} \cdot O(d_h^{-\frac{1}{4}}) \\
&= O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{d_h^{\frac{1}{4}}}\right)
\end{aligned} \tag{117}$$

where the inequality is by  $-\tilde{\ell}_n^{(t)} \leq 1$  and the property that attention is smaller than 1 (e.g.  $\frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \leq 1$ ). We also have

$$\begin{aligned}
|\alpha_{n,+}^{(t)}| &= \left| -\frac{\eta}{NM} \ell_n^{(t)} \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \right. \\
&\quad \cdot \left( -V_+^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \\
&\quad \quad \left. + \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \\
&\quad \quad \left. + V_{n,i}^{(t)} \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \quad \quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right)^2 \right) \right. \\
&\quad \quad \left. - \sum_{k \neq i} \left( V_{n,k}^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \quad \quad \left. \left. \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right) \right) \\
&\quad + \sum_{k=2}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,k}^{(t)} \rangle \cdot \left( -V_+^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \\
&\quad \quad \left. \cdot \frac{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \\
&\quad \quad \left. + V_{n,i}^{(t)} \left( \frac{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \quad \quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right)^2 \right) \right. \\
&\quad \quad \left. - \sum_{l \neq i} \left( V_{n,l}^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \quad \quad \left. \left. \cdot \frac{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,l}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,k}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right) \right) \Bigg| \\
&\leq \frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{NM} \cdot M \cdot O(d_h^{-\frac{1}{4}}) \\
&= O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{d_h^{\frac{1}{4}} N}\right),
\end{aligned}$$

where the inequality is by  $-\tilde{\ell}_n^{(t)} \leq 1$ , Lemma 13 and the property that attention is smaller than 1. Similarly, we have

$$|\alpha_{-,-}^{(t)}|, |\beta_{+,+}^{(t)}|, |\beta_{-,-}^{(t)}| = O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{d_h^{\frac{1}{4}}}\right),$$

3456  
3457  
3458  
3459  
3460  
3461  
3462  
3463  
3464  
3465  
3466  
3467  
3468  
3469  
3470  
3471  
3472  
3473  
3474  
3475  
3476  
3477  
3478  
3479  
3480  
3481  
3482  
3483  
3484  
3485  
3486  
3487  
3488  
3489  
3490  
3491  
3492  
3493  
3494  
3495  
3496  
3497  
3498  
3499  
3500  
3501  
3502  
3503  
3504  
3505  
3506  
3507  
3508  
3509

$$|\alpha_{n,-,l}^{(t)}|, |\beta_{n,+,l}^{(t)}|, |\beta_{n,-,l}^{(t)}| = O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{d_h^{\frac{1}{4}}N}\right),$$

$$|\alpha_{n,l,-}^{(t)}|, |\beta_{n,l,+}^{(t)}|, |\beta_{n,l,-}^{(t)}|, |\alpha_{n,l,n',l'}^{(t)}|, |\beta_{n,l,n',l'}^{(t)}| = O\left(\frac{\eta(\|\boldsymbol{\mu}\|\tau + \sigma_p^2 d)}{d_h^{\frac{1}{4}}N}\right)$$

for  $t \in [0, T_0]$ . Next we use induction to show that the following proposition  $\mathcal{A}(t)$  holds for  $t \in [0, T_0]$ .  $\mathcal{A}(t)$ :

$$\begin{aligned} & |\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle|, |\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| \\ & = O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}\right), \end{aligned}$$

$$\begin{aligned} & |\langle \mathbf{q}_{\pm}^{(t)}, \mathbf{q}_{\mp}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{\pm}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',j}^{(t)} \rangle| \\ & = O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}\right), \end{aligned}$$

$$\begin{aligned} & |\langle \mathbf{k}_{\pm}^{(t)}, \mathbf{k}_{\mp}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{\pm}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| \\ & = O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}\right), \end{aligned}$$

$$\|\mathbf{q}_{\pm}^{(t)}\|_2^2, \|\mathbf{k}_{\pm}^{(t)}\|_2^2 = \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h)$$

$$\|\mathbf{q}_{n,i}^{(t)}\|_2^2, \|\mathbf{k}_{n,i}^{(t)}\|_2^2 = \Theta(\sigma_p^2 \sigma_h^2 d d_h)$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N]$ .

By Lemma 12 we know that  $\mathcal{A}(0)$  is true. Now we assume  $\mathcal{A}(0), \dots, \mathcal{A}(T)$  is true, then we need to prove that  $\mathcal{A}(T+1)$  is true. We first proof  $|\langle \mathbf{q}_+^{(T+1)}, \mathbf{k}_+^{(T+1)} \rangle| = O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}\right)$ , as an example.

$$\begin{aligned} |\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle| & = \left| \alpha_{+,+}^{(t)} \|\mathbf{k}_+^{(t)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+,i}^{(t)} \langle \mathbf{k}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle \right. \\ & \quad + \beta_{+,+}^{(t)} \|\mathbf{q}_+^{(t)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+,i}^{(t)} \langle \mathbf{q}_+^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle \\ & \quad + \left( \alpha_{+,+}^{(t)} \mathbf{k}_+^{(t)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+,i}^{(t)} \mathbf{k}_{n,i}^{(t)} \right) \\ & \quad \cdot \left( \beta_{+,+}^{(t)} \mathbf{q}_+^{(t)\top} + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+,i}^{(t)} \mathbf{q}_{n,i}^{(t)\top} \right) \Big| \\ & \leq O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{d_h^{\frac{1}{4}}}\right) \cdot \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h) \\ & \quad + NM \cdot O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{d_h^{\frac{1}{4}}N}\right) \cdot O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}\right) \\ & \quad + \{\text{lower order term}\} = O\left(\eta \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h^{\frac{3}{4}}\right) \end{aligned}$$

Taking a summation, we obtain that

$$\begin{aligned}
|\langle \mathbf{q}_+^{(T+1)}, \mathbf{k}_+^{(T+1)} \rangle| &\leq |\langle \mathbf{q}_+^{(0)}, \mathbf{k}_+^{(0)} \rangle| + \sum_{t=0}^T |\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle| \\
&\leq |\langle \mathbf{q}_+^{(0)}, \mathbf{k}_+^{(0)} \rangle| + \sum_{t=0}^{T_0-1} |\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle| \\
&\leq O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}\right) \\
&\quad + O\left(\frac{1}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right) \cdot O(\eta \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h^{\frac{3}{4}}) \\
&= O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}\right) + O\left(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h^{\frac{1}{2}}\right)
\end{aligned}$$

Similarly to  $\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle$ , it is easy to know that the inner product does not change by a magnitude more than the product of  $\max\{\alpha, \beta\}$  and  $\max\{\langle q, q \rangle, \langle k, k \rangle\}$  in a single iteration, which can be expressed as follows

$$\begin{aligned}
|\langle \mathbf{q}_+^{(t+1)}, \mathbf{k}_+^{(t+1)} \rangle - \langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle| \\
&= O\left(\max\left\{\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{d_h^{\frac{1}{4}}}, \frac{\eta(\sigma_p^2 d + \|\boldsymbol{\mu}\| \tau)}{d_h^{\frac{1}{4}} N}\right\}\right) \cdot \Theta(\max\{\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h, \sigma_p^2 \sigma_h^2 d d_h\}) \\
&= O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{d_h^{\frac{1}{4}}}\right) \cdot \Theta(\max\{\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h, \sigma_p^2 \sigma_h^2 d d_h\}) \\
&= O\left((\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h^{\frac{3}{4}} \cdot \max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\}\right)
\end{aligned}$$

where the second equality is by the condition that  $N \cdot \text{SNR}^2 = \Omega(1)$ .

Taking a summation, we obtain that

$$\begin{aligned}
|\langle \mathbf{q}^{(T+1)}, \mathbf{k}^{(T+1)} \rangle - \langle \mathbf{q}^{(0)}, \mathbf{k}^{(0)} \rangle| &\leq \sum_{t=0}^{T-1} |\langle \mathbf{q}^{(t+1)}, \mathbf{k}^{(t+1)} \rangle - \langle \mathbf{q}^{(t)}, \mathbf{k}^{(t)} \rangle| \\
&\leq \sum_{t=0}^{T_0-1} O\left(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h^{\frac{3}{4}} \cdot \max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\}\right) \\
&= O\left(\frac{1}{\eta d_h^{\frac{1}{4}} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2}\right) \cdot O\left(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h^{\frac{3}{4}} \cdot \max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\}\right) \\
&= O\left(\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 d_h^{\frac{1}{4}}\right).
\end{aligned}$$

It is clear that the magnitude of  $\langle \mathbf{q}^{(T+1)}, \mathbf{k}^{(T+1)} \rangle - \langle \mathbf{q}^{(0)}, \mathbf{k}^{(0)} \rangle$  is smaller than  $\max\{\|\boldsymbol{\mu}\|_2^2, \sigma_p^2 d\} \cdot \sigma_h^2 \cdot \sqrt{d_h \log(6N^2 M^2 / \delta)}$ . Thus the magnitude of the bound for  $\langle \mathbf{q}^{(T+1)}, \mathbf{k}^{(T+1)} \rangle$  is the same as that of  $\langle \mathbf{q}^{(T)}, \mathbf{k}^{(T)} \rangle$ . The proof for  $\langle \mathbf{q}^{(T+1)}, \mathbf{q}^{(T+1)} \rangle$  and  $\langle \mathbf{k}^{(T+1)}, \mathbf{k}^{(T+1)} \rangle$  is exactly the same, and we can conclude the proof by an induction.

## I.5 LOWER BOUNDS OF $\alpha$ AND $\beta$

In this subsection, we present some bounds for  $\alpha$  and  $\beta$  which can be used in F.2 and F.3. All the calculations in this subsection are based on the precise expression for  $\alpha$  and  $\beta$  in I.1 and assume that  $\mathcal{B}(T_1), \dots, \mathcal{B}(s), \mathcal{D}(T_1), \dots, \mathcal{D}(s-1)$  hold ( $s \in [T_1, t]$ ). Then the following propositions hold:

$$V_+^{(s)} \geq 3M \cdot |V_{n,i}^{(s)}|,$$

$$\begin{aligned}
3564 \quad & V_-^{(s)} \leq -3M \cdot |V_{n,i}^{(s)}|, \\
3565 \quad & \text{softmax}(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_\pm^{(s)} \rangle), \text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_\pm^{(s)} \rangle) \geq \frac{1}{M} - o(1), \\
3566 \quad & \text{softmax}(\langle \mathbf{q}_\pm^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle), \text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle) \leq \frac{1}{M} + o(1). \\
3567 \quad & \\
3568 \quad & \\
3569 \quad & \\
3570 \quad & \\
3571 \quad & 
\end{aligned}$$

Now we give the bounds respectively for  $\alpha_{+,+}^{(s)}, \alpha_{n,+}^{(s)}, \alpha_{-,-}^{(s)}, \alpha_{n,-}^{(s)}, \alpha_{n,i,+}^{(s)}, \alpha_{n,i,-}^{(s)}, \alpha_{n,i,n'}^{(s)}$ ,  
 $\beta_{+,+}^{(s)}, \beta_{n,+}^{(s)}, \beta_{-,-}^{(s)}, \beta_{n,-}^{(s)}, \beta_{n,i,+}^{(s)}, \beta_{n,i,-}^{(s)}, \beta_{n,i,n'}^{(s)}$ .

$$\begin{aligned}
3572 \quad & \alpha_{+,+}^{(s)} = \frac{\eta}{NM} \sum_{n \in S_+} -\tilde{\ell}'_n(\theta) \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(s)} \rangle \\
3573 \quad & \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
3574 \quad & \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right) \right) \\
3575 \quad & - \sum_{i=2}^M \left( V_{n,i}^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \\
3576 \quad & \left. \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \\
3577 \quad & + \sum_{i=2}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(s)} \rangle \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
3578 \quad & \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right) \right) \\
3579 \quad & - \sum_{k=2}^M \left( V_{n,i}^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \\
3580 \quad & \left. \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \\
3581 \quad & \geq \frac{\eta}{NM} \sum_{n \in S_+} -\tilde{\ell}'_n(\theta) \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(s)} \rangle \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
3582 \quad & \cdot \left( V_+^{(s)} \left( 1 - \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \right. \\
3583 \quad & \left. - \frac{1}{2} \cdot V_+^{(s)} \sum_{i=2}^M \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \\
3584 \quad & + \frac{\eta}{NM} \sum_{n \in S_+} -\tilde{\ell}'_n(\theta) \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(s)} \rangle \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
3585 \quad & \cdot \left( V_+^{(s)} \left( 1 - \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \right) \\
3586 \quad & \\
3587 \quad & \\
3588 \quad & \\
3589 \quad & \\
3590 \quad & \\
3591 \quad & \\
3592 \quad & \\
3593 \quad & \\
3594 \quad & \\
3595 \quad & \\
3596 \quad & \\
3597 \quad & \\
3598 \quad & \\
3599 \quad & \\
3600 \quad & \\
3601 \quad & \\
3602 \quad & \\
3603 \quad & \\
3604 \quad & \\
3605 \quad & \\
3606 \quad & \\
3607 \quad & \\
3608 \quad & \\
3609 \quad & \\
3610 \quad & \\
3611 \quad & \\
3612 \quad & \\
3613 \quad & \\
3614 \quad & \\
3615 \quad & \\
3616 \quad & \\
3617 \quad & 
\end{aligned}$$

$$\begin{aligned}
& -\frac{1}{2} \cdot V_+^{(s)} \sum_{k=2}^M \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
& \geq \frac{\eta}{2NM} \sum_{n \in S_+} -\tilde{\ell}'_n^{(s)} (\|\boldsymbol{\mu}\|_2 - \tau)^2 V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
& \cdot \left( 1 - \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \\
& + \sum_{i=2}^M (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
& \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right)
\end{aligned}$$

where the first inequality is by  $V_+^{(s)} \geq 3M \cdot |V_{n,i}^{(s)}|$ , the second inequality is by the fact that the sum of attention equal to 1 and Lemma 13. Similarly, we have

$$\begin{aligned}
\beta_{+,+}^{(s)} & \geq \frac{\eta}{2NM} \sum_{n \in S_+} -\tilde{\ell}'_n^{(s)} (\|\boldsymbol{\mu}\|_2 - \tau)^2 V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + (M-1) \exp(\max_j \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
& \cdot \frac{\exp(\max_j \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + (M-1) \exp(\max_j \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} + \{\text{lower order term}\} \quad (118)
\end{aligned}$$

$$\begin{aligned}
\alpha_{-,-}^{(s)} & \geq \frac{\eta}{2NM} \sum_{n \in S_-} \tilde{\ell}'_n^{(s)} (\|\boldsymbol{\mu}\|_2 - \tau)^2 V_-^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle)}{\exp(\langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle) + (M-1) \exp(\max_j \langle \mathbf{q}_-^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
& \cdot \frac{\exp(\max_j \langle \mathbf{q}_-^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle) + (M-1) \exp(\max_j \langle \mathbf{q}_-^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} + \{\text{lower order term}\} \quad (119)
\end{aligned}$$

Similarly, by applying the update rules in Section I.3, we can derive the following bounds on  $\alpha$  and  $\beta$ .

$$\begin{aligned}
\alpha_{+,+}^{(s)}, \alpha_{-,-}^{(s)}, \beta_{+,+}^{(s)}, \beta_{-,-}^{(s)}, \alpha_{n,i,+}^{(s)}, \alpha_{n,i,-}^{(s)}, \beta_{n,+i}^{(s)}, \beta_{n,-i}^{(s)} & \geq 0, \\
\alpha_{n,+i}^{(s)}, \alpha_{n,-i}^{(s)}, \alpha_{n,i,n,j}^{(s)}, \beta_{n,i,+}^{(s)}, \beta_{n,i,-}^{(s)}, \beta_{n,j,n,i}^{(s)} & \leq 0.
\end{aligned}$$

## I.6 LOWER BOUNDS OF $\langle \mathbf{q}, \mathbf{k} \rangle$

In order to give the lower bounds for  $\langle q, k \rangle$ , we need to rewrite the bounds of  $\alpha$  and  $\beta$  in a more concise form. We first expand the equations in I.5 under the assumption that  $\mathcal{B}(s)$  and  $\mathcal{E}(s)$  holds for  $s \in [T_1, t]$ .

$$\begin{aligned}
\alpha_{+,+}^{(s)} & \geq \frac{\eta}{2NM} \sum_{n \in S_+} -\tilde{\ell}'_n^{(s)} (\|\boldsymbol{\mu}\|_2 - \tau)^2 V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \\
& \cdot \left( 1 - \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \right) \\
& + \sum_{i=2}^M (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right.
\end{aligned}$$

$$\begin{aligned}
& - \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \\
& = \frac{\eta}{2NM} \sum_{n \in S_+} -\tilde{\ell}_n^{(s)} (\|\boldsymbol{\mu}\|_2 - \tau)^2 V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \\
& \quad \cdot \frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} + \{\text{lower term}\} \\
& \geq \frac{\eta}{2NM} -\tilde{\ell}_n^{(s)} (\|\boldsymbol{\mu}\|_2 - \tau)^2 V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \\
& \quad \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} + \{\text{lower term}\} \\
& \geq \frac{\eta}{2NM} -\tilde{\ell}_n^{(s)} (\|\boldsymbol{\mu}\|_2 - \tau)^2 V_+^{(s)} \cdot \left( \frac{1}{M} - o(1) \right) \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \\
& \geq \frac{\eta}{2NM} -\tilde{\ell}_n^{(s)} (\|\boldsymbol{\mu}\|_2 - \tau)^2 V_+^{(s)} \cdot \left( \frac{1}{M} - o(1) \right) \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{C \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} + \{\text{lower term}\} \\
& \geq \frac{\eta^2 C_5 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+}^{(s)})},
\end{aligned}$$

where the first inequality is by Lemma 20 (as  $\mathcal{E}(s)$  holds) and  $\text{softmax}(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) \geq (\frac{1}{M} - o(1))$ . In the fourth inequality, by  $\langle \mathbf{q}^{(T_1)}, \mathbf{k}^{(T_1)} \rangle = o(1)$  and the monotonicity of  $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle$  ( $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle$  is increasing and  $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$  is decreasing), there exist a constant  $C$  such that  $C \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) \geq \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)$ . In the last inequality, we plugging the lower bounds of  $V_+^{(s)}$  and  $-\tilde{\ell}_n^{(s)}$  and then absorb all the constant factors. Similarly, we have

$$\beta_{+,+}^{(s)} \geq \frac{\eta^2 C_5 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+}^{(s)})}, \quad (120)$$

$$\alpha_{-,-}^{(s)} \geq \frac{\eta^2 C_5 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,-}^{(s)})}, \quad (121)$$

$$\beta_{-,-}^{(s)} \geq \frac{\eta^2 C_5 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,-}^{(s)})}. \quad (122)$$

With the concise lower bounds for  $\alpha$  and  $\beta$  above and proposition  $\mathcal{C}(s)$ , we will give the lower bounds for the dynamics of  $\langle \mathbf{q}, \mathbf{k} \rangle$ .

$$\begin{aligned}
\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle & = \alpha_{+,+}^{(s)} \|\mathbf{k}_+^{(s)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \langle \mathbf{k}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle \\
& \quad + \beta_{+,+}^{(s)} \|\mathbf{q}_+^{(s)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle \\
& \quad + \left( \alpha_{+,+}^{(s)} \mathbf{k}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \mathbf{k}_{n,i}^{(s)} \right) \cdot \left( \beta_{+,+}^{(s)} \mathbf{q}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+}^{(s)} \mathbf{q}_{n,i}^{(s)} \right)
\end{aligned}$$

$$\begin{aligned}
&= \alpha_{+,+}^{(s)} \|\mathbf{k}_+^{(s)}\|_2^2 + \beta_{+,+}^{(s)} \|\mathbf{q}_+^{(s)}\|_2^2 + \{\text{lower order term}\} \\
&\geq \frac{2\eta^2 C_5 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+}^{(s)})} \cdot \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h) \\
&\quad + \{\text{lower order term}\} \\
&\geq \frac{\eta^2 C_6 (\|\boldsymbol{\mu}\|_2 - \tau)^4 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 \sigma_h^2 d_h (s - T_1)}{N} \cdot \frac{1}{\exp(\Lambda_{n,+}^{(s)})}, \tag{123}
\end{aligned}$$

Similarly, we have the lower bounds for the dynamics of other  $\langle \mathbf{q}, \mathbf{k} \rangle$ .

## I.7 UPPER BOUNDS OF $\langle \mathbf{q}, \mathbf{k} \rangle$

In order to give the upper bounds of  $\langle \mathbf{q}, \mathbf{k} \rangle$  in stage II, we need to give the upper bounds of  $\alpha$  and  $\beta$  based on the equations in Section I.1 under the assumption that  $\mathcal{D}(T_1), \dots, \mathcal{D}(s-1)$  hold for  $s \in [T_1, t]$ .

3780  
 3781  
 3782  
 3783  
 3784  
 3785  
 3786  
 3787  
 3788  
 3789  
 3790  
 3791  
 3792  
 3793  
 3794  
 3795  
 3796  
 3797  
 3798  
 3799  
 3800  
 3801  
 3802  
 3803  
 3804  
 3805  
 3806  
 3807  
 3808  
 3809  
 3810  
 3811  
 3812  
 3813  
 3814  
 3815  
 3816  
 3817  
 3818  
 3819  
 3820  
 3821  
 3822  
 3823  
 3824  
 3825  
 3826  
 3827  
 3828  
 3829  
 3830  
 3831  
 3832  
 3833

$$\begin{aligned}
\alpha_{+,+}^{(s)} &\leq \frac{\eta}{NM} \sum_{n \in S_+} -\tilde{\ell}'_n(\theta) \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(t)} \rangle \\
&\cdot \left( V_+^{(t)} \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right)^2 \right) \right) \\
&\quad - \sum_{i=2}^M \left( V_{n,i}^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \\
&\quad \left. \cdot \frac{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right) \\
&+ \sum_{i=2}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(t)} \rangle \cdot \left( V_+^{(t)} \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \right. \\
&\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right)^2 \right) \right) \\
&\quad - \sum_{k=2}^M \left( V_{n,i}^{(t)} \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right. \\
&\quad \left. \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,k}^{(t)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_+^{(t)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{k}_{n,j}^{(t)} \rangle)} \right) \\
&\leq \frac{\eta}{NM} \sum_{n \in S_+} (\|\boldsymbol{\mu}\|_2 + \tau)^2 (V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
&\quad + \max_{i=2} |V_{n,i}^{(s)}| \cdot \frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}) \\
&\quad + \sum_{i=2}^M (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) (V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
&\quad + \max_{i=2} |V_{n,i}^{(s)}| \cdot \frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}) \\
&\leq \frac{\eta}{NM} \cdot \frac{3N}{4} \cdot (\|\boldsymbol{\mu}\|_2 + \tau)^2 \cdot \left( V_+^{(s)} \cdot \frac{C}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} + \max_i |V_{n,i}^{(s)}| \cdot \frac{C}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \right) \\
&\quad + \sum_{i=2}^M (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \cdot \left( V_+^{(s)} \cdot \frac{C}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)} + \max_i |V_{n,i}^{(s)}| \cdot \frac{C}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \right) \\
&\leq \eta^2 C_9 (\|\boldsymbol{\mu}\|_2 + \tau)^2 + \sum_{i=2}^M \frac{\eta^2 C_9 (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} + \sum_{i=2}^M \frac{\eta^2 C_9 (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}
\end{aligned} \tag{124}$$

where the first inequality is by  $-\tilde{\ell}'_n(\theta) \leq 1$  and  $\text{softmax}(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) \leq 1$ . For the second inequality, we first consider  $\frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \leq \frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}$ .

Then by the monotonicity of  $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$  and  $\langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_{n,j}^{(T_1)} \rangle = o(1)$  we have  $\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle) \leq C$  for  $s \in [T_1, t]$ . The last inequality is by  $V_+^{(s)}, V_{n,i}^{(s)} = o(1)$  for  $s \in [T_1, t]$  and absorbing the constant factors. Similarly, we have

Similar to Section I.6, we apply the bounds of  $\alpha$  and  $\beta$  above to give the upper bounds for the dynamics  $\langle \mathbf{q}, \mathbf{k} \rangle$ .

$$\begin{aligned}
\langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle &= \alpha_{+,+}^{(s)} \|\mathbf{k}_+^{(s)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \langle \mathbf{k}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle \\
&+ \beta_{+,+}^{(s)} \|\mathbf{q}_+^{(s)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{q}_{n,i}^{(s)} \rangle \\
&+ \left( \alpha_{+,+}^{(s)} \mathbf{k}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \mathbf{k}_{n,i}^{(s)} \right)^\top \\
&\cdot \left( \beta_{+,+}^{(s)} \mathbf{q}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n,+}^{(s)} \mathbf{q}_{n,i}^{(s)} \right) \\
&+ \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \|\mathbf{k}_{n,i}^{(s)}\|_2^2 \\
&= \alpha_{+,+}^{(s)} \|\mathbf{k}_+^{(s)}\|_2^2 + \beta_{+,+}^{(s)} \|\mathbf{q}_+^{(s)}\|_2^2 + \{\text{lower order terms}\} \\
&\leq \left( \frac{\eta^2 C_9 (\|\boldsymbol{\mu}\|_2 + \tau)^2}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} + \sum_{i=2}^M \frac{\eta^2 C_9 (\|\boldsymbol{\mu}\|_{2\tau} + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \right) \cdot \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h) + \{\text{lower order terms}\} \\
&\leq \frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \sigma_h^2 d_h}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} + \sum_{i=2}^M \frac{\eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_{2\tau} + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \sigma_h^2 d_h}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}.
\end{aligned} \tag{125}$$

Similarly, we have the upper bounds for the dynamics of other  $\langle \mathbf{q}, \mathbf{k} \rangle$ .

## I.8 BOUNDS FOR THE SUM OF $\alpha$ AND $\beta$

The gradients of the inner products of  $q$  and  $k$  contain a lot of coefficients  $\alpha$  and  $\beta$ , and in order to conveniently give the upper bounds of some lower order inner products, we will give upper bounds for the summation of  $\alpha$  and  $\beta$  (e.g.  $\sum_{s=T_1}^t |\alpha_{+,+}^{(s)}|$ ).

Note that in the Jacobi matrix of the Softmax function, the elements on the diagonal are  $\text{softmax}(a_i) \cdot (1 - \text{softmax}(a_i))$  and the elements on the off-diagonal are  $\text{softmax}(a_i) \cdot \text{softmax}(a_j)$ . In Stage II, the attentions on signals  $\boldsymbol{\mu}_\pm$  increase and the attentions on noises  $\boldsymbol{\xi}$  decrease, then we can consider the following cases:

- if  $a_i = \langle \mathbf{q}_+, \mathbf{k}_+ \rangle$  or  $a_i = \langle \mathbf{q}_i, \mathbf{k}_+ \rangle$ ,  $\text{softmax}(a_i)$  has a constant upper bound 1,  $(1 - \text{softmax}(a_i))$  decreases as  $\text{softmax}(a_i)$  increases. So the upper bound of  $\text{softmax}(a_i) \cdot (1 - \text{softmax}(a_i))$  decreases as  $\text{softmax}(a_i)$  increases.
- if  $a_i = \langle \mathbf{q}_+, \mathbf{k}_j \rangle$  or  $a_i = \langle \mathbf{q}_i, \mathbf{k}_j \rangle$ ,  $(1 - \text{softmax}(a_i))$  has a constant upper bound 1. So the upper bound of  $\text{softmax}(a_i) \cdot (1 - \text{softmax}(a_i))$  decreases as  $\text{softmax}(a_i)$  decreases.
- if  $a_j = \langle \mathbf{q}_+, \mathbf{k}_j \rangle$  or  $a_j = \langle \mathbf{q}_i, \mathbf{k}_j \rangle$ ,  $\text{softmax}(a_i)$  has a constant upper bound 1. So the upper bound of  $\text{softmax}(a_i) \cdot \text{softmax}(a_j)$  decreases as  $\text{softmax}(a_j)$  decreases.

Based on the above cases, we first study the bounds of the following terms

- $1 - \text{softmax}(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)$

- 3888 •  $1 - \text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)$   
 3889  
 3890 •  $\text{softmax}(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)$   
 3891  
 3892 •  $\text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)$

3893  
 3894 Note that  $1 - \text{softmax}(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) = \sum_j \text{softmax}(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)$  and  $1 - \text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) =$   
 3895  $\sum_j \text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)$ , we only need to give the upper bounds for  $\text{softmax}(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)$  and  
 3896  $\text{softmax}(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)$ .  
 3897

3898 Assume that the propositions  $\mathcal{B}(T_1), \dots, \mathcal{B}(s), \mathcal{D}(T_1), \dots, \mathcal{D}(s-1)$  hold ( $s \in [T_1, t]$ ), we have

$$3899 \quad |V_{\pm}^{(s)}|, |V_{n,i}^{(s)}| \leq O(d_h^{-\frac{1}{4}}) + \eta C_4 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 (s - T_1), \quad (126)$$

$$3900 \quad \Lambda_{n,\pm,j}^{(s)} \geq \log \left( \exp(\Lambda_{n,\pm,j}^{(T_1)}) + \frac{\eta^2 C_8 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1) \right), \quad (127)$$

$$3901 \quad \Lambda_{n,i,\pm,j}^{(s)} \geq \log \left( \exp(\Lambda_{n,i,\pm,j}^{(T_1)}) + \frac{\eta^2 C_8 (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1) \right), \quad (128)$$

3902 for  $i, j \in [M] \setminus \{1\}, n \in [N], s \in [T_1, t]$ .

3903 Then we have

$$3904 \quad \frac{\exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \quad (129)$$

$$3905 \quad \leq \frac{\exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{C \exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle)}$$

$$3906 \quad = \frac{1}{C \exp(\Lambda_{n,\pm,j}^{(s)})}$$

$$3907 \quad \leq \frac{1}{C \exp(\Lambda_{n,\pm,j}^{(T_1)}) + \frac{\eta^2 C_8 C (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1)}$$

$$3908 \quad \leq \frac{1}{C_{13} + \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1)}$$

3909 For the first inequality, by  $\langle \mathbf{q}^{(T_1)}, \mathbf{k}^{(T_1)} \rangle = o(1)$  and the monotonicity of  $\langle \mathbf{q}^{(s)}, \mathbf{k}^{(s)} \rangle$  ( $\langle \mathbf{q}^{(s)}, \mathbf{k}^{(s)} \rangle$  is  
 3910 increasing and  $\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$  is decreasing), there exists a constant  $C$  such that  $C \exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) \geq$   
 3911  $\exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_{\pm}^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)$ . The second inequality is by plugging (127). For the  
 3912 last inequality, by  $\Lambda_{n,\pm,j}^{(T_1)} = o(1)$ , there exist a constant  $C_{13}$  such that  $C_{13} \leq C \exp(\Lambda_{n,\pm,j}^{(T_1)})$  and  
 3913  $C_{13} \leq C_8 C$ . Similarly, we have

$$3914 \quad \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} \quad (130)$$

$$3915 \quad \leq \frac{1}{C \exp(\Lambda_{n,i,+j}^{(s)})}$$

$$3916 \quad \leq \frac{1}{C_{13} + \frac{\eta^2 C_{13} (\sigma_p^2 d + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{1/2}}{N (\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1)}$$

Plugging above equations into the expressions of  $\alpha, \beta$  we have

$$\begin{aligned}
|\alpha_{+,+}^{(s)}| &\leq \left| \frac{\eta}{NM} \sum_{n \in S_+} -\tilde{\ell}'_n(\theta) \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(s)} \rangle \right. \\
&\quad \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
&\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right) \right. \\
&\quad \left. - \sum_{i=2}^M \left( V_{n,i}^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
&\quad \left. \left. \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \right) \\
&\quad + \sum_{i=2}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(s)} \rangle \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
&\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right) \right. \\
&\quad \left. - \sum_{k=2}^M \left( V_{n,i}^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
&\quad \left. \left. \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \right) \Big| \\
&\leq \frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 \cdot 3N}{NM} \cdot \left( O(d_h^{-\frac{1}{4}}) + \eta C_4 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 (s - T_1) \right) \\
&\quad \cdot O \left( \frac{1}{C_{13} + \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1)} \right) \\
&= O \left( \frac{\eta \|\boldsymbol{\mu}\|_2^2 d_h^{-\frac{1}{4}}}{C_{13} + \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1)} \right) \\
&\quad + O \left( \frac{\eta^2 (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{C_{13} + \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1)} \right) \\
&= O \left( \eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 d_h^{-\frac{1}{4}} \right) + O \left( \frac{\eta^2 (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{C_{13} + \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2 / \delta))^2} \cdot (s - T_1)(s - T_1 - 1)} \right).
\end{aligned}$$

(131)

3996 where the third equality is by  $\frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2} \cdot (s - T_1)(s - T_1 - 1) \geq 0$  for  $s \in [T_1, t]$ .

3997 Next, we give an upper bound for  $\frac{\eta^2 (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{C_{13} + \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2} \cdot (s - T_1)(s - T_1 - 1)}$  as follows:

3998

3999

4000

4001 
$$\frac{\eta^2 (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2 (s - T_1)}{C_{13} + \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2} \cdot (s - T_1)(s - T_1 - 1)}$$

4002

4003

4004

4005

4006

4007

4008

4009

4010

4011

4012

4013

4014

4015

4016

4017

4018

4019

4020

4021

4022

4023

4024

4025

4026

4027

4028

4029

4030

4031

4032

4033

4034

4035

4036

4037

4038

4039

4040

4041

4042

4043

4044

4045

4046

4047

4048

4049

$$\begin{aligned} &= \frac{\eta^2 (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2}{\frac{C_{13}}{(s - T_1)} + \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2} \cdot (s - T_1) - \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2}} \\ &\leq \frac{\eta^2 (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2}{2\sqrt{\frac{\eta^2 C_{13}^2 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2} - \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2}}} \\ &= \frac{\eta^2 (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2}{\frac{2\eta C_{13} (\|\boldsymbol{\mu}\|_2 - \tau) \|\boldsymbol{\mu}\|_2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{4}}}{N^{\frac{1}{2}} \log(6N^2 M^2/\delta)} - \frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 - \tau)^2 \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2}} \\ &= \frac{\eta^2 (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2}{\Theta\left(\frac{\eta (\|\boldsymbol{\mu}\|_2 - \tau) \|\boldsymbol{\mu}\|_2 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{4}}}{N^{\frac{1}{2}} \log(6N^2 M^2/\delta)}\right)} \\ &= O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right), \end{aligned} \tag{132}$$

where the inequality is by  $ax + \frac{b}{x} \geq 2\sqrt{ab}$  for  $x > 0$ , the third equality is by absorbing the lower order term  $\frac{\eta^2 C_{13} (\|\boldsymbol{\mu}\|_2 + \tau)^4 \|\mathbf{w}_O\|_2^2 d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2/\delta))^2}$ , the last equality is by  $\|\mathbf{w}_O\|_2 = \Theta(1)$ . Plugging this into (131) and get

4025

4026

4027

4028

$$\begin{aligned} |\alpha_{+,+}^{(s)}| &= O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 d_h^{-\frac{1}{4}}\right) + O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right) \\ &= O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right). \end{aligned} \tag{133}$$

Similarly, we have

4030

4031

4032

4033

4034

4035

4036

4037

4038

$$\begin{aligned} |\alpha_{-,-}^{(s)}|, |\beta_{+,+}^{(s)}|, |\beta_{-,-}^{(s)}| &= O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right), \\ |\alpha_{n,+}^{(s)}|, |\alpha_{n,-}^{(s)}| &= O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 N^{-\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right), \\ |\beta_{n,+}^{(s)}|, |\beta_{n,-}^{(s)}| &= O\left(\frac{\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\boldsymbol{\mu}\|_2 \log(6N^2 M^2/\delta)}{\sigma_p d^{\frac{1}{2}} N^{\frac{1}{2}} d_h^{\frac{1}{4}}}\right) = \\ &= O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 \cdot \text{SNR} \cdot N^{-\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right), \end{aligned} \tag{134}$$

for  $i \in [M] \setminus \{1\}, n \in S_+$ .

4041

4042

4043

4044

4045

4046

$$|\alpha_{n,i,+}^{(s)}|, |\alpha_{n,i,-}^{(s)}| = O\left(\frac{\eta \|\boldsymbol{\mu}\|_2 \sigma_p d^{\frac{1}{2}} \log(6N^2 M^2/\delta)}{N^{\frac{1}{2}} d_h^{\frac{1}{4}}}\right) = O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right), \tag{135}$$

for  $i \in [M] \setminus \{1\}, n \in S_+$ , the last equality is by  $N \cdot \text{SNR}^2 \geq \Omega(1)$ .

4047

4048

4049

$$|\beta_{n,i,+}^{(s)}|, |\beta_{n,i,-}^{(s)}| = O\left(\frac{\eta \sigma_p^2 d \log(6N^2 M^2/\delta)}{N^{\frac{1}{2}} d_h^{\frac{1}{4}}}\right) = O\left(\eta (\|\boldsymbol{\mu}\|_2 + \tau)^2 N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right), \tag{136}$$

4050 for  $i \in [M] \setminus \{1\}, n \in S_+$ , the last equality is by  $N \cdot \text{SNR}^2 \geq \Omega(1)$ .

$$4051 \quad |\alpha_{n,i,n,j}^{(s)}|, |\beta_{n,j,n,i}^{(s)}| = O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)\sigma_p d^{\frac{1}{2}} \log(6N^2 M^2/\delta)}{N^{\frac{1}{2}} d_h^{\frac{1}{4}}}\right) \quad (137)$$

$$4052 \quad = O\left(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right),$$

4053 for  $i, j \in [M] \setminus \{1\}, n \in [N]$ , the last equality is by  $N \cdot \text{SNR}^2 \geq \Omega(1)$ .

$$4054 \quad |\alpha_{n,i,n',j}^{(s)}|, |\beta_{n,j,n',i}^{(s)}| = O\left(\frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)\sigma_p \log(6N^2 M^2/\delta) \log(6N^2 M^2/\delta)}{N^{\frac{1}{2}} d_h^{\frac{1}{4}}}\right) \quad (138)$$

$$4055 \quad = O\left(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 d^{-\frac{1}{2}} d_h^{-\frac{1}{4}} (\log(6N^2 M^2/\delta))^2\right),$$

4056 for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], n \neq n'$ , the last equality is by  $N \cdot \text{SNR}^2 \geq \Omega(1)$ . Taking a summation we obtain that

$$4057 \quad \sum_{s=T_1}^t |\alpha_{+,+}^{(s)}| = O\left(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 \|\mathbf{w}_O\|_2^2 \log(6N^2 M^2/\delta)\right) \cdot O\left(\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2 N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right)$$

$$4058 \quad = O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}}\right), \quad (139)$$

4059 where the last equality is by  $\|\mathbf{w}_O\| = \Theta(1)$ . Similarly, we have

$$4060 \quad \sum_{s=T_1}^t |\alpha_{-,-}^{(s)}|, \sum_{s=T_1}^t |\beta_{+,+}^{(s)}|, \sum_{s=T_1}^t |\beta_{-,-}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,i,+}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,i,-}^{(s)}| = O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}}\right), \quad (140)$$

4061 for  $i \in [M] \setminus \{1\}, n \in S_+$ .

$$4062 \quad \sum_{s=T_1}^t |\alpha_{n,+,i}^{(s)}|, \sum_{s=T_1}^t |\alpha_{n,-,i}^{(s)}| = O\left(N^{-\frac{1}{2}} d_h^{-\frac{1}{4}}\right), \quad (141)$$

4063 for  $i \in [M] \setminus \{1\}, n \in S_+$ .

$$4064 \quad \sum_{s=T_1}^t |\beta_{n,+,i}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,-,i}^{(s)}| = O\left(\text{SNR} \cdot N^{-\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \quad (142)$$

4065 for  $i \in [M] \setminus \{1\}, n \in S_+$ .

$$4066 \quad \sum_{s=T_1}^t |\alpha_{n,i,+}^{(s)}|, \sum_{s=T_1}^t |\alpha_{n,i,-}^{(s)}|, \sum_{s=T_1}^t |\alpha_{n,i,n,j}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,j,n,i}^{(s)}| = O\left(d_h^{-\frac{1}{4}}\right) \quad (143)$$

4067 for  $i, j \in [M] \setminus \{1\}, n \in S_+$ .

$$4068 \quad \sum_{s=T_1}^t |\alpha_{n,i,n',j}^{(s)}|, \sum_{s=T_1}^t |\beta_{n,j,n',i}^{(s)}| = O\left(d^{-\frac{1}{2}} d_h^{-\frac{1}{4}} \log(6N^2 M^2/\delta)\right) \quad (144)$$

4069 for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], n \neq n'$ .

4070 With these sums of  $\alpha$  and  $\beta$  above, we can easily prove Claim 3 and Claim 4.

### 4071 I.9 PROOF OF CLAIM 3

4072 In this subsection, we assume that  $\mathcal{E}(T_1), \dots, \mathcal{E}(t)$  hold, and then proof that  $\mathcal{C}(t+1)$  is true with the result of I.8.

$$4073 \quad \left\| \mathbf{q}_+^{(t+1)} \right\|_2^2 - \left\| \mathbf{q}_+^{(T_1)} \right\|_2^2 \leq \sum_{s=T_1}^t \left| \left\| \mathbf{q}_+^{(s+1)} \right\|_2^2 - \left\| \mathbf{q}_+^{(s)} \right\|_2^2 \right| \quad (145)$$

$$\begin{aligned}
& \leq \sum_{s=T_1}^t \left| 2\alpha_{+,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle + 2 \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle \right. \\
& \quad \left. + (\alpha_{+,+}^{(s)} \mathbf{k}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \mathbf{k}_{n,i}^{(s)}) \cdot (\alpha_{+,+}^{(s)} \mathbf{k}_+^{(s)\top} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \mathbf{k}_{n,i}^{(s)\top}) \right| \\
& \leq 2 \sum_{s=T_1}^t |\alpha_{+,+}^{(s)}| |\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle| + 2 \sum_{n \in S_+} \sum_{i=2}^M \sum_{s=T_1}^t |\alpha_{n,+}^{(s)}| |\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle| \\
& \quad + \{\text{lower order term}\} \\
& = O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \cdot O(\log(d_h^{\frac{1}{2}})) + N \cdot M \cdot O\left(N^{-\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \cdot O(\log(d_h^{\frac{1}{2}})) \\
& = O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(d_h^{\frac{1}{2}})\right)
\end{aligned} \tag{146}$$

where the first inequality is by triangle inequality. Since  $\sigma_h^2 \geq (\max\{\sigma_p^2 d, \|\boldsymbol{\mu}\|_2^2\})^{-1} \cdot d_h^{-\frac{1}{2}} (\log(6N^2 M^2 / \delta))^{-2}$  and  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$ , we have  $N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(d_h^{\frac{1}{2}}) = o(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h)$ , so  $\|\mathbf{q}_+^{(t+1)}\|_2^2 = \|\mathbf{q}_+^{(T_1)}\|_2^2 + o(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h) = \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h)$ . Similarly, we have

$$\begin{aligned}
& \left| \|\mathbf{q}_-^{(t+1)}\|_2^2 - \|\mathbf{q}_-^{(T_1)}\|_2^2 \right| = O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(d_h^{\frac{1}{2}})\right) = o(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h), \\
& \left| \|\mathbf{k}_+^{(t+1)}\|_2^2 - \|\mathbf{k}_+^{(T_1)}\|_2^2 \right| = O\left((1 + \text{SNR}) N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(d_h^{\frac{1}{2}})\right) = o(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h), \\
& \left| \|\mathbf{q}_{n,i}^{(t+1)}\|_2^2 - \|\mathbf{q}_{n,i}^{(T_1)}\|_2^2 \right| = O\left(d_h^{-\frac{1}{4}} \log(d_h^{\frac{1}{2}})\right) = o(\sigma_p^2 \sigma_h^2 d_h d_n), \\
& \left| \|\mathbf{k}_{n,i}^{(t+1)}\|_2^2 - \|\mathbf{k}_{n,i}^{(T_1)}\|_2^2 \right| = O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(d_h^{\frac{1}{2}})\right) = o(\sigma_p^2 \sigma_h^2 d_h d_n),
\end{aligned} \tag{147}$$

so we have

$$\begin{aligned}
& \|\mathbf{q}_\pm^{(t+1)}\|_2^2, \|\mathbf{k}_\pm^{(t+1)}\|_2^2 = \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h), \\
& \|\mathbf{q}_{n,i}^{(t+1)}\|_2^2, \|\mathbf{k}_{n,i}^{(t+1)}\|_2^2 = \Theta(\sigma_p^2 \sigma_h^2 d_h d_n)
\end{aligned} \tag{148}$$

for  $i \in [M] \setminus \{1\}, n \in [N]$ .

$$\begin{aligned}
& |\langle \mathbf{q}_+^{(t+1)}, \mathbf{q}_-^{(t+1)} \rangle| \leq |\langle \mathbf{q}_+^{(T_1)}, \mathbf{q}_-^{(T_1)} \rangle| + \sum_{s=T_1}^t |\langle \mathbf{q}_+^{(s+1)}, \mathbf{q}_-^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{q}_-^{(s)} \rangle| \\
& \leq |\langle \mathbf{q}_+^{(T_1)}, \mathbf{q}_-^{(T_1)} \rangle| \\
& \quad + \sum_{s=T_1}^t \left| \alpha_{+,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle \right. \\
& \quad \left. + \alpha_{-,-}^{(s)} \langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle + \sum_{n \in S_-} \sum_{i=2}^M \alpha_{n,-}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle \right. \\
& \quad \left. + \left( \alpha_{+,+}^{(s)} \mathbf{k}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+}^{(s)} \mathbf{k}_{n,i}^{(s)} \right) \right. \\
& \quad \left. \cdot \left( \alpha_{-,-}^{(s)} \mathbf{k}_-^{(s)} + \sum_{n \in S_-} \sum_{i=2}^M \alpha_{n,-}^{(s)} \mathbf{k}_{n,i}^{(s)} \right)^\top \right|
\end{aligned} \tag{149}$$

4158  
4159  
4160  
4161  
4162  
4163  
4164  
4165  
4166  
4167  
4168  
4169  
4170  
4171  
4172  
4173  
4174  
4175  
4176  
4177  
4178  
4179  
4180  
4181  
4182  
4183  
4184  
4185  
4186  
4187  
4188  
4189  
4190  
4191  
4192  
4193  
4194  
4195  
4196  
4197  
4198  
4199  
4200  
4201  
4202  
4203  
4204  
4205  
4206  
4207  
4208  
4209  
4210  
4211

$$\begin{aligned}
&\leq |\langle \mathbf{q}_+^{(T_1)}, \mathbf{q}_-^{(T_1)} \rangle| \\
&\quad + \sum_{s=T_1}^t |\alpha_{+,+}^{(s)}| |\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle| + \sum_{n \in S_+} \sum_{i=2}^M \sum_{s=T_1}^t |\alpha_{n,+}^{(s)}| |\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle| \\
&\quad + \sum_{s=T_1}^t |\alpha_{-,-}^{(s)}| |\langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle| + \sum_{n \in S_-} \sum_{i=2}^M \sum_{s=T_1}^t |\alpha_{n,-}^{(s)}| |\langle \mathbf{q}_-^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle| \\
&\quad + \{\text{lower order term}\} \tag{150} \\
&\leq |\langle \mathbf{q}_+^{(T_1)}, \mathbf{q}_-^{(T_1)} \rangle| \\
&\quad + O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \cdot o(1) + N \cdot M \cdot O\left(N^{-\frac{1}{2}} d_h^{-\frac{1}{4}}\right) \cdot \log(d_h^{\frac{1}{2}})\right) \\
&= |\langle \mathbf{q}_+^{(T_1)}, \mathbf{q}_-^{(T_1)} \rangle| + O\left(N^{\frac{1}{2}} d_h^{-\frac{1}{4}} \log(d_h^{\frac{1}{2}})\right) \\
&= o(1),
\end{aligned}$$

where the first inequality is triangle inequality, the last equality is by  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$ . Similarly, we can prove:

$$\begin{aligned}
\|\mathbf{q}_\pm^{(t)}\|_2^2, \|\mathbf{k}_\pm^{(t)}\|_2^2 &= \Theta(\|\mu\|_2^2 \sigma_h^2 d_h), \\
\|\mathbf{q}_{n,i}^{(t)}\|_2^2, \|\mathbf{k}_{n,i}^{(t)}\|_2^2 &= \Theta(\sigma_p^2 \sigma_h^2 d d_h), \\
|\langle \mathbf{q}_+^{(t)}, \mathbf{q}_-^{(t)} \rangle|, |\langle \mathbf{q}_\pm^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',j}^{(t)} \rangle| &= o(1), \\
|\langle \mathbf{k}_+^{(t)}, \mathbf{k}_-^{(t)} \rangle|, |\langle \mathbf{k}_\pm^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| &= o(1),
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}, n, n' \in [N], i \neq j$  or  $n \neq n'$ .

## I.10 UPPER BOUNDS OF $\langle q, k \rangle$

In order to give the upper bounds for  $\langle \mathbf{q}, \mathbf{k} \rangle$  in stage III, we need to give the upper bounds of  $\alpha$  and  $\beta$  based on the equations in I.1. The main difference between this subsection and I.7 is that the bounds of  $|V_\pm|, |V_{n,i}|$  is  $\log(O(\frac{1}{\epsilon}))$  in this subsection, while the bounds of  $|V_\pm|, |V_{n,i}|$  is  $\log(O(\frac{1}{\epsilon}))$  in

4212  
4213  
4214  
4215  
4216  
4217  
4218  
4219  
4220  
4221  
4222  
4223  
4224  
4225  
4226  
4227  
4228  
4229  
4230  
4231  
4232  
4233  
4234  
4235  
4236  
4237  
4238  
4239  
4240  
4241  
4242  
4243  
4244  
4245  
4246  
4247  
4248  
4249  
4250  
4251  
4252  
4253  
4254  
4255  
4256  
4257  
4258  
4259  
4260  
4261  
4262  
4263  
4264  
4265

I.7, resulting in different bounds for  $\alpha$  and  $\beta$ . Now we take  $\alpha_{+,+}^{(s)}$  as an example

$$\begin{aligned}
\alpha_{+,+}^{(s)} &\leq \frac{\eta}{NM} \sum_{n \in S_+} -\tilde{\ell}_n(\theta) \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(s)} \rangle \\
&\quad \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
&\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right) \right) \\
&\quad - \sum_{i=2}^M \left( V_{n,i}^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \\
&\quad \left. \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \\
&\quad + \sum_{i=2}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(s)} \rangle \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\
&\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right) \right) \\
&\quad - \sum_{k=2}^M \left( V_{n,i}^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \\
&\quad \left. \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \\
&\leq \frac{\eta}{NM} \sum_{n \in S_+} (\|\boldsymbol{\mu}\|_2 + \tau)^2 (V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
&\quad + \max_{i=2}^M |V_{n,i}^{(s)}| \cdot \frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}) \\
&\quad + \sum_{i=2}^M (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) (V_+^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \\
&\quad + \max_{i=2}^M |V_{n,i}^{(s)}| \cdot \frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}) \\
&\leq \frac{\eta}{NM} \cdot \frac{3N}{4} \cdot (\|\boldsymbol{\mu}\|_2 + \tau)^2 \cdot \left( V_+^{(s)} \cdot \frac{C}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} + \max_i |V_{n,i}^{(s)}| \cdot \frac{C}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \right) \\
&\quad + \sum_{i=2}^M (\|\boldsymbol{\mu}\|_2 \tau + \sigma_p \tau \sqrt{2 \log(4NM/\delta)} + \tau^2) \cdot \left( V_+^{(s)} \cdot \frac{C}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)} + \max_i |V_{n,i}^{(s)}| \cdot \frac{C}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \right) \\
&\leq \frac{\eta C_9 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \log(O(\frac{1}{\epsilon}))}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)},
\end{aligned} \tag{151}$$

where the first inequality is by  $-\tilde{\ell}_n^{(s)} \leq 1$  and  $\text{softmax}(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) \leq 1$ . For the second inequality, we first consider  $\frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \leq \frac{\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}$ , then by the monotonicity of  $\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle$  and  $\langle \mathbf{q}_+^{(T_1)}, \mathbf{k}_{n,j}^{(T_1)} \rangle = o(1)$  we have  $\sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle) \leq C$  for

4266  $t \in [T_1, T_3]$ . The last inequality is by  $V_+^{(s)}, |V_{n,i}^{(s)}| \leq 2 \log(O(\frac{1}{\epsilon}))$  for  $t \in [T_2, T_3]$  and absorbing  
 4267 the constant factors. Similar to I.7, we can give the bounds for the other  $\alpha$  and  $\beta$  as follows:  
 4268

$$4269 \alpha_{-, -}^{(s)} \leq \frac{\eta C_9 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \log(O(\frac{1}{\epsilon}))}{\exp(\langle \mathbf{q}_-^{(s)}, \mathbf{k}_-^{(s)} \rangle)},$$

$$4270$$

$$4271 \beta_{+, +}^{(s)} \leq \frac{\eta C_9 (\|\boldsymbol{\mu}\|_2 - \tau)^2 \log(O(\frac{1}{\epsilon}))}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)},$$

$$4272$$

$$4273$$

$$4274$$
(152)

4275 Similar to I.6, we apply the bounds of  $\alpha$  and  $\beta$  above to give the upper bounds for the dynamics  
 4276  $\langle q, k \rangle$ .  
 4277

$$4278 \langle \mathbf{q}_+^{(s+1)}, \mathbf{k}_+^{(s+1)} \rangle - \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle$$

$$4279 = \alpha_{+, +}^{(s)} \|\mathbf{k}_+^{(s)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n, +, i}^{(s)} \langle \mathbf{k}_+^{(s)}, \mathbf{k}_{n, i}^{(s)} \rangle$$

$$4280$$

$$4281 + \beta_{+, +}^{(s)} \|\mathbf{q}_+^{(s)}\|_2^2 + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n, +, i}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{q}_{n, i}^{(s)} \rangle$$

$$4282$$

$$4283 + \left( \alpha_{+, +}^{(s)} \mathbf{k}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n, +, i}^{(s)} \mathbf{k}_{n, i}^{(s)} \right)$$

$$4284$$

$$4285 \cdot \left( \beta_{+, +}^{(s)} \mathbf{q}_+^{(s)\top} + \sum_{n \in S_+} \sum_{i=2}^M \beta_{n, +, i}^{(s)} \mathbf{q}_{n, i}^{(s)\top} \right)$$

$$4286$$

$$4287$$

$$4288$$

$$4289$$

$$4290 = \alpha_{+, +}^{(s)} \|\mathbf{k}_+^{(s)}\|_2^2 + \beta_{+, +}^{(s)} \|\mathbf{q}_+^{(s)}\|_2^2 + \{\text{lower order term}\}$$

$$4291$$

$$4292 \leq \frac{2\eta C_9 (\|\boldsymbol{\mu}\|_2 + \tau)^2 \log(O(\frac{1}{\epsilon}))}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)} \cdot \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h) + \{\text{lower order term}\}$$

$$4293$$

$$4294 \leq \eta C_{10} \|\boldsymbol{\mu}\|_2^2 (\|\boldsymbol{\mu}\|_2 + \tau) \sigma_h^2 d_h \log\left(O\left(\frac{1}{\epsilon}\right)\right) \frac{1}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)},$$

$$4295$$

$$4296$$

$$4297$$

$$4298$$

$$4299$$
(153)

4300 Similarly, we have the upper bounds for the dynamics of other  $\langle q, k \rangle$ .  
 4301

#### 4302 I.11 BOUNDS FOR THE SUM OF $\alpha$ AND $\beta$

4303 Assume that the propositions  $\mathcal{F}(T_2), \dots, \mathcal{F}(s), \mathcal{H}(T_2), \dots, \mathcal{H}(s-1)$  hold ( $s \in [T_1, t]$ ), we have  
 4304

$$4305 |V_{\pm}^{(s)}| \leq 2 \log\left(O\left(\frac{1}{\epsilon}\right)\right),$$

$$4306$$
(154)

$$4307 |V_{n, i}^{(s)}| = O(1),$$

$$4308$$
(279)

$$4309 \Lambda_{n, \pm, j}^{(s)} \geq \Lambda_{n, \pm, j}^{(T_2)} \geq \log\left(\exp(\Lambda_{n, \pm, j}^{(T_1)}) + \Theta\left(\frac{d_h^{\frac{1}{2}}}{N(\log(6N^2 M^2 / \delta))^3}\right)\right),$$

$$4310$$
(155)

$$4311 \Lambda_{n, i, \pm, j}^{(s)} \geq \Lambda_{n, i, \pm, j}^{(T_2)} \geq \log\left(\exp(\Lambda_{n, i, \pm, j}^{(T_1)}) + \Theta\left(\frac{\sigma_p^2 d d d_h^{\frac{1}{2}}}{N \|\boldsymbol{\mu}\|_2^2 (\log(6N^2 M^2 / \delta))^3}\right)\right)$$

$$4312$$
(156)

4313 for  $i, j \in [M] \setminus \{1\}, n \in [N], s \in [T_2, t]$ . Similar to (57) and (58), we have  
 4314

$$4315 \frac{\exp(\langle q_{\pm}^{(s)}, \mathbf{k}_{n, j}^{(s)} \rangle)}{\exp(\langle q_{\pm}^{(s)}, \mathbf{k}_{\pm}^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle q_{\pm}^{(s)}, \mathbf{k}_{n, j'}^{(s)} \rangle)} = O\left(\frac{N(\log(6N^2 M^2 / \delta))^3}{d_h^{\frac{1}{2}}}\right)$$

$$4316$$
(157)

$$\frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j'=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j'}^{(s)} \rangle)} = O\left(\frac{N\|\boldsymbol{\mu}\|_2^2 (\log(6N^2M^2/\delta))^3}{\sigma_p^2 d d_h^{\frac{1}{2}}}\right) \quad (158)$$

Plugging above equations into the expressions of  $\alpha, \beta$  and letting  $O(\log(O(\frac{1}{\epsilon})))$  be the upper bound for  $|V_{\pm}^{(s)}|, |V_{n,i}^{(s)}|$  we have

$$\begin{aligned} \alpha_{+,+}^{(s)} &\leq \frac{\eta}{NM} \sum_{n \in S_+} -\tilde{\ell}'_n(\theta) \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\mu}}_+^{(s)} \rangle \\ &\cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\ &\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right) \right) \\ &- \sum_{i=2}^M \left( V_{n,i}^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \\ &\quad \left. \cdot \frac{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \\ &+ \sum_{i=2}^M \langle \tilde{\boldsymbol{\mu}}_+, \tilde{\boldsymbol{\xi}}_{n,i}^{(s)} \rangle \cdot \left( V_+^{(s)} \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \right. \\ &\quad \left. \left. - \left( \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right)^2 \right) \right) \\ &- \sum_{k=2}^M \left( V_{n,i}^{(s)} \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right. \\ &\quad \left. \cdot \frac{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,k}^{(s)} \rangle)}{\exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_+^{(s)} \rangle) + \sum_{j=2}^M \exp(\langle \mathbf{q}_{n,i}^{(s)}, \mathbf{k}_{n,j}^{(s)} \rangle)} \right) \\ &\leq \frac{\eta(\|\boldsymbol{\mu}\|_2 + \tau)^2}{NM} \cdot \frac{3N}{4} \cdot O\left(\log\left(O\left(\frac{1}{\epsilon}\right)\right)\right) \cdot O\left(\frac{N(\log(6N^2M^2/\delta))^3}{d_h^{\frac{3}{2}}}\right) \\ &= O\left(\frac{\eta N(\|\boldsymbol{\mu}\|_2 + \tau)^2 (\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{d_h^{\frac{3}{2}}}\right) \end{aligned} \quad (159)$$

Taking a summation we obtain that

$$\begin{aligned} \sum_{s=T_2}^t |\alpha_{+,+}^{(s)}| &= O\left(\frac{1}{\eta \epsilon \|\boldsymbol{\mu}\|_2^2 \|\mathbf{w}_O\|_2^2}\right) \cdot O\left(\eta \|\boldsymbol{\mu}\|_2^2 N (\log(6N^2M^2/\delta))^3 \log\left(O\left(\frac{1}{\epsilon}\right)\right) \frac{\log(O(\frac{1}{\epsilon}))}{d_h^{\frac{3}{2}}}\right) \\ &= O\left(\frac{N(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{3}{2}}}\right), \end{aligned} \quad (160)$$

where the last equality is by  $\|\mathbf{w}_O\| = \Theta(1)$ . Similarly, we have the bounds for other sum of  $\alpha$  and  $\beta$ . Thus, we can easily prove Claim 7 and Claim 8.

## I.12 PROOF OF CLAIM 7

In this subsection, we assume that  $\mathcal{I}(T_2), \dots, \mathcal{I}(t)$  hold, and then proof that  $\mathcal{G}(t+1)$  is true with the result of I.11.

$$\begin{aligned}
\|\mathbf{q}_+^{(t+1)}\|_2^2 - \|\mathbf{q}_+^{(t)}\|_2^2 &\leq \sum_{s=T_2}^t \left| \|\mathbf{q}_+^{(s+1)}\|_2^2 - \|\mathbf{q}_+^{(s)}\|_2^2 \right| \\
&\leq \sum_{s=T_2}^t \left| 2\alpha_{+,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle + 2 \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+,i}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle \right. \\
&\quad \left. + \langle \alpha_{+,+}^{(s)} \mathbf{k}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+,i}^{(s)} \mathbf{k}_{n,i}^{(s)} \rangle \cdot \left( \alpha_{+,+}^{(s)} \mathbf{k}_+^{(s)} + \sum_{n \in S_+} \sum_{i=2}^M \alpha_{n,+,i}^{(s)} \mathbf{k}_{n,i}^{(s)} \right)^\top \right| \\
&\leq 2 \sum_{s=T_2}^t \left| \alpha_{+,+}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_+^{(s)} \rangle \right| + 2 \sum_{n \in S_+} \sum_{i=2}^M \sum_{s=T_2}^t \left| \alpha_{n,+,i}^{(s)} \langle \mathbf{q}_+^{(s)}, \mathbf{k}_{n,i}^{(s)} \rangle \right| \\
&\quad + \{\text{lower order term}\} \\
&= O\left( \frac{N(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{3}{2}}} \right) \cdot \log(\epsilon^{-1} d_h^{\frac{1}{2}}) \\
&\quad + N \cdot M \cdot O\left( \frac{(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon}))}{\epsilon d_h^{\frac{3}{2}}} \right) \cdot \log(\epsilon^{-1} d_h^{\frac{1}{2}}) \\
&= O\left( \frac{N(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon})) \log(\epsilon^{-1} d_h^{\frac{1}{2}})}{\epsilon d_h^{\frac{3}{2}}} \right)
\end{aligned} \tag{161}$$

where the first inequality is by triangle inequality, the second inequality is by the update rules in D.2, the third inequality is by  $t \leq T_3$ . Since  $\sigma_h^2 \geq (\max\{\sigma_p^2 d, \|\boldsymbol{\mu}\|_2^2\})^{-1} \cdot d_h^{-\frac{1}{2}} (\log(6N^2M^2/\delta))^{-2}$  and  $d_h = \tilde{\Omega}(\max\{\text{SNR}^4, \text{SNR}^{-4}\} N^2 \epsilon^{-2})$ , we have

$$\frac{N(\log(6N^2M^2/\delta))^3 \log(O(\frac{1}{\epsilon})) \log(\epsilon^{-1} d_h^{\frac{1}{2}})}{\epsilon d_h^{\frac{3}{2}}} = o(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h), \tag{162}$$

so  $\|\mathbf{q}_+^{(t+1)}\|_2^2 = \|\mathbf{q}_+^{(T_2)}\|_2^2 + o(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h) = \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h)$ . Similarly, we have

$$\begin{aligned}
\|\mathbf{q}_\pm^{(t)}\|_2^2, \|\mathbf{k}_\pm^{(t)}\|_2^2 &= \Theta(\|\boldsymbol{\mu}\|_2^2 \sigma_h^2 d_h), \\
\|\mathbf{q}_{n,i}^{(t)}\|_2^2, \|\mathbf{k}_{n,i}^{(t)}\|_2^2 &= \Theta(\sigma_p^2 \sigma_h^2 d d_h), \\
|\langle \mathbf{q}_+^{(t)}, \mathbf{q}_-^{(t)} \rangle|, |\langle \mathbf{q}_\pm^{(t)}, \mathbf{q}_{n,i}^{(t)} \rangle|, |\langle \mathbf{q}_{n,i}^{(t)}, \mathbf{q}_{n',j}^{(t)} \rangle| &= o(1), \\
|\langle \mathbf{k}_+^{(t)}, \mathbf{k}_-^{(t)} \rangle|, |\langle \mathbf{k}_\pm^{(t)}, \mathbf{k}_{n,i}^{(t)} \rangle|, |\langle \mathbf{k}_{n,i}^{(t)}, \mathbf{k}_{n',j}^{(t)} \rangle| &= o(1)
\end{aligned}$$

for  $i, j \in [M] \setminus \{1\}$ ,  $n, n' \in [N]$ ,  $i \neq j$  or  $n \neq n'$ .