# The Mechanism of Prediction Head in Non-contrastive Selfsupervised Learning

Zixin Wen
zixinw@andrew.cmu.edu
Carnegie Mellon University

Yuanzhi Li yuanzhil@andrew.cmu.edu Carnegie Mellon University

May 13, 2022

#### **Abstract**

Recently the surprising discovery of *Bootstrap Your Own Latent* (BYOL) method by Grill et al. shows the negative term in contrastive loss can be removed if we add the so-called *prediction head* to the network architecture, which breaks the symmetry between the positive pairs. This initiated the research of *non-contrastive self-supervised learning*. It is mysterious why even when trivial collapsed *global optimal* solutions exist, neural networks trained by (stochastic) gradient descent can still learn competitive representations and avoid collapsed solutions. This phenomenon is one 14 Mayof the most typical examples of implicit bias in deep learning optimization, and its underlying mechanism remains little understood to this day.

In this work, we present our empirical and theoretical discoveries about the mechanism of prediction head in non-contrastive self-supervised learning methods. Empirically, we find that when **the prediction head is initialized as an identity matrix with only its off-diagonal entries being trained**, the network can learn competitive representations even though the trivial optima still exist in the training objective. Moreover, we observe a consistent rise and fall trajectory of off-diagonal entries during training. Our evidence suggests that understanding the identity-initialized prediction head is a good starting point for understanding the mechanism of the trainable prediction head.

Theoretically, we present a framework to understand the behavior of the trainable, but identity-initialized prediction head. Under a simple setting, we characterized the **substitution effect** and **acceleration effect** of the prediction head during the training process. The substitution effect happens when learning the stronger features in some neurons can substitute for learning these features in other neurons through updating the prediction head. And the acceleration effect happens when the substituted features can accelerate the learning of other weaker features to prevent them from being ignored. These two effects together enable the neural networks to learn all the features rather than focus only on learning the stronger features, which is likely the cause of the dimensional collapse phenomenon. To the best of our knowledge, this is also the first end-to-end optimization guarantee for non-contrastive methods using nonlinear neural networks with a trainable prediction head and normalization.

## **Contents**

1 Introduction	. 1
1.1 Comparison to Similar Studies	. 4
2 Preliminaries on Non-contrastive Learning	. 7
3 Problem Setup	. 9

3.1 Learner Network	11
3.2 Training Algorithm	12
4 Statements of Main Results	13
5 The Four Phases of the Learning Process	15
5.1 Phase I: Learning the Stronger Feature	15
5.2 Phase II: The Substitution Effect	16
5.3 Phase III: The Acceleration Effect	16
5.4 The End Phase: Convergence	17
6 Additional Related Work	18
7 Conclusion and Discussion	19
8 Experiment Details	19
A Notations and Gradients	20
A.1 Gradient Computation	21
A.2 Some Useful Bounds for Gradients	24
B Phase I: Learning the Stronger Feature	26
B.1 Induction in Phase I	27
B.2 Computing Variables at Phase I	27
B.3 Gradient Lemmas for Phase I	30
B.4 At the End of Phase I	39
C Phase II: The Substitution Effect of Prediction Head	44
C.1 Induction in Phase II	45
C.2 Gradient Lemmas for Phase II	45
C.3 At the End of Phase II	50
D Phase III: The Acceleration Effect of Prediction Head	57
D.1 Induction in Phase III	57
D.2 Gradient Lemmas for Phase III	58
D.3 At the End of Phase III	63
E The End Phase: Convergence	72
E.1 Proof of Convergence	77

F Learning Without Prediction Head	. 79
G Tensor Power Method Bounds	. 80

#### 1 Introduction

Self-supervised learning is about learning representations of real-world vision or language data without human supervision, and contrastive learning [66, 45, 43, 24, 20, 34] is one of the most successful self-supervised learning approaches. It has been known that the behavior of contrastive learning depends critically on the minimization of the negative term, which corresponds to contrasting the representations of negative pairs, i.e., pairs of different data points. However, the surprising finding of the Bootstrap Your Own Latent (BYOL) method by Grill et al. [39] initiated the research of non-contrastive self-supervised learning, which refers to contrastive learning methods without using the negative pairs. BYOL achieved state-of-the-art results in various computer vision benchmarks and there are plenty of follow-up works [41, 26, 21, 17, 33, 91, 46, 65] making improvements in this direction.

one wishes to learn a network  $\varphi$  such that  $\varphi(x)$  aligns in direction with  $\varphi(x^0)$ , where x and  $x^0$  are called the *positive pair*, generated by random augmentations from the same sample. Without

On a high level, in non-contrastive self-supervised learning,

contrasting the negative pairs, it is extremely easy for neural networks to cheat the learning task by learning certain inferior representations. One trivial solution known as the complete *collapse* is when  $\varphi(\cdot)$  is a constant vector whose variance is zero. Another trivial *global optimal* 

[46]. Nevertheless, adding a trainable prediction head on top of (one branch of)  $\varphi(x)$  magically avoids learning such solutions, solution, typically learned by the neural network after training, is when all the coordinates  $\varphi_i(\cdot)$  are

exactly aligned, which is named as dimensional collapse by Hua et al.

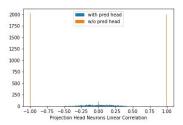
even though the prediction head can possibly learn the identity mapping and render itself useless. It is mysterious why even if the network can minimize the training objective by learning an identity prediction head and a collapsed encoder network  $\varphi(\cdot)$ , it still optimizes for a non-collapsed state-of-the-art representation instead when trained by (variants of) stochastic gradient descent (SGD).

Since the proposition of BYOL, there have been lots of empirical(b) Histograms of the correlations of encoder network neurons (before studies trying to understand non-contrastive learning. The projection head). SimSiam method by Chen and He [26] shows the exponential mov-

ing average (EMA) is not necessary for avoiding collapsed solutions while **stop-gradient** is necessary. Richemond et al. [72] empirically disproved the conjecture that information leakage from batch normalization (BN) is the reason why BYOL can avoid collapse. DINO [21] further explored replacing the normalized '2-loss by a cross-entropy loss. Zhang et al. [92] gives empirical evidence that using a single bias layer as a prediction head is capable of avoiding collapsed solutions. All the methods above

Figure 1: Dimensional Collapse.

Network trained without prediction head will learn extremely correlated neurons.



(a) Histograms of the correlations of projection head neurons.

3500 3000 2000 1500 use loss functions that are asymmetric with respect to the positive pair. If one wishes to work without both asymmetry and the negative pairs, one must add extra diversity-enforcing structures say neuron-wise regularization in *Barlow Twins* [91] or a more complicated output normalization scheme than BN [33, 46]. The seminal works [91, 46] provide empirical evidence that the prediction head encourages the network to learn more diversified features. But in theory, the question of how the prediction head helps in learning those diverse features is still unanswered.

Despite the great empirical effort put to investigate these non-contrastive learning methods, there is very little theoretical progress towards explaining them. Most of existing theories focus on contrastive learning, especially from the statistical learning perspective [83, 85, 14, 84, 42, 86, 13, 15, 50, 47, 63]. The theoretical tools used in these paper rely heavily on the properties of the minima of loss function. However, due to the existence of trivial dimensional collapsed *global optimal* solutions (even with the prediction head) of the non-contrastive methods, to the best of our knowledge, *there is no well-established statistical framework for those methods yet.* To explain the non-contrastive learning, it is inevitable to study how the solutions are chosen during the optimization. Therefore, we consider understanding the optimization process to be crucial for understanding these methods. Our research questions are:

#### Our theoretical questions: the role of prediction head

Why do most non-contrastive self-supervised methods learn collapsed solutions when the socalled prediction head is absent in the network architecture? How does the *trainable* prediction head help **optimizing** the neural network to learn more diversified representations in noncontrastive self-supervised learning?

Theoretical challenges of our questions. Due to the existence of trivial collapsed optimal solutions of the non-contrastive learning objective, we need to understand the **implicit bias in optimization** posed by the prediction head. However, to the best of our knowledge, all of the previous implicit biases theories focus only on the supervised learning tasks, and thus cannot be applied to our question. Even though [89] has characterized the training trajectory of contrastive learning, its analysis cannot incorporate the training of the prediction head. In theory, the optimization of nonlinear neural networks with at least two trainable layers in self-supervised learning is still intractable. A detailed explanation of our challenges will be given in Section 2.

There are already some theoretical papers [82, 87, 67] that try to address similar questions. While none of these papers studied the training process of the prediction head, our results provide a completely different perspective: **We explain why** *training* **the prediction head can encourage the network to learn diversified features and avoid dimensional collapses**, even

when the trivial collapsed optima still exist in the training objective, which is not covered by the

prior works. We defer the detailed comparison of similar works to Section 1.1. On a high level, the results in this paper are summarized as follows:

**Our empirical contributions.** In non-contrastive self-supervised learning, we obtain the following experimental results:

• We discover empirically that even when the prediction head is **linear** and initialized as an identity matrix with only off-diagonal entries being trainable, the performance of learned

representation is comparable to using the usual non-linear two-layer MLP or randomly initialized (trainable) linear prediction head. This disproves the belief that non-symmetric initialization of the online and target network is needed. See Figure 2.

• We empirically verified that even when the prediction head is an identity-initialized matrix, it does not always converge to a symmetric matrix during training. This proves the trainable prediction head does not need to behave like a symmetric matrix during most of the training process. Therefore the theories based on symmetric prediction head [82, 87] cannot fully explain the behaviors of the trainable prediction head. See Figure 3 and Figure 4.

**Our theoretical contributions.** We based our theory on a very simple setting, where the data consist of two features: the strong feature and the weak feature. Intuitively, we can think of the strong features in a dataset are the ones that show up more frequently or with large magnitude,

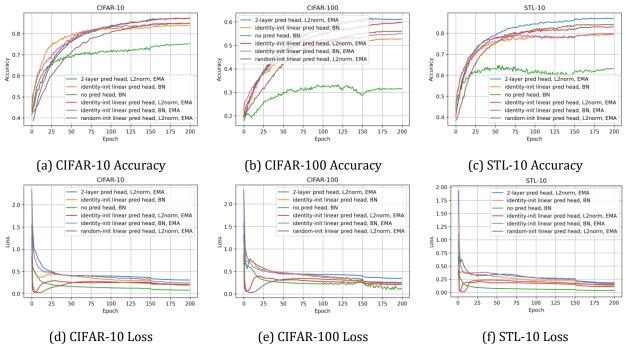


Figure 2: Performances of using different prediction heads. Here in CIFAR-10, CIFAR-100 and STL-10, identity-initialized linear prediction head can achieve good accuracies comparable to commonly used two-layer non-linear MLP or randomly-initialized linear head. All the prediction heads are trainable, while for identity initialized prediction head only the off-diagonal entries are trainable. Here BN or L2norm represents the output normalization, and EMA represents using exponential moving average to update the target network as in BYOL [41]. More details of these experiments can be seen in Section 8.

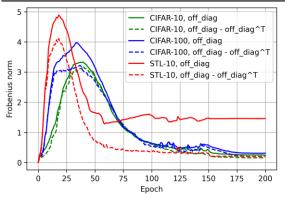
and weak features as those that show up rarely or with small magnitude. We consider learning with a **two-layer non-linear neural network with output normalization** using (stochastic) gradient descent. Under this setting, we obtain the following results.

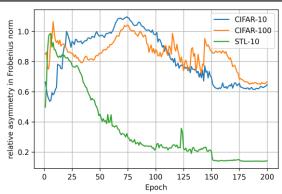
• We prove that without a prediction head, even with BN on the output to avoid complete collapse, the networks will still converge to dimensional collapsed solutions, which provides a theoretical explanation to the dimensional collapse phenomenon observed in [46]. • We prove that the trainable prediction head, combined with suitable output normalization and stop-gradient

operation, can learn diversified features to avoid the dimensional collapse problem. We characterize two effects of prediction head: the **substitution effect** and the **acceleration effect**. The intuitions of these two effects are summarized below:

### The mechanism of the trainable prediction head

In our setting, we prove that (1) without the prediction head, all the neurons will only learn the strongest feature in the data set thus causing dimensional collapses; (2) the trainable prediction head can help to learn weak features by leveraging two effects: the **substitution effect** and the **acceleration effect**. The substitution effect happens when by learning the prediction head, the learned stronger features in some neurons can substitute for learning the same features in other neurons, which decreases the learning speed of strong features in those neurons. And the acceleration effect happens when the strong features substituted via the prediction head can further accelerate the learning of weaker features in those substituted neurons.





(a) koff-diag( $E^{(t)}$ ) $k_F$  and  $kE^{(t)} - (E^{(t)}) k_F$ 

(b)  $kE^{(t)} - (E^{(t)}) k_F / koff-diag(E^{(t)}) k_F$ 

Figure 3: Trajectories of the identity-initialized prediction head. off-diag(E) is obtained by setting the diagonal of E to be zero. In (a), we discover that over all three datasets considered here, the Frobenius norm of our identity-initialized prediction head's off-diagonal matrix clearly display a two stage separation, more precisely, a rise and fall pattern; In (b), The off-diagonal matrix of the prediction head is not symmetric in CIFAR-10 and CIFAR-100. Since the diagonal entries are fixed to one, our measure is more accurate in measuring the symmetricity of the prediction head matrix.

Besides the above effects, we also explain, in our setting, how the two common components in non-contrastive learning: *stop-gradient* operation and *output normalization*, can assist the prediction head in creating those effects during the training process. We point out it is the <u>interactions</u> between these components, rather than their individual effects, that ensure the success of the training. We shall discuss this in more detail in Section 5.3.

#### 1.1 Comparison to Similar Studies

In this section, we will clarify the differences between our results and some similar studies. Especially the theoretical papers by Tian et al. [82] and Wang et al. [87]. Pokle et al. [67] compared the landscapes between contrastive and non-contrastive learning and points out the existence of non-collapsed bad minima for non-contrastive learning without a prediction head.

We point out that all the claims below are derived **only in our theoretical setting** and are partially verified in experiments over datasets such as CIFAR-10, CIFAR-100, and STL-10.

Can eigenspace alignment explain the effects of training the prediction head? The paper [82] presented a theoretical statement that (symmetric) linear prediction head will converge to a matrix that commutes with the covariance matrix of linear representations at the end of training, and they provided experiments to support their theory. However, our theory suggests that the intermediate stage of training the prediction head matters more to the feature learning of the base network than the convergence stage. Indeed, as shown in Figure 3, in many cases, the trainable projection head will **converge back to identity** after training, which commutes with any covariance matrix. However, simply setting the prediction head to identity without training leads to significantly worse results. Therefore, we believe that it is critical to study the entire learning process to understand the role of the prediction head. We prove that in our setting, the substitution effect and the acceleration effect happen during the stage when the networks are trying to learn the weaker features, and after that, the prediction head will converge back to the identity matrix at the end of training (see Proposition 5.4). Again, we emphasize that our characterization of the prediction head trajectory is partially verified by the experiments in Figure 3a: the training trajectory of the prediction head displays a clear two-stage separation, which demonstrates that the convergence result (e.g., the eigenspace alignment result in [82]) is not sufficient to characterize the training process of prediction head. We conjecture the result in [82] on the prediction head is due to a similar convergence result we obtain at the end of training.

Can the symmetric prediction head explain the trainable prediction head? In the paper [82], experiments over the STL-10 dataset showed that the linear prediction head tends to converge to a symmetric matrix during training. And the follow-up paper [87] established a theory under the symmetric prediction head (which is not trained but manually set at each iteration). However, similar to the reason why eigenspace alignment cannot fully explain the effects of the prediction head, the symmetric prediction head given in [87] might not explain the trainable prediction head as well. Under their linear network setting, where W is the weight matrix of the base encoder, they manually set the prediction head  $W_p$  at iteration t to be

$$W_p^{(t)} \leftarrow W^{(t)} \mathbb{E}_{x_1} x_1 x_1^{\top} (W^{(t)})^{\top}$$
(1.1)

and the outputs of both online and target network are not normalized. Under this manual update rule of the prediction head, they proved a subspace learning result under gaussian data setting.

Nevertheless, our experiments in Figure 2 and Figure 3b show that even if we initialize the prediction head using a symmetric matrix (identity), the trainable prediction head can be very asymmetric at the early training stage when the encoder network learn most of its features. Moreover, Figure 3b demonstrates that the prediction heads in CIFAR-10 and CIFAR-100 experiments do not converge to a symmetric matrix. In accord with these experiments, our theory suggests that the prediction head cannot converge to a symmetric matrix before the encoder network has successfully learned all the features. Moreover, the theory in [87] cannot distinguish between learning complete collapsed (zero) solutions and learning dimensional collapsed ones, therefore cannot explain why the prediction can help avoid the dimensional collapse. Actually, in the presence of feature imbalance (e.g.,  $\mathbb{E}_{x_1} x_1 x_1^{\top}$  has huge eigen-gap), the symmetric prediction head in (1.1) is also likely to collapse into a rank-one matrix where W focus on learning the largest eigenvector of the covariance  $\mathbb{E}_{x_1} x_1 x_1^{\top}$ .

The differences between our results and [87]'s are in that we are based on nonlinear network architecture and a trainable prediction head. Indeed, our theory and experiments in Figure 7 show

that when feature imbalance happens (which is very common in vision datasets, see [25]), training a nonlinear network would cause discrepancies in the learning pace between different neurons. We proved that by learning to become asymmetric, the trainable prediction head can leverage such discrepancies and creates the substitution effect (see Lemma 5.2) and the acceleration effect (see Theorem 5.3). We believe this proves that asymmetry is the key to explaining the implicit bias of the trainable prediction head and our results establish the symmetry-breaking mechanism of the prediction head in non-contrastive learning.

The role of stop-gradient and output-normalization. The seminal work [26] gave empirical results showing that stop-gradient operation is essential for avoiding the collapsed solutions. It is discussed in the theory of Tian et al. [82] that without the stop-gradient, the linear network will learn the zero (constant) solution. [87] also incorporated the stop-gradient into their theory, but they did not explain why stop-gradient is necessary for their setting. We provide a different perspective about why stop-gradient and output normalization (together) are necessary for noncontrastive learning. We proved in our setting, that the stop-gradient and output-normalization

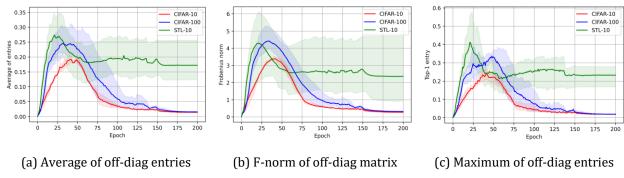


Figure 4: Trajectories of the identity-initialized prediction head with a (min,max) confidence band, average over 3 runs. In all three datasets, we observe a consistent rise and fall trajectory pattern.

together can turn the features substituted via the prediction head into a factor in the gradient of the slower learning neurons, thereby creating the acceleration effect. If either one of these components is missing, the acceleration effect of the prediction head will not happen and all neurons in the network will focus on learning the strongest feature. Formal arguments will be given in Section 5.3.

In contrast, [82, 87] did not incorporate the output normalization into their theory, even though their experiments have used certain forms of normalizations. We believe their method is closely related to the whitening method in [33]. To the best of our knowledge, our paper is the first to explain the effects of output-normalization in optimizing nonlinear neural networks in self-supervised learning.

**Dimensional collapse** Currently the only theoretical investigation on the dimensional collapse is by Jing et al. [52], where they focus on the contrastive learning setting. We believe their result on the role of the projection head is meaningful to understanding non-contrastive learning. But we emphasize that the objective (2.2) suffer from much more extreme dimensional collapse, as shown in Figure 1. Thus the causes described in Jing et al. [52] such as strong data augmentations cannot fully explain the dimensional collapse in the non-contrastive setting.

# 2 Preliminaries on Non-contrastive Learning

In this section, we formally define what is non-contrastive self-supervised learning. To do this, we first introduce contrastive learning following [24, 89] as background. We use [N] as a shorthand for the index set  $\{1,...,N\}$ .

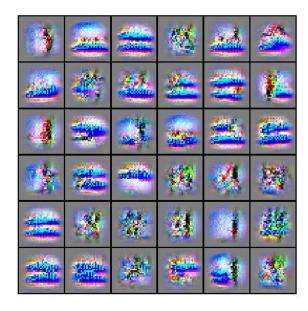
**Background on contrastive learning.** Letting  $\varphi_W(\cdot)$  be the neural networks, contrastive learning aims to learn good representations  $\varphi_W$  via contrasting representations of similar data samples to those of dissimilar ones. Usually we are given a batch of data points  $\{X_i\}_{i\in[N]}$ , and we construct for each  $i\in[N]$  a positive pair  $(X_i^{(1)},X_i^{(2)})$  (which are assumed to be simmilar) by applying ran-

 $X_i$ , and collect negative pairs  $(X_i, X_j)$  for i 6=  $j \in [N]$  (which are assumed to be dissimilar). Now given the representations  $z_i = \phi_W(X_i^{(1)}), \ z_i' = \phi_W(X_i^{(2)}), \ i \in [N]$ , we train the network  $\varphi_W$  to minimize the following contrastive loss:

$$L_{\text{contrastive}}(\phi_W) := \frac{1}{N} \sum_{i \in [N]} \underbrace{-\mathbf{sim}(z_i, z_i')/\tau}_{\text{positive term}} + \underbrace{\log \left[\sum_{j \in [N]} \exp\left(\mathbf{sim}(z_i, z_j')/\tau\right)\right]}_{\text{negative term}}$$
(2.1)

{z}





- (a) Features learned with prediction head
- (b) Features learned without prediction head

Figure 5: Feature visualization of deep neural network. We visualized the features of an Wide-ResNet-16x5 following the BYORL method by Gowal et al. [38], a adversarial robust version of BYOL. Features learned with prediction head obviously have more variety than features learned without the prediction head. Our feature visualization technique follows from [5].

where  $sim(\cdot, \cdot)$  is the similarity metric, often defined as the cosine similarity, and  $\tau$  is the so-called temperature hyper-parameter. Intuitively, minimizing the contrastive loss can be roughly viewed as trying to classify the representation  $z_i$  as  $z_i$ 0 instead of  $z_j$ 0, j 6= i. It is a common belief that in order for the network  $\varphi_W$  to be able to "distinguish" data points  $X_i$  from  $X_j$ , j 6= i, merely minimizing the positive term of contrastive loss is not sufficient.

As shown by the papers [25, 89], the performance of contrastive learning depends critically on the negative term. But the BYOL method [41] managed to remove the negative term without harm, by adding a trainable prediction head to the network architecture, which opened the new direction of non-contrastive self-supervised learning.

**Non-contrastive self-supervised learning.** We choose the SimSiam method [26] as our primary framework, whose difference with BYOL is a EMA component that is proven inessential in

[26]. Following the same notations as above, except that  $z_i^0 = \text{StopGrad}[\varphi_W(X_i^{(2)})]$  is detached from gradient computation, the loss objective become: (the symmetric network version)

$$= \frac{1}{N} \sum_{i \in [N]} - \sum_{i \in [N]} \sin(z_i, z_i')$$
 (2.2)

which is just the positive term in contrastive loss (2.1) (not divided by  $\tau$ ). Removing the negative term results in the existence of plenty trivial **global optimal** solutions. For example, the *complete collapse* refers to when  $\varphi_W(\cdot)$  is some constant vector function with zero variance. Another trivial solution called **dimensional collapse** [46], which is when all the coordinates  $[\varphi_W(\cdot)]_i$  has correlation  $\pm 1$ , meaning  $\varphi_W(\cdot)$  lies in a one-dimensional subspace of the representation space. The dimensional collapsed solution can minimize the objective (2.2) even when the network output  $\varphi_W(\cdot)$  is normalized by BN to avoid converging to a constant vector [46, 92].

However, by adding a *trainable prediction head* on top of  $z_i$ , the training miraculously succeeds and outputs a state-of-the-art feature extractor. Let  $g(\cdot)$  be a shallow feed-forward network (often one or two-layer, or even simply linear), we train g and  $\phi_W$  simultaneously on the following objective:

$$= \frac{1}{N} \sum_{i \in [N]} - L_{\text{SimSiam}} = \frac{1}{N} \sum_{i \in [N]} sim(g(z_i), z_{i0})$$
(2.3)

where  $z_i'$  is still detached from gradient computation. The  $g(z_i) = g \circ \phi_W(X_i^{(1)})$  and the detached part  $z_i^0$  = StopGrad  $[\phi_W(X_i^{(2)})]$  are often called the *online* network and the *target* network respectively following [41], known as two branches of non-contrastive learning. Even when such a trainable prediction head is able to represent identity function, the network can still avoid the common collapsed solutions, which presents challenges in understanding their training process and the underlying mechanism of trainable prediction head.

Challenges of understanding non-contrastive learning. Although the non-contrastive losses (2.2) and (2.3) seems just a term of the contrastive loss (2.1), their behaviors are vastly different. As established in [89], the negative pairs are needed for learning all the discriminative features. Without the negative term, the learner has no explicit incentive to learn all the discriminative features from the objective (2.3), especially when the trainable prediction head can possibly be an identity map.

Indeed, by setting  $g(\cdot)$  to the identity map, problem (2.3) immediately turn back into (2.2) and has the same trivial collapsed global optima. It is one of the most typical examples of *implicit bias* of optimization in deep learning.

Empirically, the seminal paper [26] discovered that *even with trainable linear prediction head which can possibly learn identity mapping*, neural networks trained by SGD still avoid such collapsed solutions. Moreover, as we show in this paper, even with an identity-initialized linear prediction head, as long as we train the prediction head via SGD, it still produces results comparable to when using other types of prediction head. Our empirical evidence in Figure 2 suggests that understanding the asymmetry provided by the off-diagonal entries in the identity-initialized linear prediction head suffices to explain (most of) the mechanisms of the prediction head. This observation significantly simplifies the theoretical problem and makes the complete characterization of the training dynamics of the prediction head possible.

Nevertheless, understanding the trainable prediction head urges us to go beyond the traditional statistical framework and optimization landscape analysis. The recent development of the *feature learning theory* of neural networks [48, 5, 3, 89, 49] showed it is possible to directly analyze the training dynamics of neural networks in various supervised or self-supervised tasks. Inspired by this line of research and our observations, we consider understanding the optimization of identityinitialized prediction head the key to understanding the underlying mechanism of these methods, and the characterization of the training dynamics of the full network the major technical challenges.

## 3 Problem Setup

In this section, we present the setting of our theoretical results. We first define the data distribution.

**Notations.** We use  $O,\Omega,\Theta$  notations to hide universal constants with respect to d and  $O,\Theta$  O0 notations to hide polynomial factors of logd. We denote a = o(1) if  $a \to 0$  when  $d \to \infty$ . We use

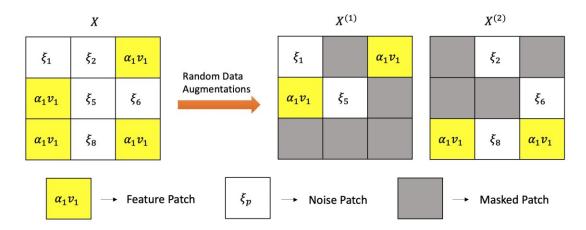


Figure 6: Illustration of the data distribution and data augmentations. Each data is equipped with a feature, either  $v_1$  or  $v_2$ , and contains a lot of noise patches. After the data augmentations, the positive pair  $(X^{(1)}, X^{(2)})$  is constructed by randomly masking out half of non-overlapping patches for each positive sample. The reason for

constructing positive pair with non-overlapping patches is because of the strong noise assumption we made in Assumption 3.3 and the *feature decoupling* principle in [89].

the notations poly(d), polylog(d) to represent large constant degree polynomials of d or log d. We use  $N(\mu,\Sigma)$  to denote standard normal distribution in with mean  $\mu$  and covariance matrix  $\Sigma$ . We use the bracket  $h\cdot$ ,  $\cdot$  i to denote the inner product and  $k\cdot k_2$  the  $\cdot$ 2-norm in Euclidean space. And for a subspace  $V \subset \mathbb{R}^d$ , we denote  $V^{\perp}$  as its orthogonal complement. We use  $\mathbf{1}_B$  to denote the indicator function of event B.

Following the standard structure of image datasets, we consider data divided into patches, where each patch can contain either features or noises.

**Definition 3.1** (data distribution and features). Let  $X \sim D$  be  $X = (X_1,...,X_P) \in \mathbb{R}^{d \times P}$  where each  $X_i \in \mathbb{R}^d$  is a patch. We assume that there are two feature vectors  $v_1,v_2$  such that  $kv \cdot k_2 = 1$ , = 1,2 and are orthogonal to each other. To generate a sample X, we uniformly sampled  $\in [2]$  and generate for each  $p \in [P]$ :

$$X_p = Z_p(X)v + \xi_p \mathbf{1}_{zp=0}, \text{ E} X \sim D[Z_p(X)] = 0, \qquad \forall p \in [P]$$

We denote  $S(X) = \{p : z_p(X) \in B\} \subseteq [P]$  as the set of feature patches and assume  $z_p(X) = z_{p^0}(X) \in \{0,\pm\alpha'\}, \forall p,p^0 \in S(X)$ , i.e., all feature patches have the same direction of v within the same X. We assume P = polylog(d),  $S(X) \equiv P_0 = \Theta(\log d)$  for every X. The assumption of  $\xi_p$  will be given in Assumption 3.3. An intuitive illustration is given in Figure 6.

**Strong and weak features.** We pick  $\alpha_1 = 2^{\text{polyloglog}(d)}$  and  $\alpha_2 = \alpha_1/\text{polylog}(d)$ . Hence  $v_1$  is the *strong feature* and  $v_2$  is the *weak feature*, and we want the learner network to learn both  $v_1, v_2$  (but by different neurons) as their learning goal. This is a simplification of the real scenario where features show up more consistently across multiple patches of the images, while noises are local and roughly independent across different patches. Intuitively, we can think of the strong features in a dataset are the ones that show up more frequently or with larger magnitude, and weak features as those that show up rarely or with smaller magnitude, which is the common case in any practical dataset.

Remark 3.2. Our analysis can be easily generalized to settings of either (1) when  $\alpha_1 = \alpha_2$  but the sampling of ` $\in$  [2] is of non-equal probability (i.e., dataset imbalance setting); or (2) when the two features always co-occur in the same sample but not of the same strength. But we still require  $\alpha_1, \alpha_2 \gg \mathsf{polylog}(d)$  to simplify the analysis.

**Assumption 3.3** (noise). Denoting  $V = \text{span}(v_1, v_2)$ , we assume  $\xi_p \in V^{\perp}$  is independent for each  $p \in [P] \setminus S(X)$ , where  $X = (X_p)_{p \in [P]} \sim D$ , and:

- (a) For any unit vector  $u \in V^{\perp}$ ,  $E[h\xi_p,ui] = 0$ , and  $E[h\xi_p,ui^6] = \sigma^6$  for some  $\sigma = \Theta(1)$ ;
- (b) It holds for some  $\varrho \in [0, \frac{1}{d^{\Omega(1)}}]$  it holds  $|\mathsf{E}[\mathsf{h}u_1, \xi_p \mathsf{i}^3 \mathsf{h}u_2, \xi_p \mathsf{i}^3]| \le \%$  and  $|\mathsf{E}[\mathsf{h}u_1, \xi_p \mathsf{i}^5 \mathsf{h}u_2, \xi_p \mathsf{i}]| \le \%$  for any two vectors  $u_1, u_2 \in \mathsf{R}^d$  that are orthogonal to each other.

*Remark* 3.4. A simple example of our noise  $\xi_p$  is the spherical Gaussian noise in  $V^{\perp}$ . Our Assumption 3.3b ensures that the prediction head cannot be used to cancel the noise correlation between different

neurons. We point out that the features in our data can be learned via clustering, but we emphasize that we do not intend to compare our algorithm with any clustering method in this setting since our goal is to study how the prediction head helps in learning the features.

#### 3.1 Learner Network

Following the SimSiam framework, the online and target network share the same encoder network in our setting, as explained in Section 2. We consider the base encoder network f as a simple convolutional neural network: Let  $W = (w_1,...,w_m) \in \mathbb{R}^{d \times m}$  be the weight matrix, where  $w_i \in \mathbb{R}^d$ , the **encoder network** f is defined by

$$f_j(X) := {\operatorname{P}}_{p \in [P]} \sigma(\operatorname{h} w_j, X_p \mathrm{i}), \quad \forall j \in [m]$$

Here we use the cubic activation function  $\sigma(z) = z^3$ , as polynomial activations are standard in literatures of deep learning theory [9, 35, 54, 2, 56, 23] and also has comparable performance in practice [2]. The (identity initialized) prediction head is defined as a matrix  $E = [E_{i,j}]_{(i,j) \in [m]^2}$  with  $E_{i,i} \equiv 1, i \in [m]$ , where only the the off-diagonals  $E_{i,j}$   $i \in [m]$  are trainable parameters. The **online network**  $F_e$  is defined by: given  $j \in [m]$ , we let  $F_j(X) := f_j(X) + \frac{P}{r_{6=j}}E_{j,r}f_r(X)$ , and

$$F_{\mathbf{e}_{j}}\!\!\left(\!\boldsymbol{X}\!\right) \coloneqq \mathsf{BN}\!\!\left[\sum_{p \in [P]} \left(\sigma(\langle w_{j}, X_{p} \rangle) + \sum_{r \neq j} E_{j,r} \sigma(\langle w_{r}, X_{p} \rangle\right)\right]$$

where the batch normalization BN here is defined as follows: Given a batch of inputs  $\{z_i\}_{i\in[N]}$ ,

$$(z_i) := \frac{z_i - \frac{1}{N} \sum_{i \in [N]} z_i}{\sqrt{\frac{1}{N} \sum_{i \in [N]} z_i^2 - \left(\frac{1}{N} \sum_{i \in [N]} z_i\right)^2}}$$
(3.1)

And the **target network** G is defined as follows: Given  $i \in [m]$ 

$$G_{\mathbf{e}_{j}}(X) := \mathsf{BN}(G_{j}(X)) = \mathsf{BN}^{\left[\sum_{p \in [P]} \sigma(\langle w_{j}, X_{p} \rangle)\right]}$$

#### Algorithm 1 Training Algorithm

**Require:** data distribution D, objective  $L_S$  (3.3), networks  $F_{e}$   $G_{e}$ , hyper-parameters  $T_{e}$ ,  $N_{e}$ ,  $N_{e}$ ,  $N_{e}$ ,  $N_{e}$ , and a bool variable TrainPredHead = True.

1: Initialize 
$$w_j^{(0)} \sim \mathcal{N}(0, I_d/d) \ \forall j \in [m]$$
 i.i.d., and  $E^{(0)} = I_m$ ; 2: **for**  $t \in \{0, 1, 2, \dots, T-1\}$  **do**

3: Sample  $X^{(t,i)} \leftarrow (X_p^{(t,i)})_{p \in [P]} \sim D, \forall i \in [N] \text{ i.i.d.};$ 

4: Sample  $\{P^{(t,i)}\}_{i \in [N]}$  i.i.d., and obtain  $S_t \leftarrow \{X^{(t,i,1)}, X^{(t,i,2)}\}_{i \in [N]}$  via data augmentations

$$X^{(t,i,1)} \leftarrow (X_p^{(t,i)} \mathbb{1}_{p \in \mathcal{P}^{(t,i)}})_{p \in [P]}, \qquad X^{(t,i,2)} \leftarrow (X_p^{(t,i)} \mathbb{1}_{p \notin \mathcal{P}^{(t,i)}})_{p \in [P]};$$

5: Perform stochastic gradient descent step to  $W^{(t)} = (w_j^{(t)})_{j \in [m]}$  by

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \nabla_{w_j} L_{\mathcal{S}_t}(W^{(t)}, E^{(t)})$$

6: **if** TrainPredHead = True **then** update the off diagonal of prediction head  $E^{(t)}$  by

$$E_{i,i}^{(t+1)} \leftarrow 1, \quad E_{i,j}^{(t+1)} \leftarrow E_{i,j}^{(t)} - \eta_E \nabla_{E_{i,j}} L_{\mathcal{S}_t}(W^{(t)}, E^{(t)}), \quad \forall j \neq i, \ i, j \in [m_1];$$

- 7: **else** keep  $E^{(t+1)} = I_m$ .
- 8: **end if**
- 9: end for

#### 3.2 Training Algorithm

**Data augmentation.** We use a very simple data augmentation: for each data  $X = (X_p)_{p \in [P]}$ , we randomly and uniformly sample half of the patches  $P \subseteq [P]$  to generate two samples (which is the so-called *positive pair* in contrastive learning):

$$X^{(1)} = (X_p \mathbf{1}_{p} \in P)_p \in [P], \qquad X^{(2)} = (X_p \mathbf{1}_{p/e} P)_p \in [P]$$
(3.2)

An intuitive illustration is given in Figure 6. Our data augmentation approach is similar to the common cropping augmentation used in contrastive learning [22, 80] and the patch masking strategy in generative pretraining [16, 44] and NLP pretraining [30]. It is also analogous to the data augmentations being studied in theoretical literatures [89, 50, 62] of self-supervised learning, especially the RandomMask augmentation in [89].

**Non-contrastive loss function.** Now we define the loss function as follows: we sample N data points  $\{X_i\}_{i\in[N]_*}X_i^{i,\sim}D_{i.d.}$  and apply our data augmentation (3.2) to obtain  $S=\{X^{(i,1)},X^{(i,2)}\}_{i\in[N]}$ . Now we define

$$\begin{split} L_{\mathcal{S}}(W,E) &:= \frac{1}{N} \sum_{i \in [N]} \left\| \widetilde{F}(X^{(i,1)}) - \underset{\text{StopGrad}}{\left[ \widetilde{G}(X^{(i,2)}] \right)} \right\|_{2}^{2} \\ &= 2 - \frac{1}{N} \sum_{i \in [N]} \langle \widetilde{F}(X^{(i,1)}) , \underset{\text{StopGrad}}{\text{StopGrad}} [G_{\mathbf{e}}(X^{(i,2)})] \mathbf{i} \end{split}$$
(3.3)

where the StopGrad operator detach gradient computation of the target network  $G_e(\cdot)$ . This form of objective (3.3) is first defined in Guo et al. [41] and is equivalent to (2.3) in Chen and He [26] when

 $F_{\rm e}$  and  $G_{\rm e}$  share the same encoder network  $f(\cdot)$  and their outputs are normalized.

**Intuition of the data augmentation.** Our data augmentation is an analog of the standard cropping data augmentation. In Definition 3.1, the features  $v_1, v_2$  appear in multiple patches, but the

noises are independent across different patches (see Figure 6). As our data augmentation produces positive pairs with non-overlapping patches, learning to emphasize noises cannot align the representations of the positive pair, but learning **either one of** the features  $\varphi(X) = {}^{P}_{p} \sigma(hv_1, X_p i)$  or  $\varphi(X) = {}^{P}_{p} \sigma(hv_2, X_p i)$  is sufficient. **We consider learning the same feature**  $v_i$  **in** all the neurons  $f_j$  in the encoder network f **as the dimensional collapsed solution.** 

Initialization and hyper-parameters. At t=0, we initialize W and E as  $W_{i,j}^{(0)} \sim \mathcal{N}(0,\frac{1}{d})$  and  $E^{(0)} = I_m$  and we only train the off-diagonal entries of  $E^{(t)}$ . For the simplicity of analysis, we let m=2, which suffices to illustrate our main message. For the learning rates, we let  $\eta \in (0,\frac{1}{\mathsf{poly}(d)}]$  be sufficiently small and  $\eta_E \in [\eta_1^{\eta}, \eta_2^{\eta}]$ , which is smaller than  $\eta^1$ .

**Optimization algorithm** Given the data augmentation and the loss function, we perform (stochastic) gradient descent on the training objective (3.3) as follows: at each iteration t = 0,..., T - 1, we sample a new batch of augmented data  $S_t = \{X^{(t,i,1)}, X^{(t,i,2)}\}_{i \in [N]}$  and update

$$W^{(t+1)} = W^{(t)} - \eta \nabla_W L_{\mathcal{S}_t}(W^{(t)}, E^{(t)}), \quad E^{(t+1)}_{i,j} = E^{(t)}_{i,j} - \eta_E \nabla_{E_{i,j}} L_{\mathcal{S}_t}(W^{(t)}, E^{(t)}), \quad \forall i \neq j, i, j \in [m]$$

If we do not train the prediction head, we just simply keep  $E^{(t)} \equiv I_m$ . We summarize our algorithm in Algorithm 1.

#### 4 Statements of Main Results

In this section, we shall present our main theoretical results on the mechanism of learning the prediction head in non-contrastive learning. To measure the correlation between neurons, we introduce the following notion: letting

$$Var(\psi(X)) := E_{X \sim D}[(\psi(X) - E[\psi(X)])^2]$$

be the variance of any function  $\psi$  of  $X \sim D$ , we denote the correlation  $\mathbf{Corr}(\psi(X), \psi^0(X))$  of any two function  $\psi, \psi^0$  over D as

$$\mathbf{Corr}^{(\psi(X),\psi'(X))} := \frac{\mathbb{E}[(\psi(X) - \mathbb{E}[\psi(X)])(\psi'(X) - \mathbb{E}[\psi'(X)])]}{\sqrt{\mathbf{Var}(\psi(X))}\sqrt{\mathbf{Var}(\psi'(X))}}$$

Now we present the main theorem of training with a prediction head, and set m = 2.

**Theorem 4.1** (learning with prediction head and BN, see Theorem E.2). For every d > 2, let

 $N \ge \mathsf{poly}^{(d)}, \ \eta \in (0, \frac{1}{\mathsf{poly}(d)}]$  be sufficiently small, and  $\eta_E \in [\frac{\eta}{\alpha_1^{O(1)}}, \frac{\eta}{\mathsf{polylog}(d)}]$  \_\_\_\_\_\_. Then with probability

1−o(1), after runing Algorithm 1 for  $T = poly(d)/\eta$  many iterations, we shall have for some `∈ [2]:

<sup>&</sup>lt;sup>1</sup> We conjecture that by modifying certain assumptions for the noise (especially by allowing the noise to span the feature subspace V), one can prove a similar result for the case  $\eta_E = \eta$ .

$$\begin{split} w_1^{(T)} &= \beta_1 v_\ell + \varepsilon_1, \quad w_2^{(T)} = \beta_2 v_{3-\ell} + \varepsilon_2 & \text{with} \\ & Furthermore, the objective converges: } \mathsf{E}_{\mathsf{S} \sim \mathsf{D}^N}[L_\mathsf{S}(\mathit{W}^{(T)}, \mathit{E}^{(T)})] \leq \mathsf{OPT}^{\,+\, \frac{1}{\mathsf{poly}(d)}} \leq O(\frac{1}{\log d})^2. \end{split}$$

Theorem 4.1 clearly shows the network learn all the desired features, even under huge imbalance between  $v_1$  and  $v_2$ . This leads to the following corollary.

**Corollary 4.2.** Under the same hyper-parameter in Theorem 4.1, with probability 1 - o(1), after runing Algorithm 1 for  $T = \text{poly}(d)/\eta$  many iterations, we shall have that the learning avoids dimensional collapse:

$$|\mathbf{Corr}^{(f_1(X), f_2(X))}| \le O(\frac{1}{\sqrt{d}}).$$

In contrast, learning without the prediction head will result in learning only the strong feature  $v_1$  in both neurons, which creates strong correlations between any two neurons. To emphasize that this problem cannot be alleviated by having more neurons, we let the number of neurons m be any positive integer in the following theorem.

**Theorem 4.3** (learning without prediction head but with BN, see Theorem F.1). Let  $N \ge \text{poly}(d)$ ,  $\eta = o(1)$  and the number of neurons m > 0 be any positive integer. Suppose we freeze  $E^{(t)} = I_m$  for all t, then with probability 1 - o(1), after runing Algorithm 1 with TrainPredHead = False for  $T = \text{poly}(d)/\eta$  many iterations, we shall have:

$$w_j^{(T)} = \beta_j v_1 + \varepsilon_j \qquad \text{with} \qquad |\beta_j| = \Theta(1), \ \|\varepsilon_j\|_2 \le \widetilde{O}(\frac{1}{\sqrt{d}}) \qquad \qquad \text{for all } j \in [m]$$

Furthermore, the objective converges:  $\mathsf{E}_{\mathsf{S}\sim\mathsf{D}^N}[L_\mathsf{S}(W^{(T)},E^{(T)})] \leq \mathsf{OPT}^{+\frac{1}{\mathsf{poly}(d)}} \leq O(\frac{1}{\log d})$ . This means the collapsed solution also reaches the global minimum of the objective.

Note that since we have used BN as our output normalization instead of `2-norm, the learner is immune to complete collapse and must have a certain variance in the outputs. Immediately, we have the following corollary.

**Corollary 4.4.** Under the same hyper-parameter in Theorem 4.3, with probability 1 - o(1), after runing Algorithm 1 with TrainPredHead = False for  $T = \text{poly}(d)/\eta$  many iterations, we shall have dimensional collapse:

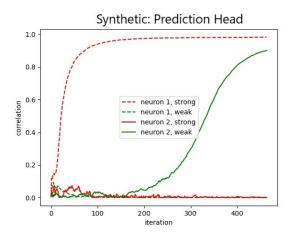
$$|\mathbf{Corr}^{(f_i(X),f_j(X))}| \ge 1 - O(\frac{1}{\sqrt{d}})$$
 for all  $i,j \in [m]$ .

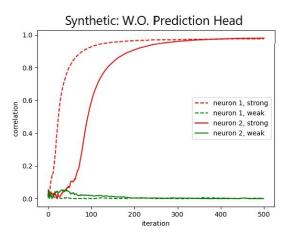
Remark 4.5. Note that since we have used BN as our output normalization instead of `2-norm, the learner is regularized to avoid complete collapse and must have a certain variance in its neurons. It is easier to obtain a complete collapse result when the network has `2-normalized outputs and there is

 $<sup>\</sup>begin{tabular}{ll} $^2$ Under our data model Definition 3.1, non-overlapping data augmentation (3.2) and learner network definition, the $$ global minimum of our objective $$ 3.2 \text{in population} & $\mathbb{E}[|S(X) \cap \mathcal{P}|:|S(X) \setminus \mathcal{P}|]$ $$ = $\Theta(\frac{1}{\log d})$ $$$ 

a low-variance feature (but not of smaller magnitude) in the data set, which we refrain from proving here.

How does using the prediction head or not create such a difference in features learned by the non-contrastive methods? We shall give some intuitions by digging through the training process and separately discuss the four phases of the training process.





- (a) Identity-initialized (trainable) prediction head
- (b) Learning without prediction head

Figure 7: The feature learning process over synthetic data. When trained with the prediction head, after the strong feature is learned in the faster learning neuron, the weak feature can be learned in the slower learning neuron. When trained without the prediction head, both neurons will learn the strong feature and ignore the weak feature.

# 5 The Four Phases of the Learning Process

We divide the complete training process into four phases: phase I for learning the stronger feature, phase II for the substitution effect, phase III for the acceleration effect, and the end phase for convergence. The first three phases explain how the prediction head can help learn the base encoder network, and the last phase of the training explains why the off-diagonal entries often shrink in the later stage of training.

#### 5.1 Phase I: Learning the Stronger Feature

At the beginning of training, the stronger feature  $v_1$  enjoys a much larger gradient as opposed to the weaker feature  $v_2$ , so naturally,  $v_1$  will be learned first. However, if for both neurons  $f_1, f_2$  the speed of learning  $v_1$  is the same, then we cannot argue the difference between them and will not be able to show the substitution from either one to another. Indeed, let us assume at initialization, the neuron  $w_1^{(0)}$  won the jackpot of having larger signal-to-noise ratio of  $hw_j^{(0)}, v_1$  i between  $f_j, j \in [2]$ , then we can show the following result under our setting.

**Lemma 5.1** (learning the stronger feature, see Lemma B.13). After some  $t \ge T_1 = d^{2+o(1)}/\eta$ , the feature  $v_1$  in neuron  $f_1$  will be learn to  $\langle w_1^{(t)}, v_1 \rangle = \Omega(1)$ , while all other features  $\langle w_j^{(t)}, v_\ell \rangle = o(1)$  for (j, ') 6= (1,1) are small. And the prediction head  $kE^{(t)} - I_2k_2 \le d^{-\Omega(1)}$  is still close to the initialization.

In this phase, the prediction head has not come into play. The substitution effect can only happen after the feature  $v_1$  in neuron  $f_1$  is learned to a certain degree, and neuron  $f_2$  remains largely unlearned.

#### 5.2 Phase II: The Substitution Effect

To illustrate the substitution effect, let us keep assuming that neuron  $w_1^{(t)}$  has already learned some significant amount of the strong feature  $v_1$ , say  $w_1^{(t)} = \beta_1 v_1 + residual$  with  $|\beta_1| = \Omega(kresidualk)$ . When this happens, we have the following result: (recall  $f_j(\cdot), j \in [2]$  are the neurons of the base encoder network)

**Lemma 5.2** (substitution effect, formal statement see Lemma C.8). After  $|\langle w_1^{(t)}, v_1 \rangle| = \Omega(1)$  in  $O(d^{2+o(1)}/\eta)$  iterations (as shown by Lemma B.13), for much shorter time than learning  $\langle w_1^{(t)}, v_1 \rangle$ , we shall have  $|E_{2,1}^{(t)}|$  increasing until  $|E_{2,1}^{(t)}f_1(X^{(1)})| \gg |f_2(X^{(1)})|$  when X is equipped with feature  $v_1$ . In other words,  $E_{2,1}^{(t)}f_1(X^{(1)})$  is a substitute for the feature  $v_1$  that should be learned by  $f_2$ .

**Intuition of the substitution effect.** After the stronger feature is learned in neuron  $f_1$ , the optimal way to align two positive representations  $F_2(X^{(1)})$ ,  $G_2(X^{(2)})$  is no longer learning features in weight  $w_2$ , but use the prediction head to "borrow" the features in  $f_1$  and incorporate them into  $F_2$ . This is how the substitution effect happens when trained with a prediction head.

**Proof sketch for Lemma 5.2.** Indeed, let us look at the learning of. In this are roughly learned to maximize the following quantity:  $E_{2,1}^{(t)}$  phase,  $w_2$  and  $E_{2,1}^{(t)}$ 

$$\begin{split} \widetilde{F}_2(X^{(1)}) \cdot \widetilde{G}_2(X^{(2)}) &\propto \left( f_2(X^{(1)}) + E_{2,1}^{(t)} f_1(X^{(1)}) \right) \times f_2(X^{(2)}) \\ &\approx \sum_{\ell \in [2]} \alpha_\ell^6 \left( \langle w_2^{(t)}, v_\ell \rangle^6 + E_{2,1}^{(t)} \cdot \langle w_1^{(t)}, v_\ell \rangle^3 \cdot \langle w_2^{(t)}, v_\ell \rangle^3 \right) \end{split}$$

As the neuron  $f_1(\cdot)$  is already learned with feature  $v_1$ , in order to maximize the RHS, we can either try to maximize  $\sum_{\ell \in [2]} \langle w_2^{(t)}, v_\ell \rangle^6$ , or to maximize  $E_{2,1}^{(t)} \langle w_1^{(t)}, v_1 \rangle^3 \cdot \langle w_2^{(t)}, v_1 \rangle^3 \approx E_{2,1}^{(t)} \langle w_2^{(t)}, v_1 \rangle^3$ . In this case, the more efficient choice is to learn  $|E_{2,1}^{(t)}|$  to substitute for maximizing  $\langle w_2^{(t)}, v_\ell \rangle^3$ . Actually, because of the high signal-to-noise ratio of learning  $w_2^{(t)}$  than  $E_{2,1}^{(t)}$ , feature  $\langle w_2^{(t)}, v_\ell \rangle$  is learned with slower pace than  $E_{2,1}^{(t)}$ , so that Lemma 5.2 can be shown.

#### 5.3 Phase III: The Acceleration Effect

After the substitution of  $v_1$  in  $F_2$ , our concern is, whether or not  $w_2$  will learn  $v_2$  and only  $v_2$  eventually, so that we can obtain a diverse representation? The answer is yes, as we summarize in the following lemma.

**Lemma 5.3** (acceleration effect, formal statement see Lemma D.8). *After*  $E_{2,1}^{(t)}$  *is learned in Lemma 5.2, learning*  $v_2$  *in*  $w_2^{(t)}$  *will be much faster than*  $v_1$ , *until*  $||w_2^{(t)} - \beta_2 v_2|| \le o(1)$  *for some*  $\beta_2 = \Theta(1)$ .

The acceleration effect is caused by the interactions between the prediction head, the stop gradient operation, and the normalization method (which in this case is the batch normalization). We shall explain these interactions with insights from our theoretical analyses below.

What is the role of the stop-gradient? Thanks to the StopGrad operation, when we compute the gradient  $-\nabla_{w^2}F_2(X^{(1)})$  · StopGrad[ $G_2(X^{(2)})$ ] to learn  $f_2$ , this negative gradient will only try to maximize  $f_2(X^{(1)}) \cdot f_2(X^{(2)})$ , rather than to maximize  $f_2(X^{(2)}) \cdot F_2(X^{(1)})$ . This is because the stop-gradient is on G not on G: while G0 has a large component of G1 borrowed from G1 using G2 does not have this component. So the gradient of G2 is to align with the features in G2 that does not contain many G1, while the gradient of G2 is to aligned with the features in G2 that contains a lot of G1. Thus the stop gradient on G3 help ignore the feature borrowed from G1 using prediction head G2 and ensures the slower learning neuron G3 will focus on learning feature G3.

What is the role of the output normalization?

Again due to the StopGrad operation, the

gradient of  $F_{e2}$  is taken with respect to the ratio  $f_2(X^{(1)})/{}^p\mathbf{Var}[F_2(X^{(1)})]$ . As gradient descent tries to maximize this ratio, a direct computation gives

$$\nabla_{w_2} \frac{f_2(X^{(1)})}{\sqrt{\mathbf{Var}(F_2(X^{(1)}))}} = \frac{\nabla_{w_2} f_2(X^{(1)}) \cdot \mathbf{Var}(F_2(X^{(1)})) - f_2(X^{(1)}) \cdot \nabla_{w_2} \mathbf{Var}(F_2(X^{(1)}))}{\mathbf{Var}(F_2(X^{(1)}))^{3/2}}$$

From some calculation, we can obtain the above gradient is proportional to

$$\sum_{\ell \in [2]} \left( [E_{2,1}^{(t)} \langle w_1^{(t)}, v_{3-\ell} \rangle^3]^2 + \mathbf{Var}[f_2(X^{(1)})] \right) \langle \nabla_{w_2} f_2(X^{(1)}), v_{\ell} \rangle v_{\ell}$$

which borrow the *substituted feature*  $v_3$ -from  $f_1(\cdot)$  to adjust the gradient of  $v \cdot \inf_2(\cdot)$ , via the prediction head  $E_{2,1}^{(t)}$ . Without the output normalization, the learning of  $v_1$  will dominate that of  $v_2$  even when we train the prediction head.

(t)

**Proof sketch for Lemma 5.3.** At this stage, when we are updating the weights of  $w_2$ , we are simultaneuously maximizing  $f_2(X^{(1)}) \cdot f_2(X^{(2)})$  and also minimizing the normalizing constants

 ${}^{\mathbf{p}}\mathbf{Var}[F_2(X^{(1)})]$ . This two goals are in slight conflict because of the normalization, and by careful calculation the gradients are roughly given by (interpreting the expectation as empirical)

$$\langle -\nabla_{w_2} L_{\mathcal{S}}, v_{\ell} \rangle \propto \mathbb{E}\left[\left( [E_{2,1}^{(t)} \langle w_1^{(t)}, v_{3-\ell} \rangle^3]^2 + \mathbf{Var}[f_2(X^{(1)})] \right) \cdot f_2(X^{(2)}) \langle -\nabla_{w_2} f_2(X^{(1)}), v_{\ell} \rangle \right]$$

Because of the learning of  $f_1$  and the substitution effect, we now knows  $[E_{2,1}^{(t)}\langle w_1^{(t)}, v_{3-\ell}\rangle^3]^2$  is much larger when `= 2, which accelerates the learning of  $v_2$  in  $w_2^{(t)}$  to surpass that of  $v_1$  and leads to Lemma 5.3.

#### 5.4 The End Phase: Convergence

As the weak features are learned, we have already obtained a good encoder network  $f(\cdot)$  as shown in Theorem 4.1. The rest of our analysis is to understand what the prediction head converges to in polynomial time. Actually, our Theorem E.2 also contains the following result:

**Proposition 5.4** (convergence of the prediction head, see Theorem E.2c). *After some*  $t \ge T = \text{poly}(d)/\eta$  *iterations, we shall have*  $||E^{(t)} - I_2||_F \le \frac{1}{\text{poly}(d)}$ .

This result also implies that after learning the weak feature  $v_2$  is complete, the off-diagonal entries of the prediction head will reverse their trajectory and converge to zero at the end of training. While

we admit that only some of our real-world experiments show the convergence to zero for the off-diagonal entries of the prediction head, most of the experiments do display a rise and fall trajectory pattern of off-diagonal entries consistently.

#### 6 Additional Related Work

**Self-supervised learning** The area of self-supervised learning has evolved at a tremendous speed in recent years. It has created huge success in natural language processing [30, 90, 18] and established a paradigm where the networks are first trained on an unsupervised pretext task and then be finetuned in downstream applications. In vision, supervised pretraining had been the go-to choice until representations learned by contrastive learning [79, 43, 24, 20, 27, 28, 34, 68, 33] became dominant in many downstream tasks. Another type of self-supervised learning is the generative learning [69, 16, 44], which also gives promising results in downstream adaptations. Interesting applications such as [68, 70] also illustrate the power of contrastive learning in multiple domains.

Theory of self-supervised learning The theoretical side of self-supervised learning developed quickly due to the success of contrastive learning, which is closely related to the methods we are studying. Since Arora et al. [12], lots of papers have studied the properties of contrastive learning, as mentioned in the introduction. [25, 73] discussed many interesting phenomena associated with the negative term in contrastive learning. Saunshi et al. [75] provided pieces of evidence that contrastive loss is function class-specific rather than agnostic. Wen and Li [89] took a feature learning view to understand contrastive learning with neural networks, which inspired our analysis in the noncontrastive setting. For generative self-supervised learning, [55, 78] provides downstream performance guarantees for generative pretrained models. [74, 88] studied the natural language tasks, where the data are sequentially structured. Liu et al. [62] gave a recovery guarantee for tensors in generative learning under hidden Markov models. [4] analyzed multi-layer generative adversarial networks and provided an optimization guarantee for their stochastic gradient descent ascent algorithm.

**Feature learning theory of deep learning** Our theoretical results are also inspired by the recent progress of the feature learning theory of neural networks [59, 60, 5, 3, 53, 94, 48]. Li et al. [59] initiate the study of the speed difference in learning different types of features. [60] developed theory for learning two-layer neural networks over Gaussian distribution beyond the *neural tangent kernel* (NTK) [7, 8, 6, 32, 11]. Allen-Zhu and Li [5] studied the origin of adversarial examples and how adversarial training help in robustify the networks. [3] tried to explain ensemble and knowledge distillation under multi-view assumptions. Techniques in this paper are built on this line of research, as the non-convex nature of these analyses allows us to describe the interaction between neural networks, optimization algorithms, and the structures of data. [1, 2] also obtained results separating deep neural networks and shallow models such as kernel methods. Before this recent progress, [81, 93, 19, 76, 31, 57, 58] also studied how shallow neural networks can learn on certain simple data distributions, but all of them focus on the supervised learning. There are also plenty of studies [77, 40, 10, 64, 51, 71, 29] on the implicit bias of optimization in deep learning, but none of their techniques can be applied to the setting of self-supervised learning.

### 7 Conclusion and Discussion

In this paper, we showed how the prediction head can ensure the neural network learns all the features in non-contrastive learning through theoretical investigation. Our key observation is that the prediction head can leverage two effects called substitution effect and acceleration effect during the training process. We also explained how the necessary components such as output normalization and stop-gradient operation are involved and how they interact during training. Furthermore, we proved that without the prediction head, all neurons of the neural network would focus on learning the strongest feature and result in a collapsed representation. We believe our theory, although based on a very simple setup, can provide some insights into the inner workings of non-contrastive selfsupervised learning. We also believe our theoretical framework can be extended to understanding other phenomena in the practice of deep learning.

On the other hand, our results are still very preliminary, we point out the following open problems that are not addressed by this paper:

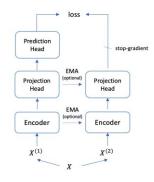
- When the output normalization is `2-norm instead of BN. Experiments in Figure 2 seem to suggest that there is still a gap between using `2-norm and BN as output normalization methods. In this case, the acceleration effect may not happen in exactly the same way as in the BN case, but we believe they share the same underlying mechanism and can be proven in theory.
- The mystery of the projection head. As our experiments in Figure 1 showed, the outputs of the projection head in the symmetric case (without the prediction head) suffer an extremely strong correlation even with batch normalization used. However, the impact on the base encoder is milder and thus the network can avoid complete collapse, shown in Figure 1 and Figure 2. It is mysterious how the projection head works in non-contrastive learning, and also how it compares to the case of contrastive learning, which has been studied by [24, 52]. Learning non-linearly features. For the simplicity of analysis, we have assumed the features in the data set are linear. It is of interest to study whether neural networks trained by non-contrastive self-supervised learning can learn non-linear representations better than traditional learning methods such as linear regression or kernel methods, as there has been a series of papers [1, 36, 37, 2, 53] trying to understand it in the supervised setting.

In the end, we also point out that theories based on a one-hidden-layer neural network and linear data composition assumption obviously cannot explain all the phenomena in deep learning. In supervised learning, the *backward feature correction* [2] process is observed and theoretically proven as a mechanism for learning hierarchical feature extractors. It is an important open direction to understand how a multi-layer network can learn the complicated features in non-contrastive selfsupervised learning.

# 8 Experiment Details

The framework we use in our experiments is shown in Figure 8. We use a modified version of the codebase shared by the authors of [33], Figure 8: Framework.

and we use the same data augmentation in their implementation. All our experiments (except for Figure 5 and Figure 7) use the following architecture and hyper-parameters: we choose standard ResNet-18 as base encoder architecture, 0.003 as the learning rate for Adam optimizer, a two-layer MLP with ReLU activation and 512 hidden neurons as the projection head, an identity-initialized but diagonally froze linear matrix (with shape (64x64)) as the prediction head and a non-tracking-stats, non-affine, non-momentum BN layer as the output normalization. Our experiments in Figure 3 use the same architecture and hyper-parameters, but some runs are trained with EMA with momentum 0.99, with output



BN replaced by '2-norm or using different prediction heads (such as a two-layer MLP or a linear head, with Pytorch default initialization). Evaluation in Figure 2 is by training a linear classifier on top of frozen encoder with no data augmentation.

# **Appendix: The Proofs**

We will be working with population gradients throughout the entire appendix. Indeed, since our algorithms use fresh random samples at each iteration, one can easily obtain from standard concentration inequalities an empirical estimate of population gradients up to poly error with N = poly(d) samples. So we can obtain the same proofs in finite sample case as long as the training ends before some  $T = poly(d)/\eta$ . Now we give some notations and warm-up calculations.

#### **Notations and Gradients** A

 $R_i := h \Pi v_{\perp} w_i, w_i$ 

In this section, we will give some useful notations and warm-up computations for the technical proofs in subsequent sections. We summarize here the notations that will also be defined in later sections:

Notations. We denote 
$$\mathcal{E}_j = \mathbb{E}[\langle w_j, \xi_p \rangle^6]$$
,  $\mathcal{E}_{j,3-j} = \mathbb{E}\left[(\langle w_j, \xi_p \rangle^3 + E_{j,3-j} \langle w_{3-j}, \xi_p \rangle^3)^2\right]$ , and  $C_0 = \frac{\mathbb{E}[|S(X) \cap \mathcal{P}| \cdot |S(X) \setminus \mathcal{P}|]}{2}$ ,  $C_1 = \frac{\mathbb{E}\left[|S(X) \cap \mathcal{P}|^2\right]}{2}$ ,  $C_2 = P - |S(X)|$ ,  $B_{j,3} = \operatorname{StopGrad}[hw_{j,\nu}i^3]$ ,  $B_{j,3} = hw_{j,\nu}i$ ,  $Q_j = \left(\operatorname{E}[\operatorname{StopGrad}[G^2_j(X^{(2)})]]\right)^{-1/2}$ . and

$$U_{j} := E[F_{j2}(X_{(1)})] = P \in [2] C_{1}\alpha'_{6}(B_{j,'3} + E_{j,3-j}B_{33-j,'})_{2} + C_{2}E_{j,3-j}$$

$$H_{j,'} := C_{1}\alpha'_{6}(B_{j,'3} + E_{j,3-j}B_{33-j,'})_{2} + C_{2}E_{j,3-j},$$

$$K_{j,'} := C_{1}\alpha'_{6}(B_{j,'3} + E_{j,3-j}B_{33-j,'})(B_{j,33-'} + E_{j,3-j}B_{33-j,3-'})$$

$$Moreover, we denote  $\Phi_{j} := Q_{j}/U_{j}$ , and  $(recall\ V := span(v_{1},v_{2}))$ 

$$\overline{R}_{1,2} := \frac{\langle \Pi_{V^{\perp}}w_{1}, w_{2} \rangle}{\|\Pi_{V^{\perp}}w_{1}\|_{2}\|\Pi_{V^{\perp}}w_{2}\|_{2}}$$$$

 $R_{1,2} := h \Pi v_{\perp} w_{1,w_{2}i}$ 

For any 
$$j \in [2]$$
, the gradient  $-\nabla_{wj}L(\textit{W,E})$  can be decomposed as  $-\nabla_{w_j}L(\textit{W},E) = \sum_{\ell \in [2]} (\Lambda_{j,\ell} + \Gamma_{j,\ell} - \Upsilon_{j,\ell})v_\ell - \sum_{(j',\ell) \in [2] \times [2]} \Sigma_{j',\ell}\nabla_{w_j}\mathcal{E}_{j',3-j'}$  
$$\Lambda_{j,\ell} := C_0\Phi_j\alpha_\ell^6B_{j,\ell}^5H_{j,3-\ell}$$
 
$$\Gamma_{j,\ell} := C_0\Phi_{3-j}E_{3-j,j}\alpha_\ell^6B_{3-j,\ell}^3B_{j,\ell}^2H_{3-j,3-\ell}$$
 
$$\Upsilon_{j,\ell} := C_0\alpha_{3-\ell}^6 \left(\Phi_jB_{j,3-\ell}^3B_{j,\ell}^2K_{j,\ell} + \Phi_{3-j}E_{3-j,j}B_{3-j,3-\ell}^3B_{j,\ell}^2K_{3-j,\ell}\right)$$
 
$$\Sigma_{j,\ell} := C_0C_2\Phi_j\alpha_\ell^6B_{j,\ell}^3(B_{j,\ell}^3 + E_{j,3-j}B_{3-j,\ell}^3)$$

Sometimes we need to decompose  $Y_{j,k} = Y_{j,k,k} + Y_{j,k,k}$  which is straightforward from its expression. In Section D, we further define

$$\Xi_{j}^{(t)} = C_{0}C_{1}\alpha_{1}^{6}\alpha_{2}^{6}\Phi_{j}^{(t)}\left((B_{1,1}^{(t)})^{6}(B_{2,2}^{(t)})^{6} + (B_{2,1}^{(t)})^{6}(B_{1,2}^{(t)})^{6}\right)$$
  
$$\Delta_{j,\ell}^{(t)} = C_{0}\Phi_{j}^{(t)}\alpha_{\ell}^{6}(B_{j,\ell}^{(t)})^{3}(B_{3-j,\ell}^{(t)})^{3}C_{2}\mathcal{E}_{j,3-j}^{(t)}$$

for the gradients of the prediction head.

#### **A.1 Gradient Computation**

Let us L(W,E) to be the population version of the objective. Because  $E[F_i(X^{(1)})]$  and  $E[G_i(X^{(2)})]$  are both zero (which can be verified easily from the zero-mean assumptions of  $z_p(X)$  and  $\xi_p$ ), a direct computation gives:

We first calculate the normalizing quantity  $E[F_i^2(X^{(1)})]$ :

$$\mathbb{E}[F_j^2(X^{(1)})] = \mathbb{E}\left[\left(\sum_{p \in [P]} \sigma(\langle w_j, X_p^{(1)} \rangle) + E_{j,3-j}\sigma(\langle w_{3-j}, X_p^{(1)} \rangle)\right)^2\right]$$

$$= \frac{1}{2} \sum_{\ell \in [2]} \mathbb{E}\left[|S(X) \cap \mathcal{P}|^2 \alpha_{\ell}^6 (\langle w_j, v_{\ell} \rangle^3 + E_{j,3-j} \langle w_{3-j}, v_{\ell} \rangle^3)^2\right]$$

(Because all signal patches has the same sign within the same data)  $+ \mathbb{E}\left[|\mathcal{P} \setminus S(X)|(\langle w_i, \xi_p \rangle^3 + E_{i,3-i}\langle w_{3-i}, \xi_p \rangle^3)^2\right]$ 

(Because noise patches are independent and have mean zero)

$$= \sum_{\ell \in [2]} \alpha_{\ell}^{6} (\langle w_{j}, v_{\ell} \rangle^{3} + E_{j,3-j} \langle w_{3-j}, v_{\ell} \rangle^{3})^{2} \frac{\mathbb{E} \left[ |S(X) \cap \mathcal{P}|^{2} \right]}{2} + (P - |S(X)|) \mathcal{E}_{j,3-j}$$

where we let

$$\mathcal{E}_{j,3-j} \stackrel{\text{def}}{=} \mathbb{E} \left[ (\langle w_j, \xi_p \rangle^3 + E_{j,3-j} \langle w_{3-j}, \xi_p \rangle^3)^2 \right] \\ = \mathbb{E} \left[ \langle w_j, \xi_p \rangle^6 + 2E_{j,3-j} \langle w_j, \xi_p \rangle^3 \langle w_{3-j}, \xi_p \rangle^3 + E_{j,3-j}^2 \langle w_{3-j}, \xi_p \rangle^6 \right]$$

On the other hand, we have

$$\begin{split} & \mathbb{E}[F_{j}(X^{(1)}) \cdot \mathsf{StopGrad}[G_{j}(X^{(2)})]] \\ &= \mathbb{E}\left[\left(\sum_{p \in [P]} \sigma(\langle w_{j}, X_{p}^{(1)} \rangle) + E_{j,3-j}\sigma(\langle w_{3-j}, X_{p}^{(1)} \rangle)\right) \times \left(\sum_{p \in [P]} \sigma(\langle w_{j}, X_{p}^{(2)} \rangle)\right)\right] \\ &= \frac{1}{2} \sum_{\ell \in [2]} \mathbb{E}\left[\sum_{p \in S(X) \cap \mathcal{P}} \alpha_{\ell}^{3}(\langle w_{j}, v_{\ell} \rangle^{3} + E_{j,3-j}\langle w_{3-j}, v_{\ell} \rangle^{3}) \times \sum_{p \in S(X) \setminus \mathcal{P}} \alpha_{\ell}^{3}\right] \end{split}$$

StopGrad[hw<sub>i</sub>,v·i<sup>3</sup>] ?

$$= {^{X}\alpha_{.6}(hw_{j}v\cdot i^{3} + E_{j,3-j}hw_{3-j}v\cdot i^{3}) \cdot StopGrad[hw_{j}v\cdot i^{3}] \cdot E}$$

$$[|S(X) \cap P| \cdot |S(X) \setminus P|]$$
2

`E[2]

Now, by denoting

$$C_0 = \frac{\mathbb{E}[|S(X) \cap \mathcal{P}| \cdot |\mathcal{S}(X) \setminus \mathcal{P}|]}{2}, \quad C_1 = \frac{\mathbb{E}\left[|S(X) \cap \mathcal{P}|^2\right]}{2}, \quad C_2 = P - |S(X)|,$$

$$\bar{B_{j,}}^3 = \mathsf{StopGrad}[hw_{j,}v \cdot \mathbf{i}^3], \quad B_{j,}^* = hw_{j,}v \cdot \mathbf{i}, \quad Q_j = (\mathsf{E}[\mathsf{StopGrad}[G^2_j(X^{(2)})]])^{-1/2}.$$

we denote  $U_i := E[F_i^2(X^{(1)})]$ , where the expanded expression is

$$U_j = E[F_{j2}(X_{(1)})] = X C_1\alpha^{6}(B_{j,3} + E_{j,3-j}B_{33-j,3})_2 + C_2E_{j,3-j}$$
  
 $\in [2]$ 

and we can rewrite the objective as follows

$$L(W, E) = 2 - \sum_{j \in [2]} \sum_{\ell \in [2]} \frac{Q_j C_0 \alpha_\ell^6 \bar{B}_{j,\ell}^3 (B_{j,\ell}^3 + E_{j,3-j} B_{3-j,\ell}^3)}{U_j^{1/2}}$$
(A.1)

Now denote

$$H_{j,`} = C_1 \alpha \cdot 6(B_{j,`3} + E_{j,3-j}B_{33-j,`})_2 + C_2 E_{j,3-j,}$$

$$K_{j,`} = C_1 \alpha \cdot 6(B_{j,`3} + E_{j,3-j}B_{33-j,`})(B_{j,33-`} + E_{j,3-j}B_{33-j,3-`})$$

It is easy to calculate

$$Q^{-_{j}2} = \mathsf{E}[\mathsf{StopGrad}[G^{2_{j}}(X^{(2)})]]$$

$$\begin{split} &= \mathbb{E}\left[\left(\sum_{p \in [P]} \sigma(\langle w_j, X_p^{(2)} \rangle)\right)^2\right] \\ &= \frac{1}{2} \sum_{\ell \in [2]} \alpha_\ell^6 \langle w_j, v_\ell \rangle^6 \mathbb{E}\left[|S(X) \cap \mathcal{P}|^2\right] + \mathbb{E}\left[|\mathcal{P} \setminus S(X)| \langle w_j, \xi_p \rangle^6\right] \\ &= \sum_{\ell \in [2]} C_1 \alpha_\ell^6 B_{j,\ell}^6 + C_2 \mathcal{E}_j \end{split}$$

where  $E_j = E[hw_{ji}\xi_p i^6]$ . And thus the gradient can be computed as (notice  $B_{ji}^{-3} = B_{ji}^{-3}$ )

$$\begin{split} -\nabla_{w_{j}}L(W,E) &= \sum_{\ell \in [2]} \left( \frac{C_{0}Q_{j}\alpha_{\ell}^{6}H_{j,3-\ell}B_{j,\ell}^{5}}{U_{j}^{3/2}} \right) v_{\ell} + \sum_{\ell \in [2]} \left( \frac{C_{0}Q_{3-j}E_{3-j,j}\alpha_{\ell}^{6}B_{3-j,\ell}^{3}B_{j,\ell}^{2}H_{3-j,3-\ell}}{U_{3-j}^{3/2}} \right) v_{\ell} \\ &- \sum_{\ell \in [2]} \left( \frac{C_{0}Q_{j}\alpha_{3-\ell}^{6}B_{j,3-\ell}^{3}B_{j,\ell}^{2}K_{j,\ell}}{U_{j}^{3/2}} + \frac{C_{0}Q_{3-j}E_{3-j,j}\alpha_{3-\ell}^{6}B_{3-j,3-\ell}^{3}B_{j,\ell}^{2}K_{3-j,\ell}}{U_{3-j}^{3/2}} \right) v_{\ell} \\ &- \sum_{j' \in [2]} \sum_{\ell \in [2]} \frac{C_{0}C_{2}Q_{j'}\alpha_{\ell}^{6}B_{j',\ell}^{3}(B_{j',\ell}^{3} + E_{j',3-j'}B_{3-j',\ell}^{3})}{U_{j'}^{3/2}} \nabla_{w_{j}}\mathcal{E}_{j',3-j'} \\ &= \sum_{\ell \in [2]} (\Lambda_{j,\ell} + \Gamma_{j,\ell} - \Upsilon_{j,\ell})v_{\ell} - \sum_{(j',\ell) \in [2] \times [2]} \Sigma_{j',\ell}\nabla_{w_{j}}\mathcal{E}_{j',3-j'} \end{aligned} \tag{A.2}$$

 $\nabla_{w_j} \mathbf{E}_{j,3-j} = 6 \mathbf{E} [\mathbf{h} w_j, \xi_p \mathbf{i} \mathbf{5} \xi_p + E_{j,3-j} \mathbf{h} w_j, \xi_p \mathbf{i} \mathbf{2} \mathbf{h} w_{3-j}, \xi_p \mathbf{i} \mathbf{3} \xi_p ]$ 

 $\nabla_{w_i} E_{3-j,j} = 6 E[E_{32-j,j} h w_i, \xi_{p_i} \xi_p + E_{3-j,j} h w_{3-j,} \xi_{p_i} i 3 h w_i, \xi_{p_i} i 2 \xi_p]$ 

As for the gradient of the prediction head, we can calculate

$$\begin{split} -\nabla_{E_{j,3-j}}L(W,E) &= \sum_{\ell \in [2]} \frac{C_0 Q_j \alpha_\ell^6 B_{j,\ell}^3 B_{3-j,\ell}^3 U_j}{U_j^{3/2}} \\ &- \sum_{\ell \in [2]} \frac{C_0 Q_j \alpha_\ell^6 B_{j,\ell}^3 (B_{j,\ell}^3 + E_{j,3-j} B_{3-j,\ell}^3) \sum_{\ell' \in [2]} C_1 \alpha_{\ell'}^6 (B_{j,\ell'}^3 + E_{j,3-j} B_{3-j,\ell'}^3) B_{3-j,\ell'}^3}{U_j^{3/2}} \\ &- \sum_{\ell \in [2]} \frac{C_0 C_2 Q_j \alpha_\ell^6 B_{j,\ell}^3 (B_{j,\ell}^3 + E_{j,3-j} B_{3-j,\ell}^3)}{U_j^{3/2}} \nabla_{E_{j,3-j}} \mathcal{E}_{j,3-j} \\ &= \sum_{\ell \in [2]} \frac{C_0 Q_j \alpha_\ell^6 B_{j,\ell}^3 (B_{3-j,\ell}^3 H_{j,3-\ell} - B_{3-j,3-\ell}^3 K_{j,3-\ell})}{U_j^{3/2}} \\ &- \sum_{\ell \in [2]} \sum_{j,\ell} \mathbb{E} \left[ 2 \langle w_j, \xi_p \rangle^3 \langle w_{3-j}, \xi_p \rangle^3 + 2 E_{j,3-j} \langle w_{3-j}, \xi_p \rangle^6 \right] \end{split}$$

where  $\Sigma_{j,\cdot}$  is defined in (A.2). In fact, all the above gradient expressions can be simplified by letting  $\Phi_j := Q_j/U_j$  for  $j \in [2]$ , which is what we shall do in later sections.

**Summarizing the notations.** We shall define some useful notations to simplify the proof. We define  $V = \text{span}(v_1, v_2)$ . Let  $\Pi_A$  be the projection operator to subspace  $A \subset \mathbb{R}^d$ , then

$$\overline{R}_{1,2} := h \Pi v_{\perp} w_{1}, w_{j} i \qquad \qquad \overline{R}_{1,2} := h \Pi v_{\perp} w_{1}, w_{2} i \qquad \qquad \overline{R}_{1,2} := \frac{\langle \Pi_{V^{\perp}} w_{1}, w_{2} \rangle}{\|\Pi_{V^{\perp}} w_{1}\|_{2} \|\Pi_{V^{\perp}} w_{2}\|_{2}}$$

#### A.2 Some Useful Bounds for Gradients

In this section we use the superscript  $^{(t)}$  to denote the Below we present  $\nabla_{w_j}\mathcal{E}_{j',3-j'}^{(t)}$  iteration t during training. claim which comes from

direct calculations of  $\Sigma_{i}$ , and, which is very useful in the following sections.

Claim A.1 (on  $\Sigma_{j}$  and.  $\nabla_{w_j} \mathcal{E}_{j',3-j'}^{(t)}$   $R_j, R_{1,2}^{(t)}$  Letbe defined as above, then we have

(a) 
$$\Sigma_{j,\ell}^{(t)} = O(\Sigma_{1,1}^{(t)}) \frac{(B_{j,\ell}^{(t)})^6 + E_{j,3-j}^{(t)} (B_{3-j,\ell}^{(t)})^3 (B_{j,\ell}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}}.$$

$$(b) \ \langle \nabla_{w_j} \mathcal{E}_{i,3-j}^{(t)}, \Pi_{V^\top} w_i^{(t)} \rangle = \Theta([R_i^{(t)}]^3) \pm \Theta(E_{i,3-j}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2}.$$

(c) 
$$\langle \nabla_{w_j} \mathcal{E}_{3-j,j}^{(t)}, w_j^{(t)} \rangle = \Theta((E_{3-j,j}^{(t)})^2) [R_j^{(t)}]^3 \pm O(E_{3-j,j}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2}$$

$$(d) \ \langle \nabla_{w_j} \mathcal{E}_{i,3-j}^{(t)}, w_{3-j}^{(t)} \rangle = (\Theta(\overline{R}_{1,2}^{(t)}) \pm \varrho) [R_j^{(t)}]^{5/2} [R_{3-j}^{(t)}]^{1/2} + O(E_{j,3-j}^{(t)}) R_j^{(t)} [R_{3-j}^{(t)}]^2.$$

(e) 
$$\langle \nabla_{w_j} \mathcal{E}_{3-j,j}^{(t)}, w_{3-j}^{(t)} \rangle = ((E_{3-j,j}^{(t)})^2 (\Theta(\overline{R}_{1,2}^{(t)}) \pm \varrho) [R_j^{(t)}]^{5/2} [R_{3-j}^{(t)}]^{1/2} + O(E_{3-j,j}^{(t)}) R_j^{(t)} [R_{3-j}^{(t)}]^2)$$

*Proof.* The part on  $\Sigma_{j,\cdot}$  is trivial from its expression, we shall focus on proving (b) – (d).  $\mathbf{On}(\nabla_{w_j}\mathcal{E}_{j',3-j'}^{(t)},w_j^{(t)})$ : If j=j', then

$$h\nabla_{w_j}E_{j,(t3)-j,w_j(t)}\mathbf{i} = \Theta(1)E[hw_j(t),\xi_{p}\mathbf{i}\mathbf{6} + E_{j,(t3)-j}hw_j(t),\xi_{p}\mathbf{i}\mathbf{3}hw_3(t-)_j,\xi_{p}\mathbf{i}\mathbf{3}]$$

$$=\Theta\big(1\big)\mathsf{E}\big[\mathsf{h} w j(t),\xi p \mathbf{i} 6\big] + O\big(Ej,(t3)-j\big)\mathsf{E}\big[\mathsf{h} w j(t),\xi p \mathbf{i} 3\big(\mathsf{h} w 3(t-)j,\xi p \mathbf{i} 3\big)\big]$$

$$- \left. \mathsf{h} \big( I - w^{\!\top} \! j_{,t} w^{\!\top} \! j_{,t} \! > \big) w_3(t-) j_{,} \xi_p \mathrm{i}_3 \big) \right]$$

$$+ O\big(E_{j,}(t3)-j\big) \mathsf{E}\big[\mathsf{h} w_j(t), \xi_p \mathsf{i} \mathsf{3} \mathsf{h}\big(I-w^-j,tw^-j,t>\big) w_3(t-)j, \xi_p \mathsf{i} \mathsf{3}\big]$$

Write 
$$w_{j,t}=rac{\Pi_{V^\perp}w_j^{(t)}}{\|\Pi_{V^\perp}w_j^{(t)}\|_2}$$
, we can derive

$$\begin{split} & \mathbb{E}[\langle w_{j}^{(t)}, \xi_{p} \rangle^{3} (\langle w_{3-j}^{(t)}, \xi_{p} \rangle^{3} - \langle (I - \bar{w}_{j,t} \bar{w}_{j,t}^{\top}) w_{3-j}^{(t)}, \xi_{p} \rangle^{3})] \\ &= \mathbb{E}[\langle w_{j}^{(t)}, \xi_{p} \rangle^{3} \langle \bar{w}_{j,t} \bar{w}_{j,t}^{\top} w_{3-j}^{(t)}, \xi_{p} \rangle O(\langle w_{3-j}^{(t)}, \xi_{p} \rangle^{2})] \\ &= O(\frac{R_{1,2}^{(t)}}{\|\Pi_{V^{\perp}} w_{j}^{(t)}\|_{2}^{2}}) \mathbb{E}[\langle w_{j}^{(t)}, \xi_{p} \rangle^{4} \langle w_{3-j}^{(t)}, \xi_{p} \rangle^{2}] \\ &\leq O(\frac{R_{1,2}^{(t)}}{\|\Pi_{V^{\perp}} w_{j}^{(t)}\|_{2}^{2}}) \mathbb{E}[\langle w_{j}^{(t)}, \xi_{p} \rangle^{6}]^{\frac{2}{3}} \mathbb{E}[\langle w_{3-j}^{(t)}, \xi_{p} \rangle^{6}]^{\frac{1}{3}} \\ &\leq O(\overline{R}_{1,2}^{(t)}) \|\Pi_{V^{\perp}} w_{j}^{(t)}\|_{2}^{3} \|\Pi_{V^{\perp}} w_{3-j}^{(t)}\|_{2}^{3} \end{split} \tag{by H`older's inequality)}$$

and by our assumption on noise  $\xi_p$ , we also have

$$\mathbb{E}[\langle w_j^{(t)}, \xi_p \rangle^3 \langle (I - \bar{w}_{j,t} \bar{w}_{j,t}^\top) w_{3-j}^{(t)}, \xi_p \rangle^3] \le O(\varrho) \|\Pi_{V^\perp} w_j^{(t)}\|_2^3 \|\Pi_{V^\perp} w_{3-j}^{(t)}\|_2^3$$

Combined with the fact that 
$$\mathbb{E}[\langle w_j^{(t)}, \xi_p \rangle^6] = O(\|\Pi_{V^\perp} w_j^{(t)}\|_2^3)$$
, we can get  $\langle \nabla_{w_j} \mathcal{E}_{j,3-j}^{(t)}, w_j^{(t)} \rangle = O(\|\Pi_{V^\perp} w_j^{(t)}\|_2^6) \pm O(E_{j,3-j}^{(t)}) (R_{1,2}^{(t)} + \varrho) \|\Pi_{V^\perp} w_j^{(t)}\|_2^3 \|\Pi_{V^\perp} w_{3-j}^{(t)}\|_2^3$ 

when  $j^0 = 3 - j$ , we also have

 $h\nabla_{w_j}E_3(-t)_{j,j},w_j(t)i = \Theta(1)E[(E_3(t-)_{j,j})_2hw_j(t),\xi_pi_6 + E_3(t-)_{j,j}hw_j(t),\xi_pi_3hw_3(t-)_{j,\xi_pi_3}]$ 

$$=O((E_{3(t-)j,j})_2)k\Pi v_{\perp}w_j(t)k62\pm O(E_{3(t-)j,j})(R_{1(t,)2}+\%)k\Pi v_{\perp}w_j(t)k23k\Pi v_{\perp}w_3(t-)_jk23$$

$$\mathbf{On}^{\langle \nabla_{w_j} \mathcal{E}_{j',3-j'}^{(t)}, w_{3-j}^{(t)} \rangle}$$
; when  $i^0 = i$ , we have

 $h\nabla_{w_{j}}E_{j,(t3)-j,W3(t-)j}i = O(1)E[hw_{j}(t),\xi_{p}i5hw_{3}(t-)j,\xi_{p}i + E_{j,(t3)-j}hw_{j}(t),\xi_{p}i2hw_{3}(t-)j,\xi_{p}i4]$ 

$$= O(1) E[hw_j(t), \xi_p i sh(I - w_j, tw_j, t > + w_j, tw_j, t > w_j, t > w$$

+  $O(1)E[E_{i,(t3)-i}hw_{i(t)},\xi_{pi2}hw_{3(t-)i},\xi_{pi4}]$ 

Using H"older's inequality and our assumpsion on  $\xi_p$ , we have

$$\mathbb{E}[\langle w_i^{(t)}, \xi_p \rangle^5 \langle (I - \bar{w}_{j,t} \bar{w}_{j,t}^\top) w_{3-j}^{(t)}, \xi_p \rangle] \lesssim \varrho \|\Pi_{V^\perp} w_j^{(t)}\|_2^5 \|\Pi_{V^\perp} w_{3-j}^{(t)}\|_2$$

In the meantime, we also have

$$\mathbb{E}[\langle w_i^{(t)}, \xi_p \rangle^5 \langle \bar{w}_{j,t} \bar{w}_{i,t}^\top w_{3-j}^{(t)}, \xi_p \rangle] = \Theta(\overline{R}_{1,2}^{(t)}) \mathbb{E}[\langle w_i^{(t)}, \xi_p \rangle^6] [R_i^{(t)}]^{-1/2} [R_{3-j}^{(t)}]^{1/2} = \Theta(\overline{R}_{1,2}^{(t)}) [R_i^{(t)}]^{5/2} [R_{3-j}^{(t)}]^{1/2}$$

for the last term in (A.3), we can also use H"older's inequality to get

$$E_{j,3-j}^{(t)}\mathbb{E}[\langle w_j^{(t)}, \xi_p \rangle^2 \langle w_{3-j}^{(t)}, \xi_p \rangle^4] \lesssim E_{j,3-j}^{(t)}\mathbb{E}[\langle w_j^{(t)}, \xi_p \rangle^6]^{1/3}\mathbb{E}[\langle w_{3-j}^{(t)}, \xi_p \rangle^6]^{2/3} \lesssim E_{j,3-j}^{(t)}R_j^{(t)}[R_{3-j}^{(t)}]^2$$

Therefore, we can combine above analysis to get

$$\langle \nabla_{w_j} \mathcal{E}_{i,3-j}^{(t)}, w_{3-j}^{(t)} \rangle = (\Theta(\overline{R}_{1,2}^{(t)}) \pm \varrho) [R_i^{(t)}]^{5/2} [R_{3-j}^{(t)}]^{1/2} + O(E_{i,3-j}^{(t)}) R_i^{(t)} [R_{3-j}^{(t)}]^2$$

When  $j^0 = 3 - j$ , we also have

$$\langle \nabla_{w_j} \mathcal{E}_{3-j,j}^{(t)}, w_{3-j}^{(t)} \rangle = 6\mathbb{E}[(E_{3-j,j}^{(t)})^2 \langle w_j^{(t)}, \xi_p \rangle^5 \langle w_{3-j}^{(t)}, \xi_p \rangle + E_{3-j,j}^{(t)} \langle w_j^{(t)}, \xi_p \rangle^2 \langle w_{3-j}^{(t)}, \xi_p \rangle^4]$$

$$= 6(E_{3-j,j}^{(t)})^2 (\Theta(\overline{R}_{1,2}^{(t)}) \pm \varrho) [R_j^{(t)}]^{5/2} [R_{3-j}^{(t)}]^{1/2} + E_{3-j,j}^{(t)} R_j^{(t)} [R_{3-j}^{(t)}]^2$$

which proves the claim.

# **B** Phase I: Learning the Stronger Feature

In this section, we shall discuss the initial phase of learning the stronger feature. Firstly, we establish some properties at the initialization for our induction afterwards.

**Initialization properties.** We prove the following properties for our network at initialization.

 $(0) = I_2.$ 

Recall our initialization is  $w_j \sim N(0, I_d/d), \forall j \in [2]$  and E

**Lemma B.1** (properties at initialization). Recall that without loss of generality we let  $|B_{1,1}^{(0)}| = \max_{j \in [2]} |B_{j,1}^{(0)}|$ . With probability 1 - o(1), the following holds:

$$(a) \ \|w_j^{(0)}\|_2^2 = 1 \pm \widetilde{O}(\frac{1}{\sqrt{d}}) \text{ for all } j \in [2], \text{ and } |\langle w_1^{(0)}, w_2^{(0)} \rangle| \leq \widetilde{O}(\frac{1}{\sqrt{d}}) ;$$

(b) 
$$\max_{j,\cdot} |B_{j,\cdot}^{(0)}| \le O(p \log d/d)$$
 and  $\min_{j,\ell} |B_{j,\ell}^{(0)}| \ge \Omega(\frac{1}{\log d}) \max_{j,\ell} |B_{j,\ell}^{(0)}|$ ,

(c) 
$$|B_{1,1}^{(0)}| \ge |B_{2,1}^{(0)}|(1 + \frac{1}{\log d})$$

(d) 
$$\mathcal{E}_{j}^{(0)} = (1 - O(\frac{1}{d^3}))\sigma^6 \|w_j^{(0)}\|_2^6 = \Theta(1)$$
 for all  $j \in [2]$ ;

(e) 
$$H_{j,\ell}^{(0)} = C_2 \mathcal{E}_j^{(0)} (1 + \widetilde{O}(\frac{1}{\sqrt{d}}))$$
 for all (j,")  $\in$  [2]  $\times$  [2];

$$(f) \ \ U_{j}^{(0)} = C_{2} \mathcal{E}_{j}^{(0)} (1 + \widetilde{O}(\frac{\alpha_{1}^{6}}{\sqrt{d}})) \quad \textit{for} \quad \textit{all} \quad j \quad \in \quad \text{[2]};$$

$$(g) \ \ (Q_j^{(0)})^{-2} = C_2 \mathcal{E}_j^{(0)} (1 + \widetilde{O}(\frac{\alpha_1^6}{\sqrt{d}})) \text{ for all } j \in \textbf{[2]};$$

(0) 
$$6/d^3$$
 for all  $(j, `) \in [2] \times [2]$ .

(h)  $K_{i,\cdot} \leq O_{\mathbf{e}}(\alpha)$ 

Let us first introduce a fact about Gaussian ratio distribution without proof.

**Fact B.2** (Gaussian ratio distribution). If X and Y are two independent standard Gaussian variables, then the probability density of Z = X/Y is  $p(z) = \frac{1}{\pi(1+z^2)}, z \in (-\infty, \infty)$ .

Proof of Lemma B.1. a. Norm bound comes from simple  $\chi^2$  concentration inequality and our initialization  $w_j^{(0)} \sim \mathcal{N}(0, \frac{I_d}{d})$ . The inner product bound comes from Gaussian concentration.

b. It is from a direct calculation under our initialization, and some application of Gaussian c.d.f. and a union bound.

- c. It is from a probability distribution of Gaussian ratio distribution from Fact B.2 to bound the probability of  $|B_{1,1}^{(0)}|/|B_{2,1}^{(0)}| \leq (1+\frac{1}{\log d})$  (WLOG we let  $|B_{1,1}^{(0)}| = \max_{j \in [2]} |B_{j,1}^{(0)}|$ ).
- d. It can be directly proven from our assumption on noise  $\xi_p$  in the subspace  $V^{\perp}$  and (a).

(0) – 
$$\underline{\mathbf{1}}$$
  $E_{j,3-j}^{(0)}$  \_\_\_\_\_\_  $\in$  [2] and = 0, it is easy to directly

- e. Since at the initialization we have  $B_{j,} = O_{\mathbf{e}}(\sqrt{d}), j$ , upper bound the errors.
- f. Again from  $B_{j,\ell}^{(0)} = \widetilde{O}(\frac{1}{\sqrt{d}}), \forall j,\ell \in$  [2] at initialization and a direct upper bound.
- g. Proof is similar to (e).
- h. Directly from a naive upper bound using (b).

#### **B.1** Induction in Phase I

We define phase I as all iterations  $t \le T_1$ , where  $T_1 := \min\{t : B_{1,1}^{(t)} \ge 0.01\}$ , we will prove the existence of  $T_1$  at the end of this section. We state the following induction hypotheses, which will hold throughout the phase I:

**Inductions B.3.** For each  $t \le T_1$ , all of the followings hold:

(a). 
$$\|w_j^{(t)}\|_2 = \|w_j^{(0)}\|_2 \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$
 for each  $j \in [2]$ ;

(b). 
$$|B_{1,2}^{(t)}|, |B_{2,1}^{(t)}|, |B_{2,2}^{(t)}| = \widetilde{\Theta}(\frac{1}{\sqrt{d}})$$

(c). 
$$|B_{1,1}^{(t)}| \ge \Omega(\frac{1}{\log d}) \max(|B_{1,2}^{(t)}|, |B_{2,2}^{(t)}|, |B_{2,1}^{(t)}|)$$

$$(d). \ |E_{1,2}^{(t)}| \leq \widetilde{O}(\varrho + \tfrac{1}{\sqrt{d}}) \tfrac{\eta_E}{\eta} |B_{1,1}^{(t)}|_{\textit{and}} |E_{2,1}^{(t)}| \leq \widetilde{O}(\tfrac{1}{d}).$$

(e). 
$$R_1^{(t)}, R_2^{(t)} = \Theta(1), |R_{1,2}^{(t)}| \le \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

*Remark* B.4. Since we have chosen  $\eta_E \leq \eta$  and  $\varrho \leq \frac{1}{d^{\Omega(1)}}$ , Induction B.3d implies  $|E_{j,3-j}^{(t)}| = o(1)$  throughout  $t \leq T_1$ .

We shall prove the above induction holds in later sections, but first we need some useful claims assuming our induction holds in this phase.

# B.2 Computing Variables at Phase I

Firstly we establish a claim controlling the noise terms  $E_{j}$ ,  $E_{j,3-j}$  during this phase. **Claim B.5.** 

At each iteration  $t \le T_1$ , if Induction B.3 holds, then

(a) 
$$\mathcal{E}_1^{(t)} = \mathcal{E}_2^{(t)} \pm O(\sum_{\ell \in [2]} |B_{j,\ell}^{(t)}| + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}))$$

(b) 
$$\mathcal{E}_{j}^{(t)} = \mathcal{E}_{j}^{(0)} \pm O(\sum_{\ell \in [2]} |B_{j,\ell}^{(t)}| + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}))$$

(c) 
$$\mathcal{E}_{j,3-j}^{(t)} = \mathcal{E}_{j}^{(t)} \pm \widetilde{O}(E_{j,3-j}^{(t)}(\varrho + \frac{1}{\sqrt{d}}) + (E_{j,3-j}^{(t)})^{2})$$

Proof. For (a), we can simply write down

$$\mathcal{E}_{j}^{(t)} = \mathbb{E}[\langle w_{j}, \xi_{p} \rangle^{6}] = \sigma^{6} \|\Pi_{V^{\perp}} w_{j}^{(t)}\|_{2}^{6}$$

Note that by Induction B.3a we always have  $\|w_j^{(t)}\|_2 = \|w_j^{(0)}\|_2 \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$ , and by Lemma B.1a we also have  $\|w_j^{(0)}\|_2 = (1 \pm \widetilde{O}(\frac{1}{\sqrt{d}}))\|w_j^{(0)}\|_2$ , which implies

$$\begin{split} \|\Pi_{V^{\perp}}w_{j}^{(t)}\|_{2} - \|\Pi_{V^{\perp}}w_{3-j}^{(t)}\|_{2} &= \|w_{j}^{(t)}\|_{2} - \|w_{3-j}^{(t)}\|_{2} \pm O(\sum_{j,\ell \in [2]^{2}} B_{j,\ell}^{(t)}) \\ &= \|w_{j}^{(0)}\|_{2} - \|w_{3-j}^{(0)}\|_{2} \pm O(\sum_{j,\ell \in [2]^{2}} B_{j,\ell}^{(t)}) \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \\ &= \widetilde{O}(\frac{1}{\sqrt{d}}) \pm O(\sum_{j,\ell \in [2]^{2}} B_{j,\ell}^{(t)}) \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \end{split}$$

By the elementary equality  $x^n - y^n = (x - y)^{\mathbf{P}}_{0 \le i \le n-1} x^i y^{n-1-i}$ , we can obtain (a). The proof of (b) is almost the same as (a), and the proof of (c) is just direct calculation.  $\Box$ 

Equipped with Claim B.5, we can establish the following lemma, which will be frequently applied to bound the gradient in our induction argument.

**Lemma B.6** (variables control in phase I). Suppose Induction B.3 holds at some iteration  $t \le T_1$ , then we have:

$$(a) \ \ \textit{if} \ \forall \ell \in [2], \alpha_{\ell} | B_{j,\ell}^{(t)} | \leq O(1), \\ \textit{then} \ \Phi_{j}^{(t)} = (C_2 \mathcal{E}_{j}^{(t)})^{-2} (1 \pm_{\text{polylog}} 1 \pm_{\text{poly$$

$$\begin{array}{ll} \textit{(b)} \;\; \textit{if} \; \exists \ell \in [2], |B_{j,\ell}^{(t)}| \geq \Omega(\frac{1}{\alpha_{\ell}}), \\ \textit{then} \; \Phi_{j}^{(t)} = O((C_{2}\mathcal{E}_{j}^{(t)} + \sum_{\ell \in [2]} C_{1}\alpha_{\ell}^{6}(B_{j,\ell}^{(t)})^{6})^{-2}), \\ \textit{(c)} \;\; \textit{if} \; \alpha_{\ell}|B_{j,\ell}^{(t)}| \leq O(1), \; H_{j,\ell}^{(t)} = C_{2}\mathcal{E}_{j}^{(t)}(1 + \\ & \text{polylog}_{-----}(\textit{d})) = \Theta(\textit{C2}), \\ \textit{otherwise} \end{array}$$

(d) 
$$|K_{i,\ell}^{(t)}| \leq \widetilde{O}(\alpha_{\ell}^6/d^{3/2})$$

 $\textit{Proof.} \qquad \text{(a) From our assumptions that } |B_{1,2}^{(t)}|, |B_{2,1}^{(t)}|, |B_{2,2}^{(t)}| \leq \widetilde{O}(\frac{1}{\sqrt{d}}) \text{ and } \quad \alpha_1 B_{1,1}^{(t)} \leq O_{\text{(1)}}, \text{ and also the fact that } \mathcal{E}_j^{(t)} = \Omega(\sigma^6) = \Omega(1), \ C_2 = \Theta(\operatorname{polylog}(d)) \gg C_1, \text{ we can calculate }$ 

$$=O(C_{1})+C_{2}\mathcal{E}_{j}^{(t)}+\widetilde{O}(\varrho+\frac{1}{\sqrt{d}}) \quad U_{j(t)}=\text{X }C_{1}\alpha'_{6}((B_{j,\ (t)})_{3}+\\ (t) \qquad \qquad E_{j,(t3)-j}(B_{3(t-1)j,\ )}_{3})_{2}+C_{2}\text{E}_{j,(t3)-j}$$

$$=C_{2}\text{E}_{j}\left(1\pm\operatorname{polylog}(d)\right)$$

Meanwhile, we can also compute similarly

$$Q_{j(t)} = X C_{1}\alpha'_{6}(B_{j,'(t)})_{6} + C_{2}E_{j} = C_{2}E_{j(t)}(1 \pm \text{polylog}_{1} 1 (d))$$
' $\in [2]$ 

Therefore  $\Phi_j^{(t)} = Q_j^{(t)}/(U_j^{(t)})^{3/2} = (C_2 \mathcal{E}_j^{(t)} (1 \pm_{\text{polylog}} 1 \text{ (d)}))$ -2 as desired.

- (b) The proof is similar to that of (a).
- (c) when  $\alpha_1 B_{1,1}^{(t)} \leq O$  (1), the proof is similar to (a). When  $\alpha_1 B_{1,1}^{(t)} \geq O$  (1), we have from Induction

B.3a and  $H_{j,\ell}^{(t)}$ 's expression that

$$H_{j,\ell}^{(t)} = C_1 \alpha_{\ell}^6 ((B_{j,\ell}^{(t)})^3 + E_{j,3-j}^{(t)} (B_{3-j,\ell}^{(t)})^3)^2 + C_2 \mathcal{E}_{j,3-j}^{(t)} \le \widetilde{O}(\alpha_{\ell}^6)$$

And since  $T_1 := \min\{t: B_{1,1}^{(t)} \geq 0.01\}$ , so for  $t \leq T_1$ , we have

$$C$$
  $\mathcal{E}^{(t)}$ , so for  $t \leq T_1$ , we have 
$$C \mathcal{E}^{(t)} - |E_{j,3-j}^{(t)}| \geq \Omega(C_2)$$

$$H_{j,1} \geq C_2 E_{j,3-j} \geq 2_j$$

where ¬ is from Claim B.5b and - is from Induction B.3d.

(d) Since we have assumed  $|B_{1,2}^{(t)}|, |B_{2,1}^{(t)}|, |B_{2,2}^{(t)}| \leq \widetilde{O}(\frac{1}{\sqrt{d}})$ , it is direct to bound  $|K_{j,\ell}^{(t)}| \leq \widetilde{O}(\alpha_{\ell}^6/d^{1.5})$ .

Claim B.7 (about  $\Sigma_{j}$  and. If  $\nabla_{w_{j}} \mathcal{E}_{j',3-j'}^{(t)}$  Induction B.3 holds at iteration  $t \leq T_{1}$ , then

$$(a) \ \ \Sigma_{j,\ell}^{(t)} = O(\Lambda_{1,1}^{(t)}B_{1,1}^{(t)}) \frac{(B_{j,\ell}^{(t)})^6 + E_{j,3-j}^{(t)}(B_{3-j,\ell}^{(t)})^3(B_{j,\ell}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}}.$$

(b) 
$$\langle \nabla_{w_j} \mathcal{E}_{j,3-j}^{(t)}, w_j^{(t)} \rangle = O(1) \pm O(E_{j,3-j}^{(t)}) (R_{1,2}^{(t)} + \varrho),$$

(c) 
$$\langle \nabla_{w_j} \mathcal{E}_{3-j,j}^{(t)}, w_j^{(t)} \rangle = O((E_{3-j,j}^{(t)})^2) \pm O(E_{3-j,j}^{(t)})(R_{1,2}^{(t)} + \varrho)$$

(d) 
$$|\langle \nabla_{w_j} \mathcal{E}_{j,3-j}^{(t)}, w_{3-j}^{(t)} \rangle| = O(R_{1,2}^{(t)} + \varrho) + O(E_{j,3-j}^{(t)}),$$

(e) 
$$|\langle \nabla_{w_j} \mathcal{E}_{3-j,j}^{(t)}, w_{3-j}^{(t)} \rangle| = O(R_{1,2}^{(t)} + \varrho)(E_{3-j,j}^{(t)})^2 + O(E_{3-j,j}^{(t)})$$

*Proof.* Notice that  $\|\Pi_{V^{\perp}}w_j^{(t)}\|_2 = \Theta(1), \forall j \in [2] \text{ for } t \leq T_1$ , which is because of  $\|w_j^{(t)}\|_2 = \sqrt{2} \pm o(1)$  from Induction B.3a and  $\max_{j,\ell} |B_{j,\ell}^{(t)}| < 0.02^3$ . Now we can apply Claim A.1 to obtain the bounds.  $\square$ 

due to our choice of s small, we can make sure when  $T_1 = \min\{t: B_{1,1}^{(t)} \geq 0.01\}, \ B_{1,1}^{(T_1)} < 0.02$ .

 $_{3}\,_{3}\eta=\frac{1}{\mathsf{poly}(d)}$ 

#### **B.3** Gradient Lemmas for Phase I

We first present an interesting lemma regarding the effects of Batch-Normalization on the gradients of weights. The following lemma allow us maintain the norm of weights to above a constant throughout phase I.

**Lemma B.8** (effects of BN on gradients). For any  $W = (w_1, w_2)$  and E, it holds

(a) 
$$P_{i \in [2]} h \nabla_{wi} L(W,E), w_i i = 0;$$

Further, if Induction B.3 holds for each  $t \le T_1$ , we have

(b) 
$$|\langle \nabla_{w_j} L(W^{(t)}, E^{(t)}), w_j^{(t)} \rangle| \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) |\Lambda_{1,1}| \sum_{j \in [2]} |E_{j,3-j}^{(t)}|$$
 for each  $j \in [2]$ .

*Proof.* **Proof of (a):** We first calculate the gradient term as follows:

$$\begin{split} \nabla_W & L(\textit{W,E}) &= \nabla_W \quad j \texttt{X} \in [2] \quad q \quad \\ &= \texttt{EE}[F[F_{j2j}((\textit{XX}_{(1)}^{(1)})]) \textbf{q} \cdot \texttt{StopGradE}[\texttt{StopGrad}[G_j[(\textit{GX}_{2j}^{(2)}(\textit{X})]]_{(2)})]] \\ &= \sum_{j \in [2]} \frac{\mathbb{E}[(\nabla_W F_j(X^{(1)})) \cdot [G(X^{(2)})]_j] \cdot \mathbb{E}[F_j^2(X^{(1)})]}{(\mathbb{E}[F_j^2(X^{(1)})])^{3/2} \sqrt{\mathbb{E}[G_j^2(X^{(2)})]}} \\ &- \sum_{j \in [2]} \frac{\mathbb{E}[(\nabla_W F_j(X^{(1)})) \cdot F_j(X^{(1)})] \cdot \mathbb{E}[[F_j(X^{(1)}) \cdot [G(X^{(2)})]_j]}{(\mathbb{E}[F_j^2(X^{(1)})])^{3/2} \sqrt{\mathbb{E}[G_j^2(X^{(2)})]}} \end{split}$$

Since by our definition  $h\nabla_W F_j(X^{(1)}), Wi = P_{i \in [2]} h\nabla_{wi} [F_j(X^{(1)}), w_i = 3[F_j(X^{(1)}), w_i = m_{i \in [2]} h\nabla_{wi} L(W,E), w_i = 0.$ 

**Proof of (b):** Firstly we define a new notion

$$\nabla_{i,j} = \nabla_{w_i} \mathbf{q} \underline{\hspace{1cm}}$$

$$\mathsf{EE}[F[F2j((XX_{(1)(1)})])\mathbf{q} \cdot \mathsf{StopGradE}[\mathsf{StopGrad}[Gj[(GX2j(2)(X))]](2))]]$$

$$j$$

Then it is straightforward to verify that  $P_{i\in[2]}h\nabla_{i,j},w_i$  i = 0 for any  $j\in[2]$ , which implies that  $|h\nabla_{j^0,j},w_j$  i| =  $|h\nabla_{3-j^0,j},w_{3-j^0}|$ . So in order to obtain an upper bound for  $|h\nabla_{w_j}L(W,E),w_j| = |P_{j^0\in[2]}h\nabla_{j,j^0},w_j|$ , we only need to upper bound  $|h\nabla_{j,j^0},w_{3-j^0}|$ , each of which can be calculated as (ignoring all time superscript  $P_{j^0\in[2]}h\nabla_{j,j^0}$ )

$$\begin{split} |\langle \nabla_{3-j,j}, w_{3-j} \rangle| &= \frac{\mathbb{E}\left[\sum_{p \in [P] \cap \mathcal{P}} E_{j,3-j} \sigma(\langle w_{3-j}, X_p \rangle) \cdot [G(X^{(2)})]_j\right] \cdot \mathbb{E}[F_j^2(X^{(1)})]}{(\mathbb{E}[F_j^2(X^{(1)})])^{3/2} \sqrt{\mathbb{E}[G_j^2(X^{(2)})]}} \\ &- \frac{\mathbb{E}\left[\sum_{p \in [P] \cap \mathcal{P}} E_{j,3-j} \sigma(\langle w_{3-j}, X_p \rangle) \cdot F_j(X^{(1)})\right] \cdot \mathbb{E}[[F_j(X^{(1)}) \cdot [G(X^{(2)})]_j]}{(\mathbb{E}[F_j^2(X^{(1)})])^{3/2} \sqrt{\mathbb{E}[G_j^2(X^{(2)})]}} \end{split}$$

Now we compute

$$\mathbb{E}\left[\sum_{p\in[P]\cap\mathcal{P}}E_{j,3-j}\sigma(\langle w_{3-j},X_p\rangle)[G(X^{(2)})]_j\right] = \mathbb{E}\left[\sum_{p\in[P]\cap\mathcal{P}}E_{j,3-j}\sigma(\langle w_{3-j},X_p\rangle)\sum_{p\in[P]\setminus\mathcal{P}}\sigma(\langle w_j,X_p\rangle)\right] = X E_{j,3-j}C_0\alpha G_{33-j,B_{j,3}}$$

$$\stackrel{\cdot}{\in}[2]$$

?

and

?

E X 
$$hw_3$$
- $j$ , $Xpi$ ) ·  $Fj(X_{(1)})$   $\mathbb{Z}$   $E_{j,3}$ - $j\sigma(p \in [P] \cap P$ 

$$=\mathsf{E} \, \boxed{ \mathbb{Z} \, X \, E_{j,3-j} \sigma \big( hw_{3-j}, X_{p} \mathrm{i} \big) \cdot X \, \big( \sigma \big( hw_{j}, X_{p} \mathrm{i} \big) + E_{j,3-j} \sigma \big( hw_{3-j}, X_{p} \mathrm{i} \big) \big) } \boxed{ \mathbb{Z} }$$

$$p \in [P] \cap \mathsf{P}$$

$$= X E_{j,3-j}C_1\alpha `6B_{33-j,`}(B_{j,`3} + E_{j,3-j}B_{33-j,`}) + C_2E_{j,3-j}\mathsf{E}[hw_{j,}\xi_{p}i3hw_{3-j,}\xi_{p}i3 + E_{j,3-j}hw_{3-j,}\xi_{p}i6] ``\in [2]$$

So we can further obtain the nominator in the expression of  $|h\nabla_{3-j,i}w_{3-j}i|$  as

E② X 
$$E_{j,3-j}\sigma(hw_{3-j,Xp}i) \cdot [G(X_{(2)})]j$$
②  $\stackrel{\cdot}{=} E[F_{j2}(X_{(1)})]$ 
 $p \in [P] \cap P$ 
②

$$- \mathsf{E} \mathbb{Z} \; X \; E_{j,3-j} \sigma \big( \mathsf{h} w_{3-j,X_P} \mathsf{i} \big) \cdot F_j \big( X_{(1)} \big) \mathbb{Z} \cdot \mathsf{E} \big[ \big[ F_j \big( X_{(1)} \big) \cdot \big[ G \big( X_{(2)} \big) \big]_j \big] \\ p \in [P] \cap \mathsf{P}$$

$$= \mathbb{Z}\mathbf{X} \ E_{j,3-j}C_0\alpha `6B_{33-j,`B_{j,`3}} \mathbb{Z} \cdot \mathbb{Z}\mathbf{X} \ C_1\alpha `6(B_{j,`3} + E_{j,3-j}B_{33-j,`})_2 + C_2\mathbf{E}_{j,3-j}\mathbb{Z}$$
 ` $\in$ [2]

- ②X 
$$E_{j,3-j}C$$
1α′6 $B$ 33- $j,$ ′( $B_{j,$ ′3 +  $E_{j,3-j}B$ 33- $j,$ ′)② · ②X  $C$ 0α′6 $B_{j,$ ′3 ( $B_{j,$ ′3 +  $E_{j,3-j}B$ 33- $j,$ ′)②   
`∈[2]   
②

$$-C2E_{j,3-j}\mathsf{E}\big[\mathsf{h}w_{j,}\xi_{p}\mathsf{i}3\mathsf{h}w_{3-j,}\xi_{p}\mathsf{i}3+E_{j,3-j}\mathsf{h}w_{3-j,}\xi_{p}\mathsf{i}6\big]\cdot \mathbb{Z}\mathsf{X}\;C0\alpha^{\mathsf{c}}6B_{j,\,3}\;(B_{j,\,3}+E_{j,3-j}B_{33-j,\,\mathbf{b}})\mathbb{Z}$$

$$\stackrel{\mathsf{c}}{=}[2]$$

= 
$$E_{j,3-j} \times C_0 \alpha \cdot 6B_{33-j, \cdot} (B_{j, \cdot 3} H_{j,3-\cdot} - B_{j,33-\cdot} K_{j,3-\cdot})$$
  
` $\in [2]$ 

$$-C2E_{j,3-j}\mathsf{E}\big[hw_{j,}\xi_{p}\mathrm{i}3hw_{3-j,}\xi_{p}\mathrm{i}3+E_{j,3-j}hw_{3-j,}\xi_{p}\mathrm{i}6\big]\cdot \mathbb{Z}X\ C0\alpha'6B_{j,'3}\left(B_{j,'3}+E_{j,3-j}B_{33-j,'}\right)\mathbb{Z}$$

$$`\in[2]$$

Now can sum over  $j^0 \in [2]$  to get

$$|\langle \nabla_{w_i} L(W, E), w_j \rangle|$$

$$\leq \sum_{j \in [2]} \sum_{\ell \in [2]} C_0 E_{j,3-j} \left| \Phi_j \alpha_\ell^6 B_{3-j,\ell}^3 B_{j,\ell}^3 H_{j,3-\ell} \right| + \sum_{j \in [2]} \sum_{\ell \in [2]} \left| C_0 E_{j,3-j} \Phi_j \alpha_\ell^3 B_{3-j,\ell}^3 B_{j,3-\ell}^3 K_{j,3-\ell} \right| \\ + \sum_{j \in [2]} \sum_{\ell \in [2]} \left| C_2 E_{j,3-j} \Phi_j \mathbb{E}[\langle w_j, \xi_p \rangle^3 \langle w_{3-j}, \xi_p \rangle^3 + E_{j,3-j} \langle w_{3-j}, \xi_p \rangle^6] C_0 \alpha_\ell^6 B_{j,\ell}^3 (B_{j,\ell}^3 + E_{j,3-j} B_{3-j,\ell}^3) \right|$$

Next we are going to bound each term, for the first term of LHS we have

$$\sum_{j \in [2]} \sum_{\ell \in [2]} \left| C_0 E_{j,3-j} \Phi_j \alpha_{\ell}^6 B_{3-j,\ell}^3 B_{j,\ell}^3 H_{j,3-\ell} \right| \leq \sum_{j \in [2]} \sum_{\ell \in [2]} \left| E_{j,3-j} \right| |\Lambda_{j,\ell}| \left| \frac{B_{3-j,\ell}^3 B_{j,\ell}^3}{B_{j,\ell}^2} \right| \\
\leq |\Lambda_{1,1}| \sum_{j \in [2]} |E_{j,3-j}| \left| \frac{B_{3-j,\ell}^3 B_{j,\ell}^3 \Phi_j}{B_{1,1}^5 \Phi_1} \right| \\
\leq \widetilde{O}(\frac{d^{o(1)}}{\sqrt{d}}) |\Lambda_{1,1}| \sum_{j \in [2]} |E_{j,3-j}|$$

where the last inequality is because

- By Lemma B.6a,b, we have  $\Phi_j^{(t)}/\Phi_1^{(t)} \leq O(\alpha_1^O(1)) \leq d^{o(1)}$  during  $t \leq T_1$ .
- $\bullet \ (B_{3-j,\ell}^{(t)})^3 (B_{j,\ell}^{(t)})^3 \leq \widetilde{O}(\tfrac{1}{\sqrt{d}}) (B_{1,1}^{(t)})^5 \text{ from Induction B.3b,c.}$

Similarly, we can also compute

$$\sum_{j \in [2]} \sum_{\ell \in [2]} \left| C_0 E_{j,3-j} \Phi_j \alpha_\ell^3 B_{3-j,\ell}^3 B_{j,3-\ell}^3 K_{j,3-\ell} \right| \leq \sum_{j \in [2]} \sum_{\ell \in [2]} E_{j,3-j} |\Lambda_{1,1}| \left| \frac{B_{3-j,\ell}^3 B_{j,3-\ell}^3 K_{j,3-\ell}}{B_{1,1}^5 H_{j,3-\ell}} \right| \\
\leq \widetilde{O}(\frac{d^{o(1)}}{d^2}) |\Lambda_{1,1}| \sum_{j \in [2]} |E_{j,3-j}|$$

and

$$\begin{split} & \sum_{j \in [2]} \sum_{\ell \in [2]} \left| C_2 E_{j,3-j} \Phi_j \mathbb{E}[\langle w_j, \xi_p \rangle^3 \langle w_{3-j}, \xi_p \rangle^3 + E_{j,3-j} \langle w_{3-j}, \xi_p \rangle^6] C_0 \alpha_\ell^6 B_{j,\ell}^3 (B_{j,\ell}^3 + E_{j,3-j} B_{3-j,\ell}^3) \right| \\ & \stackrel{@}{\leq} \sum_{j \in [2]} \sum_{\ell \in [2]} \left| E_{j,3-j} \Lambda_{j,\ell} \right| \left| \frac{B_{j,\ell}^3 + E_{j,3-j} B_{3-j,\ell}^3}{B_{j,\ell}^2} \right| \left| \mathbb{E}[\langle w_j, \xi_p \rangle^3 \langle w_{3-j}, \xi_p \rangle^3 + E_{j,3-j} \langle w_{3-j}, \xi_p \rangle^6] \right| \\ & \stackrel{@}{\leq} \sum_{j \in [2]} \sum_{\ell \in [2]} \left| E_{j,3-j} \Lambda_{j,\ell} \right| \left| \frac{B_{j,\ell}^3 + E_{j,3-j} B_{3-j,\ell}^3}{B_{j,\ell}^2} \right| \left( O(R_{1,2} + \varrho) + O(E_{j,3-j}) \right) \\ & \leq \widetilde{O}(R_{1,2} + \varrho) |\Lambda_{1,1}| \sum_{j \in [2]} |E_{j,3-j}| \end{split}$$

where  $\neg$  is due to Lemma B.6c, - is from the same calculation in Claim B.7 for  $E[hw_b\xi_pi^3hw_{3-b}\xi_pi^3]$  and Induction B.3a. Now combining the above and Induction B.3e together we have

$$|\langle \nabla_{w_j} L(W, E), w_j \rangle| \le \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) |\Lambda_{1,1}| \sum_{j \in [2]} |E_{j,3-j}|$$

which gives the desired bound.

Next we give a lemma characterizing the gradient of feature  $v_1$  in this phase.

**Lemma B.9** (learning feature  $v_1$  in phase I). For each  $t \le T_1$ , if Induction B.3 holds at iteration t, then using notations of (A.2), we have:

(a) 
$$\langle -\nabla_{w_1} L(W^{(t)}, E^{(t)}), v_1 \rangle = (1 \pm \widetilde{O}(\frac{1}{d})) \Lambda_{1,1}^{(t)}$$

(b) 
$$\langle -\nabla_{w_2}L(W^{(t)}, E^{(t)}), v_1 \rangle = (1 \pm O(\frac{1}{\sqrt{d}}))\Lambda_{2,1}^{(t)} + \Gamma_{2,1}^{(t)} \le (1 \pm O(\frac{1}{\sqrt{d}}))\Lambda_{2,1}^{(t)} \pm \frac{(B_{2,1}^{(t)})^2}{(B_{1,1}^{(t)})^2}E_{1,2}^{(t)}\Lambda_{1,1}^{(t)}$$

*Proof.* From (A.2), we write down the gradient formula for  $B_{j,1}^{(t)}$  as follows:  $\langle -\nabla_{w_j} L_{\mathcal{D}}(W^{(t)}, E^{(t)}), v_1 \rangle = \Lambda_{j,1}^{(t)} + \Gamma_{j,1}^{(t)} - \Upsilon_{j,1}^{(t)}$ 

$$\langle -\nabla_{w_j} L_{\mathcal{D}}(W^{(t)}, E^{(t)}), v_1 \rangle = \Lambda_{j,1}^{(t)} + \Gamma_{j,1}^{(t)} - \Upsilon_{j,1}^{(t)}$$

where (ignoring the superscript (t) for the RHS)

$$\begin{split} &\Lambda_{j,1}^{(t)} = C_0 \Phi_j \alpha_1^6 B_{j,1}^5 H_{j,2} \\ &\Gamma_{j,1}^{(t)} = C_0 \Phi_{3-j} E_{3-j,j} \alpha_1^6 B_{3-j,1}^3 B_{j,1}^2 H_{3-j,2} \\ &\Upsilon_{j,1}^{(t)} = C_0 \alpha_2^6 \left( \Phi_j B_{j,2}^3 B_{j,1}^2 K_{j,1} + \Phi_{3-j} E_{3-j,j} B_{3-j,2}^3 B_{j,1}^2 K_{3-j,1} \right) \end{split}$$

We first prove (a), and we deal with each term individually:

**Comparing**  $\Lambda_{1,1}^{(t)}$  and  $\Gamma_{1,1}^{(t)}$ : When  $t \le T_{1,1}$ , we have from Lemma B.6a that

Further, by Induction B.3b,c,d and our definition of stage 1, we know  $E_{1,2}^{(t)} \leq \widetilde{O}(\frac{1}{d})$ . Now from

Induction B.3b that  $B_{2,1}^{(t)} \leq \widetilde{O}(\frac{1}{\sqrt{d}})$ , together we have

$$\Gamma_{1,1}^{(t)} = C_0 \alpha_1^6 E_{2,1}^{(t)} \Phi_2^{(t)} H_{2,2}^{(t)} (B_{2,1}^{(t)})^3 (B_{1,1}^{(t)})^2 \le \widetilde{O}(\frac{1}{d}) C_0 \alpha_1^6 \Phi_1^{(t)} H_{1,2}^{(t)} (B_{1,1}^{(t)})^5 = \widetilde{O}(\frac{\Lambda_{1,1}^{(t)}}{d})^5$$

When  $t \in [T_{1,1}, T_1]$ , by Lemma B.6b we have

$$\Phi_1^{(t)}H_{1,2}^{(t)} \geq \Omega(\frac{C_2}{(C_1\alpha_1^6(B_{1,1}^{(t)})^6 + O(C_2))^2}) \geq \omega(\frac{1}{d^{0.1}}) \qquad \text{and} \ E_{2,1}^{(t)}\Phi_2^{(t)}H_{2,2}^{(t)} \leq \widetilde{O}(\frac{1}{d})$$

Now from our definition of stage 2, it holds that  $B_{1,1}^{(t)} \geq \Omega(\frac{1}{\alpha_1})$  while  $B_{2,1}^{(t)} \leq \widetilde{O}(\frac{1}{\sqrt{d}})$  by Induction B.3b, which gives

$$\Gamma_{1,1}^{(t)} = C_0 \alpha_1^6 E_{2,1}^{(t)} \Phi_2^{(t)} H_{2,2}^{(t)} (B_{2,1}^{(t)})^3 (B_{1,1}^{(t)})^2 \le \widetilde{O}(\frac{1}{d}) C_0 \alpha_1^6 \Phi_1^{(t)} H_{1,2}^{(t)} (B_{1,1}^{(t)})^5 = \widetilde{O}(\frac{\Lambda_{1,1}^{(t)}}{d})^5$$

 $\Lambda_{1,1}^{(t)}$   $\Upsilon_{1,1}^{(t)}$  Comparingand: Now consider  $\Upsilon_{1,1}$ , by Lemma B.6, we can follow the same analysis as above to get

$$\Phi_{j}^{(t)}K_{j,\ell}^{(t)} \leq \widetilde{O}(\frac{\alpha_{1}^{O(1)}}{d^{3/2}})\Phi_{1}^{(t)}H_{1,2}^{(t)}$$
 for any  $(j, `) \in [2] \times [2]$ 

Combined with  $E_{2,1}^{(t)} \leq o$ (1), we can derive

$$\begin{split} \Upsilon_{1,1}^{(t)} &= C_0 \alpha_2^6 \left( \Phi_1^{(t)} K_{1,1}^{(t)} (B_{1,2}^{(t)})^3 (B_{1,1}^{(t)})^2 + E_{1,2}^{(t)} \Phi_2^{(t)} K_{2,1}^{(t)} (B_{2,2}^{(t)})^3 (B_{1,1}^{(t)})^2 \right) \\ &\leq \widetilde{O}(\frac{\alpha_1^{O(1)} \alpha_2^6}{d^{3/2}}) C_0 \alpha_1^6 \Phi_1^{(t)} H_{1,2}^{(t)} (B_{1,1}^{(t)})^5 \\ &= \widetilde{O}(\frac{\Lambda_{1,1}^{(t)}}{d^{3/2 - o(1)}} \qquad \qquad \text{(since $C_1$ = $O_{\textbf{e}}(1)$ and $\alpha_1, \alpha_2$ = $d^{o(1)}$)} \end{split}$$

 $\Lambda_{2,1}^{(t)}$   $\Upsilon_{2,1}^{(t)}$  Comparingand: Till now (a) is proved, we can deal with (b) by only comparing  $\Lambda_{2,1}$  with

 $\Upsilon^{(t)}_{2,1}$ . Similar to the above arguments, we have by Induction B.3b we know  $K^{(t)}_{j,1} = \widetilde{O}(\frac{C_1\alpha_1^6}{d^{3/2}}), \forall j \in [2]$ , and thus

$$\Phi_{j}^{(t)}K_{j,\ell}^{(t)} \leq \widetilde{O}(\frac{\alpha_{1}^{6}}{d^{3/2}})\Phi_{2}^{(t)}H_{2,2}^{(t)} \qquad \qquad \text{for any (j, `)} \in [2] \times [2]$$

By Induction B.3e we know  $E_{1,2}^{(t)} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$ . Also, note that from Induction B.3b we have  $\widetilde{O}((B_{1,2}^{(t)})^3/d) \leq \widetilde{O}((B_{2,1}^{(t)})^5)$ , and thus

$$E_{1,2}^{(t)}\Phi_{1}^{(t)}K_{1,1}^{(t)}(B_{1,2}^{(t)})^{3}(B_{2,1}^{(t)})^{2} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})\widetilde{O}(\frac{\alpha_{1}^{6}}{d^{5/2}})\Phi_{2}^{(t)}H_{2,2}^{(t)}\widetilde{O}(B_{1,2}^{(t)})^{3} \leq O(\frac{1}{d^{3/2}})\Phi_{2}^{(t)}H_{2,2}^{(t)}(B_{2,1}^{(t)})^{5}$$

So together we have

$$\begin{split} |\Upsilon_{2,1}^{(t)}| &= |C_0 \alpha_2^6 \left( \Phi_2^{(t)} K_{2,1}^{(t)} (B_{2,2}^{(t)})^3 (B_{2,1}^{(t)})^2 + E_{2,1}^{(t)} \Phi_1^{(t)} K_{1,1}^{(t)} (B_{1,2}^{(t)})^3 (B_{2,1}^{(t)})^2 \right) | \\ &\leq O(\frac{1}{d^{3/2}}) C_0 \alpha_1^6 \Phi_2^{(t)} H_{2,2}^{(t)} |(B_{2,1}^{(t)})^5| \\ &= O(\frac{1}{d^{3/2}}) |\Lambda_{2,1}^{(t)}| \end{split}$$

Comparing  $\Gamma_{2,1}^{(t)}$  with  $\Lambda_{1,1}^{(t)}$ : It suffices to notice that

$$|\Gamma_{2,1}^{(t)}| \le |E_{1,2}^{(t)}|C_0\alpha_1^6\Phi_1^{(t)}H_{1,2}^{(t)}|B_{1,1}^{(t)}|^3(B_{2,1}^{(t)})^2 = \frac{(B_{2,1}^{(t)})^2}{(B_{1,1}^{(t)})^2}|E_{1,2}^{(t)}||\Lambda_{1,1}^{(t)}|$$

Combining the bounds for  $\Lambda_{2,1}$  and  $\Gamma_{2,1}$ , we obtain the proof of (b).

Then we can also calculate the gradients of feature  $v_2$  in this phase.

**Lemma B.10** (learning feature  $v_2$  in phase I). For each  $t \le T_1$ , if Induction B.3 holds at iteration t, then using notations of (A.2), we have for each  $j \in [2]$ :

$$\langle -\nabla_{w_j} L(W^{(t)}, E^{(t)}), v_2 \rangle = \left(1 \pm \widetilde{O}(\alpha_1^6) (E_{3-j,j}^{(t)} + (B_{j,1}^{(t)})^3)\right) \Lambda_{j,2}^{(t)}$$
(B.1)

*Proof.* Again as in the proof of Lemma B.9, we expand the notations: (ignoring the superscript <sup>(t)</sup> for the RHS)

$$\begin{split} &\Lambda_{j,2}^{(t)} = C_0 \alpha_2^6 \Phi_j H_{j,1} B_{j,2}^5 \\ &\Gamma_{j,2}^{(t)} = C_0 \alpha_2^6 \Phi_j E_{3-j,j} B_{3-j,2}^3 B_{j,2}^2 H_{3-j,1} \\ &\Upsilon_{j,2}^{(t)} = C_0 \alpha_1^6 \left( \Phi_j B_{j,1}^3 B_{j,2}^2 K_{j,2} + \Phi_{3-j} E_{3-j,j} B_{3-j,1}^3 B_{j,2}^2 K_{3-j,2} \right) \\ &\text{(t)} &\text{(t)} \end{split}$$

We first compare  $\Lambda_{j,2}$  and  $\Gamma_{j,2}$  as follows: Lemma B.6 we have

- $B_{3-j,2}^{(t)} \leq \widetilde{O}(B_{j,2}^{(t)})$  by Induction B.3b;
- From Lemma B.6a,b we can have  $\Phi^{(t)}_{3-j} \leq \widetilde{O}(\alpha_1^{O(1)})\Phi_j^{(t)}, \forall j \in [2].$

Together they imply:

 $Co\alpha_{26}E_{3(t-)j,j}\big(B_{3(t-)j,2}\big)_3\big(B_{j,(t2)}\big)_2\Phi_{(3t-)j}H_{3(t-)j,1}\leq Oe\big(\alpha_{10(1)}E_{3(t-)j,j}\big)Co\alpha_{26}\Phi_{(jt)}H_{j,(t2)}\big(B_{j,(t2)}\big)_5$ 

$$O(1)$$
 (t) (t)  
=  $Oe(\alpha_1 \quad E_{j,3-j})\Lambda_{j,2}$  (B.2)

Now we turn to compare 
$$\Lambda_{j,2}$$
 with  $\Upsilon_{j,2}$ . We split  $\Upsilon_{j,2}$  into two terms  $\Upsilon_{j,}$   $\stackrel{(t)}{\underset{j,2,1}{}} = C_0 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{j,2}^{(t)}, \quad \Upsilon_{j,2,2}^{(t)} = C_0 \alpha_1^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{3-j,2}^{(t)}$ 

(t)

For  $\Upsilon_{i,2,1}$ , we can calculate

$$\begin{split} \Upsilon_{j,2,1}^{(t)} &= C_0 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{j,2}^{(t)} \\ & - \frac{6}{} \\ &\leq \widetilde{O}(\alpha_1^6 (B_{j,1}^{(t)})^3) C_0 \alpha_2^6 \Phi_j^{(t)} H_{j,1}^{(t)} (B_{j,2}^{(t)})^5 \qquad (\widetilde{O}(\frac{C_1}{d^{3/2}}) \leq \widetilde{O}((B_{j,2}^{(t)})^3 \\ &= \widetilde{O}(\alpha_1^6 (B_{j,1}^{(t)})^3) \Lambda_{j,2}^{(t)} \\ &\leq \widetilde{O}(\frac{C_1 \alpha_2}{d^{3/2}}) (B_{j,1}^{(t)})^3 \cdot C_0 \alpha_1^6 \Phi_j^{(t)} H_{j,1}^{(t)} (B_{j,2}^{(t)})^2 \qquad (K_{j,\ell}^{(t)} \leq \widetilde{O}(\frac{C_1 \alpha_\ell^6}{d^{3/2}}) \text{ from Lemma B.6d)} \\ &\text{) from Induction B.3b) (B.3)} \end{split}$$

(t)

And for  $\Upsilon_{j,2,1}$ , we use Induction B.3b and Lemma B.6d again to get

$$(B_{3-j,1}^{(t)})^3 (B_{3-j,2}^{(t)})^2 K_{3-j,2}^{(t)} \le \widetilde{O}(C_1 \alpha_2^6 (B_{j,2}^{(t)})^5)$$

and thus combined with  $\Phi_{3-j}^{(t)} \leq \widetilde{O}(\alpha_1^6)\Phi_j^{(t)}, \forall j \in [2]$  from Lemma B.6a,b, we can derive

$$Y_{j,(t2),2} = C0\alpha 16\Phi(3t-)jE_{j,(t3)-j}(B_{3(t-)j,1})3(B_{j,(t2)})2K_{3(t-)j,2}$$

$$\leq Oe(\alpha 16E_{3(t-)j,j})C0\alpha 26\Phi(jt)H_{j,(t1)}(B_{j,(t2)})5$$

$$= Oe(\alpha_{16}E_{3(t-)j,j})\Lambda_{(j,t2)}$$
(B.4)

Now combine the results of (B.2), (B.3) and (B.4) finishes the proof of (B.1).

**Lemma B.11** (learning prediction head  $E_{1,2}, E_{2,1}$  in phase I). *If Induction B.3 holds at iteration*  $t \le T_1$ , *then we have* 

$$(a) -\nabla_{E_{1,2}}L(W^{(t)}, E^{(t)}) = O(\Lambda_{1,1}^{(t)}B_{1,1}^{(t)}) \left( -O(E_{1,2}^{(t)}) + \widetilde{O}(\frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^3}) + O(R_{1,2}^{(t)}) \right);$$

$$(b) -\nabla_{E_{2,1}}L(W^{(t)}, E^{(t)}) = \widetilde{O}(\frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^2})\Lambda_{1,1}^{(t)} + \sum_{\ell \in [2]} C_2\Lambda_{2,\ell}^{(t)}B_{2,\ell}^{(t)} \left( -O(E_{2,1}^{(t)}) + O(R_{1,2}^{(t)}) \right);$$

*Proof.* We first write down the gradient for  $E^{(t)}_{j,3-j}$ : (ignoring the time superscript  $E^{(t)}$ )

$$-\nabla E_{j,3-j}L(W,E) = X C_0 \Phi_{j}\alpha {}^{\circ}6B_{j}, {}^{\circ}3 \left(B_{33-j}, {}^{\circ}H_{j,3-} - B_{33-j,3-} {}^{\circ}K_{j,3-} \right) - X \Sigma_{j}, {}^{\circ}\nabla E_{j,3-j}E_{j,3-j}$$

$$\text{``}\in [2]$$

$$\text{where} \nabla E_{j,3-j} \mathcal{E}_{j,3-j} = \mathbb{E}\left[2\langle w_j, \xi_p \rangle^3 \langle w_{3-j}, \xi_p \rangle^3 + 2E_{j,3-j} \langle w_{3-j}, \xi_p \rangle^6\right]. \text{ Thus we have}$$

$$\nabla E_{j,3-j} \mathcal{E}_{j,3-j}^{(t)} = O(1)E_{j,3-j}^{(t)} + O(R_{1,2}^{(t)})$$

and by Claim B.5 and Lemma B.6a,b

$$\Sigma_{j,\ell}^{(t)} = O(\Lambda_{1,1}^{(t)} B_{1,1}^{(t)}) \frac{(B_{j,\ell}^{(t)})^6 + E_{j,3-j}^{(t)} (B_{3-j,\ell}^{(t)})^3 (B_{j,\ell}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}} \le O(\Lambda_{1,1}^{(t)} B_{1,1}^{(t)})$$

Now let us look at  $\nabla_{E_{1,2}}L(W^{(t)},E^{(t)})$ , first we consider the term

$$X C_0 \Phi_{(1t)} \alpha_{\cdot 6} (B_{1(t, ')})_3 ((B_{2(t, ')})_3 H_{1(t, 3)^{-}} - (B_{2(t, 3)^{-}})_3 K_{1(t, 3)^{-}})_{\cdot \in [2]}$$

Using Lemma B.6 and Induction B.3b,c, we know

 $\bullet \ H_{1,1}^{(t)} \leq \widetilde{O}(H_{1,2}^{(t)}) \text{ at } t \leq T_{1,1} \\ \mathsf{and} \\ H_{1,1}^{(t)} \leq \widetilde{O}(\alpha_1^6 H_{1,2}^{(t)}) \text{ for } t \in [T_{1,1}, T_1];$ 

• 
$$B_{2,1}^{(t)}, B_{1,2}^{(t)}, B_{2,2}^{(t)} \le \widetilde{O}(B_{2,1}^{(t)}) \le \widetilde{O}(B_{1,1}^{(t)})$$

$$K_{1,3-\ell}^{(t)} \leq \widetilde{O}(\alpha_1^6/d^{3/2}).$$

It can be computed that

$$\begin{split} C_0\Phi_1^{(t)}\alpha_2^6(B_{1,2}^{(t)})^3(B_{2,2}^{(t)})^3H_{1,1}^{(t)} &\leq \widetilde{O}(1)\left(\frac{B_{2,1}^{(t)}}{B_{1,1}^{(t)}}\right)^3C_0\Phi_1^{(t)}\alpha_1^3(B_{1,1}^{(t)})^6H_{1,2}^{(t)} \\ &\sum_{\ell\in[2]}\left|C_0\Phi_1^{(t)}\alpha_\ell^6(B_{1,\ell}^{(t)})^3(B_{2,\ell}^{(t)})^3K_{1,3-\ell}^{(t)}\right| &\leq \widetilde{O}(\frac{\alpha_1^6}{d^{3/2}})\frac{(B_{2,1}^{(t)})^3}{(B_{1,1}^{(t)})^3}C_0\Phi_1^{(t)}\alpha_1^6(B_{1,1}^{(t)})^6H_{1,2}^{(t)} \end{split}$$

Now we turn to  $\nabla_{E2,1}L(W^{(t)},E^{(t)})$ , similarly we have

$$C_0\Phi_2^{(t)}\alpha_1^6(B_{2,1}^{(t)})^3(B_{1,1}^{(t)})^3H_{2,2}^{(t)} \leq \widetilde{O}(1)\left(\frac{B_{2,1}^{(t)}}{B_{1,1}^{(t)}}\right)^3C_0\Phi_1^{(t)}\alpha_1^6(B_{1,1}^{(t)})^6H_{1,2}^{(t)}$$

and since  $H_{2,1}^{(t)} \leq O(C_2) = O(H_{1,2}^{(t)})$  by Lemma B.6c, we can go through the same arguments again to obtain

$$\left| C_0 \Phi_2^{(t)} \alpha_2^6 (B_{1,2}^{(t)})^3 (B_{2,2}^{(t)})^3 H_{2,1}^{(t)} \right| \leq \widetilde{O}(1) \left( \frac{B_{1,2}^{(t)}}{B_{1,1}^{(t)}} \right)^3 C_0 \Phi_1^{(t)} \alpha_1^6 (B_{1,1}^{(t)})^6 H_{1,2}^{(t)} 
\left| C_0 \Phi_2^{(t)} \alpha_2^6 (B_{1,2}^{(t)})^3 (B_{2,1}^{(t)})^3 K_{2,1}^{(t)} \right| \leq \widetilde{O}(\frac{\alpha_1^6}{d^{3/2}}) \left( \frac{B_{1,2}^{(t)}}{B_{1,1}^{(t)}} \right)^3 C_0 \Phi_1^{(t)} \alpha_1^6 (B_{1,1}^{(t)})^6 H_{1,2}^{(t)}$$

Now the proof is complete.

Also, we will need the following lemma controlling gradient bounds for the noise term.

**Lemma B.12** (update of  $R_{1,2}^{(t)}$  in phase I). Suppose Induction B.3 holds at iteration  $t \leq T_1$ , then we have  $(a) |\langle -\nabla_{w_1}L(W^{(t)}, E^{(t)}), \Pi_{V^{\perp}}w_2^{(t)}\rangle| \leq \widetilde{O}(\frac{1}{\sqrt{d}} + \varrho)\Lambda_{1,1}^{(t)}B_{1,1}^{(t)}$ 

(b) 
$$|\langle -\nabla_{w_2} L(W^{(t)}, E^{(t)}), \Pi_{V^{\perp}} w_1^{(t)} \rangle| \leq \widetilde{O}(\frac{1}{\sqrt{d}} + \varrho) \Lambda_{1,1}^{(t)} B_{1,1}^{(t)}$$

Proof. Proof of (a): Firstly, by Claim B.7a, we can directly write

$$\begin{split} \langle \nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} w_2^{(t)} \rangle &= -\sum_{j,\ell} \Sigma_{j,\ell}^{(t)} \langle \nabla_{w_1} \mathcal{E}_{j,3-j}^{(t)}, w_2^{(t)} \rangle \\ &= -\Lambda_{1,1}^{(t)} B_{1,1}^{(t)} \sum_{(j,\ell) \in [2]^2} \frac{(B_{j,\ell}^{(t)})^6 + E_{j,3-j}^{(t)} (B_{3-j,\ell}^{(t)})^3 (B_{j,\ell}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}} \langle \nabla_{w_1} \mathcal{E}_{j,3-j}^{(t)}, w_1^{(t)} \rangle \end{split} \tag{B.5) Now we discuss}$$

each summand respectively: for (j, ) = (1,1), we have

$$\frac{(B_{j,\ell}^{(t)})^6 + E_{j,3-j}^{(t)}(B_{3-j,\ell}^{(t)})^3(B_{j,\ell}^{(t)})^3}{(B_{1,1}^{(t)})^6} = 1 + E_{1,2}^{(t)} \frac{(B_{2,1}^{(t)})^3}{(B_{1,1}^{(t)})^3} = 1 + o(\frac{1}{d^{3/2}(B_{1,1}^{(t)})^3})$$
(B.6)

where the last one is due to Induction B.3d. And for `= 2, we can see from Induction B.3b and d,

$$\begin{aligned} \text{that } \max^{(j,\ell)\neq(1,1)}|B_{j,\ell}^{(t)}| &= \widetilde{O}(\frac{1}{\sqrt{d}}) \text{ and } \quad E_{j,3-j}^{(t)} \leq o \text{(1) to give} \\ &\frac{(B_{j,2}^{(t)})^6 + E_{j,3-j}^{(t)}(B_{3-j,2}^{(t)})^3(B_{j,2}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}} \leq \widetilde{O}(\frac{1}{d^3}) \frac{1}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}} \end{aligned}$$

 $lpha_\ell B_{j,\ell}^{(t)} \leq O$ (1) for all (j, `)  $\in$  [2] $^2$ , so Lemma B.6a applies for On one hand, when  $t \le T_{1,1}$ , we have both  $\Phi_j$  and results in  $\Phi_2^{(t)}/\Phi_1^{(t)} \leq O$  $B_{i,2}^{(t)}/B_{1,1}^{(t)} \leq$ (1). We can also apply Induction B.3c to have

 $O_{\mathbf{e}}(1)$ . On the other hand, when  $t \in [T_{1,1}, T_1]$ , we have by Induction B.3b and Lemma B.6a,b that  $\Phi_2^{(t)}/\Phi_1^{(t)} \leq \widetilde{O}(\alpha_1^{O(1)}) = d^{o(1)}. \text{ but now} \\ B_{1,1}^{(t)} = d^{-o(1)} \gg \widetilde{O}(d^{-1/2}), \text{ therefore}$ 

$$\widetilde{O}(\frac{1}{d^3}) \frac{1}{(B_{1,1}^{(t)})^6} \frac{\Phi_2^{(t)}}{\Phi_1^{(t)}} \le \widetilde{O}(\frac{1}{d^{3/2}}) \frac{1}{(B_{1,1}^{(t)})^3}$$

So together, they imply

$$\frac{(B_{j,2}^{(t)})^6 + E_{j,3-j}^{(t)}(B_{3-j,2}^{(t)})^3(B_{j,2}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}} \le \widetilde{O}(\frac{1}{d^{3/2}(B_{1,1}^{(t)})^3}) \tag{B.7}$$

and similarly, we have

$$\frac{(B_{2,1}^{(t)})^6 + E_{2,1}^{(t)}(B_{1,1}^{(t)})^3(B_{2,1}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_2^{(t)}}{\Phi_1^{(t)}} \le \widetilde{O}(\frac{1}{d^{3/2}(B_{1,1}^{(t)})^3}) \tag{B.8}$$

Next we turn to  $\langle \nabla_{w_1} \mathcal{E}_{j,3-j}^{(t)}, w_2^{(t)} \rangle$ . When j=1, we can apply Claim B.7d to get

$$\langle \nabla_{w_1} \mathcal{E}_{1,2}^{(t)}, w_2^{(t)} \rangle = O(R_{1,2}^{(t)} + \varrho) + O(E_{1,2}^{(t)}) = O(\varrho + \frac{1}{\sqrt{d}}) + O(E_{1,2}^{(t)}) \le O(\varrho + \frac{1}{\sqrt{d}})$$
(B.9)

and when j = 2, we can apply Claim B.7e to get

$$\langle \nabla_{w_1} \mathcal{E}_{2,1}^{(t)}, w_2^{(t)} \rangle = -(E_{2,1}^{(t)})^2 O(R_{1,2}^{(t)} + \varrho) + O(E_{2,1}) = \widetilde{O}(\frac{1}{d^2})(\varrho + \frac{1}{\sqrt{d}}) + O(\frac{1}{d})$$
(B.10)

Combining (B.5), (B.6), (B.7), (B.8), (B.9), and (B.10) completes the proof of (a).

(t) 
$$\langle 
abla_{w_2} \mathcal{E}_{1,2}^{(t)}, w_1^{(t)} 
angle$$

**Proof of (b):** The  $\Sigma_{j,\cdot}$  part is the same as in the proof of (a), so we only deal with and  $\langle \nabla_{w_2} \mathcal{E}_{2,1}^{(t)}, w_1^{(t)} \rangle_{\text{here. For}} \langle \nabla_{w_2} \mathcal{E}_{2,1}^{(t)}, w_1^{(t)} \rangle_{\text{, we apply Claim B.7d to get}}$ 

$$\langle \nabla_{w_2} \mathcal{E}_{2,1}^{(t)}, w_1^{(t)} \rangle = O(R_{1,2}^{(t)} + \varrho) + O(1)E_{1,2}^{(t)}$$
(B.11)

and for  $\langle \nabla_{w_2} \mathcal{E}_{1,2}^{(t)}, w_1^{(t)} \rangle$ , we have

$$\langle \nabla_{w_2} \mathcal{E}_{1,2}^{(t)}, w_1^{(t)} \rangle = O(R_{1,2}^{(t)} + \varrho) (E_{2,1}^{(t)})^2 + O(1) E_{2,1}^{(t)}$$
(B.12)

Inserting (B.6), (B.7), (B.8) and (B.11), (B.12) into the expression of  $h-\nabla_{w^2}L(W^{(t)},E^{(t)})$ ,  $\Pi_{V^{\perp}}w_1^{(t)}$  if finishes the proof of (b).

### B.4 At the End of Phase I

**Lemma B.13** (Phase I). Suppose  $\eta \leq \frac{1}{\mathsf{poly}(d)}$  is sufficiently small, then Induction B.3 holds for at least all  $t \leq T_1 = O(\frac{d^2}{\eta})$ , and at iteration  $t = T_1$ , we have

(a) 
$$B_{1,1}^{(T_1)} = \Omega(1)$$

(b) 
$$||w_j^{(T_1)}||_2 = 1 \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

$$(c) \ \ B_{2,1}^{(T_1)} = \widetilde{\Theta}(\tfrac{1}{\sqrt{d}})_{\textit{and}} B_{j,2}^{(T_1)} = B_{j,2}^{(0)}(1 \pm o(1))_{\textit{for } j \in \texttt{[2]}};$$

$$(d) \ E_{2,1}^{(T_1)} = \widetilde{O}(\tfrac{\eta_E/\eta}{d}) \quad \text{and} \quad E_{1,2}^{(T_1)} \leq \widetilde{O}(\varrho + \tfrac{1}{\sqrt{d}}) \quad ;$$

(e) 
$$R_{1,1}^{(T_1)}, R_2^{(T_1)} = \Theta(1)_{and} R_{1,2}^{(T_1)} = \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

*Proof.* We begin by first prove the existence of  $T_1:=\min\{t:B_{1,1}^{(t)}\geq 0.01\}=O(\frac{d^2}{\eta})$  if Induction

B.3 holds whenever  $B_{1,1}^{(t)} \le 0$ :01, then we will turn back to prove Induction B.3 holds throughout  $t \le T_1$ . We split the analysis into two stages:

**Proof of**  $T_1 \leq O(\frac{d^2}{\eta})$ : By Lemma B.9a we can write down the update of  $B_{1,1}^{(t)}$  as

$$B_{1,1}^{(t+1)} = B_{1,1}^{(t)} + \eta(1 \pm \widetilde{O}(\frac{1}{d}))\Lambda_{1,1}^{(t)} = B_{1,1}^{(t)} + \eta(1 \pm \widetilde{O}(\frac{1}{d}))\Phi_1^{(t)}C_0\alpha_1^6H_{1,2}^{(t)}(B_{1,1}^{(t)})^5$$
(B.13)

When  $\alpha_1 B_{1,1}^{(t)} \leq O$ (1), by Lemma B.6a,c we have  $\Phi$  can lower bound the update as

$$H_{1}^{(t)}=\Theta(\frac{1}{C_{2}^{2}})\text{ and }H_{1,2}^{(t)}=\Omega(C_{2})\text{, this means we}$$
 
$$B_{1,1}^{(t+1)}\geq B_{1,1}^{(t)}+\Omega(\frac{\eta C_{0}\alpha_{1}^{6}}{C_{2}})(B_{1,1}^{(t)})^{5}$$

since  $\frac{C_0\alpha_1^6}{C_2}$  is a constant, we know there exist some  $t^0 \ge 0$  such that  $B_{1,1}^{(t')} \ge \Omega(\frac{1}{\alpha_1})$ . Also recall that

$$T_{1,1} := \min\{t : B_{1,1}^{(t)} \geq \Omega(\frac{1}{\alpha_1})\}. \text{ So by Lemma G.1, where } \eta = \frac{1}{\mathsf{poly}(d)}, C_t = \Omega(\frac{C_0\alpha_1^6}{C_2}) \ \delta = \underset{\mathsf{polylog}}{\underbrace{\qquad \qquad }} 1 \ (d)$$

and  $A = \Omega(\frac{1}{\alpha_1}), \log(A/B_{1,1}^{(0)}) = \widetilde{O}_{1}$ , we have

$$T_{1,1} = O\left(\frac{C_2}{\eta C_0 \alpha_1^6}\right) \sum_{x_t \le O\left(\frac{1}{\alpha_1}\right)} \eta C_t \le O\left(\frac{C_2}{\eta C_0 \alpha_1^6}\right) \left(O(1) + \frac{\widetilde{O}(\eta)}{B_{1,1}^{(0)}}\right) \frac{1}{(B_{1,1}^{(0)})^4} \le \widetilde{O}\left(\frac{1}{\eta \alpha_1^6 (B_{1,1}^{(0)})^4}\right)$$

Since  $(B_{1,1}^{(0)})^4 \geq \widetilde{\Omega}(\frac{1}{d^2})$  from our initialization, we have  $T_{1,1} \leq O(\frac{d^2}{\eta})$  and thus  $T_{1,1}$  exists. Now we consider when  $B_{1,1}^{(t)} \geq \Omega(\frac{1}{\alpha_1})$ . Now by Lemma B.6b,c, we have  $\Phi_1^{(t)} \geq \Omega((C_2 + \alpha_1^6)^{-2})$ , which gives an update:

$$B_{1,1}^{(t+1)} \ge B_{1,1}^{(t)} + \Omega\left(\frac{\eta C_0 \alpha_1^6}{(C_2 + \alpha_1^6)^2}\right) (B_{1,1}^{(t)})^5$$

so again by Lemma G.1, choosing  $C_t = \Omega(\frac{C_0 lpha_1^6}{(C_2 + lpha_1^6)^2})$ 

$$T_1 = \frac{O((C_2 + \alpha_1^6)^2)}{\eta C_0 \alpha_1^6} \sum_{x_t \in [\Omega(\frac{1}{\Omega^4}), 0.01]} \eta C_t \le \left( O(1) + \frac{\widetilde{O}(\eta)}{B_{1,1}^{(T_{1,1})}} \right) \frac{\widetilde{O}(\alpha_1^{12})}{(B_{1,1}^{(T_{1,1})})^4} \le \widetilde{O}(\frac{\alpha_1^6}{\eta (B_{1,1}^{(T_{1,1})})^4}) \le O(\frac{\alpha_1^6}{\eta})$$

where  $O(\frac{\alpha_1^6}{\eta}) \ll O(\frac{d^2}{\eta})$ , so we have proved that  $T_1$  exist. Now we begin to prove that Induction B.3 holds for all  $t \le T_1$ .

**Proof of Induction B.3:** We first prove (b)–(d), and then come back to prove (a) and (d). At t = 0, we know all induction holds from Properties B.1. Now we suppose Induction B.3 holds for all iterations  $\leq t - 1$  and prove it holds at t.

**The growth of**  $B_{2,1}^{(t)}$ : Applying Lemma B.9, we have for  $t \le T_{1,1}$ 

$$B_{1,1}^{(t+1)} \ge B_{1,1}^{(t)} + \eta(1 - \widetilde{O}(\frac{1}{d}))\Lambda_{1,1}^{(t)}$$

$$B_{2,1}^{(t+1)} \le B_{2,1}^{(t)} + \eta (1 + O(\frac{1}{\sqrt{d}})) \Lambda_{2,1}^{(t)} + \eta \frac{(B_{2,1}^{(t)})^2}{(B_{1,1}^{(t)})^2} E_{1,2}^{(t)} \Lambda_{1,1}^{(t)}$$

 $\text{For some} t_1' := \min\{t: B_{1,1}^{(t)} \geq \tfrac{\Omega(1)}{d^{0.49}}\} \text{, we have} E_{1,2}^{(t)} \leq \widetilde{O}(B_{1,1}^{(t)}\varrho) \lesssim \tfrac{1}{d^{0.49}} \operatorname{during} t \leq t_1' \text{, and}$ 

$$\frac{(B_{2,1}^{(t)})^2}{(B_{1,1}^{(t)})^2} E_{1,2}^{(t)} \Lambda_{1,1}^{(t)} \lesssim \frac{(B_{2,1}^{(t)})^2}{d^{0.49} (B_{1,1}^{(t)})^2} \Lambda_{1,1}^{(t)} \leq \widetilde{O}(\frac{1}{d^{0.49}}) \Lambda_{2,1}^{(t)}$$

which allow us to give an upper bound  $\operatorname{to}^{B_{2,1}^{(t+1)}}$  as

$$\leq (1 + O(\frac{1}{d^{0.49}}))\Phi_2^{(\cdot)}C_0\alpha_1^6C_2\mathcal{E}_2^{(\cdot)}(1 + \frac{1}{d^{0.49}})(B_{2,1}^{(t)})^5 \qquad \qquad t \leq t_1')$$
 
$$B_{2,1}^{(t+1)} \leq (1 + O(\frac{1}{\sqrt{d}}))\Lambda_{2,1}^{(t)} + \widetilde{O}(\frac{1}{d^{0.49}})\Lambda_{2,1}^{(t)}$$
 
$$\sim \underline{1} \qquad t \qquad \underline{1} \qquad \text{(when polylog)}$$

Since we also have

$$B_{1,1}^{(t+1)} \geq (1 - \widetilde{O}(\frac{1}{d}))\Lambda_{1,1}^{(t)} \geq (1 - \widetilde{O}(\frac{1}{d}))\Phi_1^{(t)}C_0\alpha_1^6\mathcal{E}_1^{(t)}(1 - \frac{1}{(d)})(B_{1,1}^{(t)})^5 - \frac{1}{(d)}$$

Since  $B_{1,1}^{(0)} \geq B_{2,1}^{(0)} (1 + \Omega(\frac{1}{\log d}))$ , we can now apply Corollary G.2 to the two sequence  $B_{1,1}^{(t+1)}$  and  $B_{2,1}^{(t+1)}$ , where  $S_t = \frac{\Phi_1^{(t)} \mathcal{E}_1^{(t)}}{\Phi_2^{(t)} \mathcal{E}_2^{(t)}} (1 + \log B_{1,1}^{(t)})$  to get

$$B_{1,1}^{(t_1')} \ge \frac{1}{d^{0.499}} \quad \text{while } B_{2,1}^{(t_1')} \le \widetilde{O}(\frac{1}{\sqrt{d}})$$

Note that here the update of  $B_{2,1}^{(t)}$  at every step satisfies  $\operatorname{sign}(B_{2,1}^{(t+1)}-B_{2,1}^{(t)})=\operatorname{sign}(B_{2,1}^{(t)})$  which implies  $B_{2,1}^{(t'_1)}=\widetilde{\Theta}(\frac{1}{\sqrt{d}})$ . Now for every  $T\in[t'_1,T_1]$ , we can apply Lemma G.3 to get that

$$\sum_{t \in [t_1',T]} \eta \frac{(B_{2,1}^{(t)})^2}{(B_{1,1}^{(t)})^2} E_{1,2}^{(t)} \Lambda_{1,1}^{(t)} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) O(\frac{1}{B_{1,1}^{(t_1')}}) \max_{t \leq T} \{(B_{2,1}^{(t)})^2\} \leq O(\frac{1}{d^{0.5 + \Omega(1)}})$$

Suppose we have proved that  $B_{2,1}^{(t)} \leq \widetilde{O}(\frac{1}{\sqrt{d}})$  for each  $t \leq T$ , we define a new sequence

$$\begin{split} \widetilde{B}_{2,1}^{(t+1)} &= \widetilde{B}_{2,1}^{(t)} + \eta(1 + \widetilde{O}(\frac{1}{d^{0.49}}))\Phi_2^{(t)}C_0\alpha_1^6C_2\mathcal{E}_2^{(t)}(1 + \mathbf{1} \\ &\qquad \qquad \text{, polylog} \\ \widetilde{B}_{2,1}^{(t_1')} &= B_{2,1}^{(t_1')} + \sum_{t \in [t_1',T]} \eta \frac{(B_{2,1}^{(t)})^2}{(B_{1,1}^{(t)})^2} E_{1,2}^{(t)} \Lambda_{1,1}^{(t)} &= (1 \pm o(1))\widetilde{B}_{2,1}^{(t_1')} \\ \text{where} \end{split}$$

It can be directly seen that  $|\widetilde{B}_{2,1}^{(t)} - \widetilde{B}_{2,1}^{(0)}| \ge |B_{2,1}^{(t)} - B_{2,1}^{(0)}|_{\mbox{for all}}$   $t \in [t_1', T]$ . Notice that now  $\widetilde{B}_{2,1}^{(t_1')} \le d^{\Omega(1)}B_{1,1}^{(t_1')}$ , we can now apply Corollary G.2 again to get

Now we deal with  $t \in [T_{1,1}, T_1]$ . During this stage, we can directly apply Corollary G.2 to  $\widetilde{B}_{2,1}^{(t)}$  and  $B_{1,1}^{(t)}$ , where  $S_t = \frac{\Phi_1^{(t)} H_{1,2}^{(t)}}{\Phi_2^{(t)} H_{2,2}^{(t)}} \leq O(\alpha_1^{O(1)})$ , to get that

$$|B_{2(T,1)} - B_{2(0),1}| \le |B_{2(T,1)} - B_{2(0),1}| \le \sqrt{\frac{1}{d}}$$
(for every  $T \le T_1$ )
$$dpolylog(d)$$

And thus by Lemma B.1, we have  $B_{2,1}^{(T)}=B_{2,1}^{(0)}(1\pm o(1)).$ 

The growth of  $B_{1,2}^{(t)}$  and  $B_{2,2}^{(t)}$ : By Lemma B.10, we can write down the update as

$$B_{j,2}^{(t+1)} = B_{j,2}^{(t)} + \eta \left( 1 \pm \widetilde{O}(\alpha_1^6) (E_{3-j,j}^{(t)} + (B_{j,1}^{(t)})^3) \right) \Lambda_{j,2}^{(t)}$$

Since  $B_{2,1}^{(t)} \leq \widetilde{O}(\frac{1}{\sqrt{d}})$  and  $E_{1,2}^{(t)} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})B_{1,1}^{(t)}, E_{2,1}^{(t)} \leq \widetilde{O}(\frac{1}{d})$  because we chose  $\eta_{\mathit{E}} \leq \eta$ , we only need to care about  $(B_{1,1}^{(t)})^3$  in the update expression. Now define  $t_2' := \min\{t : B_{1,1}^{(t)} \geq \Omega(\frac{1}{\alpha^2})\}$ , we have

• For 
$$t \leq t_2'$$
 , by Corollary G.2 and setting  $O(\frac{\alpha_2^6 \Phi_j^{(t)} H_{j,1}^{(t)}}{\alpha_1^6 \Phi_j^{(t)} H_{1,2}^{(t)}}) \leq \widetilde{O}(\frac{\alpha_2^6}{\alpha_1^6}) \ll \frac{x_t = B_{1,1}^{(t)}, \ C_t = (1 - \widetilde{O}(\frac{1}{d})) \Phi_1^{(t)} C_0 \alpha_1^6 H_{1,2}^{(t)}, \ S_t = \frac{|B_{j,2}^{(t)} - B_{j,2}^{(0)}| \leq O(\frac{\alpha_2^6}{\alpha_1^6} \sqrt[4]{a})}{\sum_{j=1}^{d_{polylog}(d)} (\text{by Lemma B.6a,c}), \text{ we have}} \frac{|B_{j,2}^{(t)} - B_{j,2}^{(0)}| \leq O(\frac{\alpha_2^6}{\alpha_1^6} \sqrt[4]{a})}{\sum_{j=1}^{d_{polylog}(d)} (\sqrt[4]{a_{polylog}(d)})} \in [\Omega(\sqrt[4]{a_{polylog}(d)}), O(\sqrt[4]{a_{polylog}(d)})$ 

Lemma B.1.

• For  $t \in [t'_2, T_1]$ , we can use Corollary G.2 again and let  $x_t = B_{1,1}^{(t)}$ , we know  $B_{1,1}^{(t'_2)} \ge d^{\Omega(1)} B_{2,1}^{(t'_2)}$ . Setting  $C_t = (1 - Q(\frac{1}{d})) \Phi_1^{(t)} C_0 \alpha_1^6 H_{1,2}^{(t)}$ ,  $S_t = O((1 + \alpha_1^6) \frac{\alpha_2^6 \Phi_j^{(t)} H_{1,1}^{(t)}}{\alpha_1^6 \Phi_1^{(t)} H_{1,2}^{(t)}}) \le O(\alpha^{O(1)})$ , we can have  $|B_{j,2}^{(t)} - B_{j,2}^{(t^0_2)}| - \sqrt[4]{\frac{1}{d \operatorname{polylog}(d)}}$   $B_{j,2}^{(t)} \in [\Omega(\sqrt[4]{\frac{1}{d \log d}}), O(\sqrt[4]{\frac{\log d}{d}})$   $t \in [t_2^0, T_1]$ . which implies

This proves Induction B.3b. Indeed, simple calculations also proves Induction B.3c, since the update of  $B_{1,1}^{(t)}$  is always larger than others' during  $t \le T_1$ .

For Induction B.3d: From Lemma B.11, we can write down the update

$$-\nabla_{E_{1,2}}L(W^{(t)},E^{(t)}) = O(\Lambda_{1,1}^{(t)}B_{1,1}^{(t)}) \left( -C_1E_{1,2}^{(t)} + \widetilde{O}(\frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^3}) + C_2(R_{1,2}^{(t)} + \varrho) \right)$$

for some constants  $C_1, C_2 = \Theta(1)$ . Applying Lemma G.3 to  $O(\Lambda_{1,1}^{(t)} B_{1,1}^{(t)}) \frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^3}$ , we can obtain

$$\sum_{t \leq T} O(\eta_E \Lambda_{1,1}^{(t)} B_{1,1}^{(t)}) \frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^3} = \frac{\eta_E}{\eta} \sum_{t \leq T} O(\eta \Lambda_{1,1}^{(t)}) \frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^2} \leq \widetilde{O}(\frac{\eta_E/\eta}{d^{3/2}}) \frac{1}{B_{1,1}^{(0)}} \leq \widetilde{O}(\frac{\eta_E/\eta}{d})$$

So here it suffices to notice that whenever  $|E_{1,2}^{(t)}| < 2\frac{C_2}{C_1}(R_{1,2}^{(t)} + \varrho)$  (which is obviously satisified at t = 0), we would have

$$O(\Lambda_{1,1}^{(t)}B_{1,1}^{(t)})\left(-O(E_{1,2}^{(t)}) + C_2(R_{1,2}^{(t)} + \varrho)\right) = -O(\Lambda_{1,1}^{(t)}B_{1,1}^{(t)})\widetilde{O}(R_{1,2}^{(t)} + \varrho) \le O(\Lambda_{1,1}^{(t)}B_{1,1}^{(t)})\widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

In that case, we will always have (since  $E_{1,2}^{(0)}=0$ )

$$E_{1,2}^{(t+1)} \le \left| \sum_{t \le T} \widetilde{O}(\eta_E \Lambda_{1,1}^{(t)} B_{1,1}^{(t)}) \frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^3} \right| + \sum_{s \le t} O(\eta_E \Lambda_{1,1}^{(s)} B_{1,1}^{(s)}) (R_{1,2}^{(s)} + \varrho) \le \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \frac{\eta_E}{\eta} B_{1,1}^{(t+1)}$$

Similarly for  $\nabla_{E2,1}L(W^{(t)},E^{(t)})$ , we can write down

$$-\nabla_{E_{2,1}}L(W^{(t)}, E^{(t)}) = \widetilde{O}(\frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^2})\Lambda_{1,1}^{(t)} + \sum_{\ell \in [2]} C_2\Lambda_{2,\ell}^{(t)}B_{2,\ell}^{(t)} \left(-O(E_{2,1}^{(t)}R_2^{(t)}) + O(R_{1,2}^{(t)})\right)$$

by Lemma G.3, we have

$$\sum_{t \le T_1} \eta_E \widetilde{O}(\frac{(B_{1,2}^{(t)})^3}{(B_{1,1}^{(t)})^2}) \Lambda_{1,1}^{(t)} \le \widetilde{O}(\frac{\eta_E/\eta}{d})$$

and since from previous comparison results we know that

$$\sum_{t \le T_1} \sum_{\ell \in [2]} \eta_E C_2 \Lambda_{2,\ell}^{(t)} B_{2,\ell}^{(t)} = \frac{\eta_E}{\eta} \sum_{t \le T_1} \sum_{\ell \in [2]} \eta C_2 \Lambda_{2,\ell}^{(t)} B_{2,\ell}^{(t)} \le \widetilde{O}(\frac{\eta_E/\eta}{d})$$

we can then prove the claim.

**For Induction B.3a:** We can write down the update of  $\|w_j^{(t)}\|_2^2$  as follows:

$$kw_{j(t+1)}k_{22} = kw_{j(t)} - \eta \nabla_{w_{j}} L(W(t), E(t))k_{22}$$

$$= kw_{j(t)}k_{22} - \eta h \nabla_{w_{j}} L(W(t), E(t)), w_{j(t)}i + \eta_{2}k \nabla_{w_{j}} L(W(t), E(t))k_{22}$$

from (A.2) and Induction B.3a,b,c at iteration t and our assumption on  $\xi_p$ , we know  $\|\nabla_{w_i}L(W^{(t)},E^{(t)})\|_2^2 \leq \widetilde{O}(d)$ 

which allow us to choose  $\eta \leq \frac{1}{\mathsf{poly}(d)}$  to be small enough so that  $\eta dT_1 \leq \frac{1}{\eta \mathsf{poly}(d)}$ . Then by Lemma B.8b, we have

$$(t+1) 2 kw_j(0)k22 \pm \eta X |h\nabla w_j L(W(s), E(s)), w_j(s)i| \pm poly1(d)$$

$$kw_j k2 =$$

$$\leq \|w_j^{(0)}\|_2^2 \pm \eta \sum_{s \leq t}^{s \leq t} \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) |\Lambda_{1,1}^{(s)}| \sum_{j \in [2]} |E_{j,3-j}^{(s)}| \pm \frac{1}{\operatorname{poly}(d)}$$

Since from the above analysis of the update of  $B_{1,1}^{(t)}$ , we know  $\sum_{t \leq T_1} \Lambda_{1,1}^{(t)} \leq O$  (1). Moreover, we also know that  $|B_{1,1}^{(t)}|$  is increasing and  $\operatorname{sign}(\Lambda_{1,1}^{(t)}) = \operatorname{sign}(\Lambda_{1,1}^{(s)})$  for any  $s,t \leq T_1$ . Thus they imply  $\sum_{s \leq t} |\Lambda_{1,1}^{(s)}| = |\sum_{s \leq t} \Lambda_{1,1}^{(s)}| = O$  (1), which can be combine with Induction B.3d to prove the claim.

**Proof of Induction B.3e:** We can write down the update of  $R_{1,2}^{(t)} = \langle \Pi_{V^\perp} w_1^{(t)}, w_2^{(t)} \rangle$  as follows

$$\begin{split} h\Pi V_{\perp} w_{1}(t+1), & w_{2}(t+1)\mathrm{i} = h\Pi V_{\perp} w_{1}(t) - \Pi V_{\perp} \eta \nabla_{w_{1}} L\big(W(t), E(t)\big), \Pi V_{\perp} w_{2}(t) - \Pi V_{\perp} \eta \nabla_{w_{2}} L\big(W(t), E(t)\big)\mathrm{i} \\ &= R_{1}(t,)_{2} - \eta h \nabla_{w_{1}} L\big(W(t), E(t)\big), \Pi V_{\perp} w_{2}(t)\mathrm{i} - \eta h \nabla_{w_{2}} L\big(W(t), E(t)\big), \Pi V_{\perp} w_{1}(t)\mathrm{i} \\ &+ \eta_{2} h \Pi V_{\perp} \nabla_{w_{1}} L\big(W(t), E(t)\big), \Pi V_{\perp} \nabla_{w_{2}} L\big(W(t), E(t)\big)\mathrm{i} \end{split}$$

By Cauchy-Schwarz inequality and the same analysis above we have

$$|h\Pi V_{\perp} \nabla w_1 L(W(t), E(t)), \Pi V_{\perp} \nabla w_2 L(W(t), E(t)) \mathbf{i}| \leq k \nabla w_1 L(W(t), E(t)) k 2 k \nabla w_2 L(W(t), E(t)) k 2$$

$$\leq O_{e}(d)$$

so by our choice of  $\eta$ 

 $X \eta_2 |h\Pi V_{\perp} \nabla_{w_1} L(W(t), E(t)), \Pi V_{\perp} \nabla_{w_2} L(W(t), E(t)) i| \leq \underline{\qquad} 1 \text{ poly}(d) t \leq T_1$ 

and by Lemma B.12 we have

$$\left| -\eta \langle \nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^{\perp}} w_2^{(t)} \rangle - \eta \langle \nabla_{w_2} L(W^{(t)}, E^{(t)}), \Pi_{V^{\perp}} w_1^{(t)} \rangle \right| \leq \eta \widetilde{O}(\Lambda_{1,1}^{(t)} B_{1,1}^{(t)}) (\varrho + \frac{1}{\sqrt{d}})$$

which implies

$$\leq \widetilde{O}(\frac{1}{\sqrt{d}}) + \sum_{s \leq t} \eta \widetilde{O}(\Lambda_{1,1}^{(s)} B_{1,1}^{(s)}) + \cdots \qquad \qquad 1$$
 
$$\leq \widetilde{O}(\frac{1}{\sqrt{d}}) + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) B_{1,1}^{(t+1)}$$
 
$$\leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

which completes the proof of Induction B.3. As for (a) – (e) of Lemma B.13, they are just direct corrolary of our induction at  $t = T_1$ .

# C Phase II: The Substitution Effect of Prediction Head

In this phase, As  $B_{1,1}^{(t)}$  is learned to become very large  $(B_{1,1}^{(t)} \gtrsim \|w_1^{(t)}\|_2)$ . The focus now shift to grow  $E_{2,1}^{(t)}$ , because we want  $C_1\alpha_1^6((B_{2,1}^{(t)})^3 + E_{2,1}^{(t)}(B_{1,1}^{(t)})^3)^2$  in  $H_{2,1}^{(t)}$  to dominate  $\mathcal{E}_{2,1}^{(t)}$ . We can write down the gradient of  $E_{2,1}^{(t)}$  as

$$-\nabla_{E_{2,1}}L(W(t),E(t)) = X C_0\Phi(2t)\alpha^{\epsilon}(B_{2}(t,))3((B_{1}(t,))3H_{2}(t,3)-(B_{2}(t,3$$

Now let us define

$$T_2 := \min\{t : R_2^{(t)} < \frac{1}{\log d} |E_{1,2}^{(t)}|\}$$
 (C.1)

We will prove that  $E_{2,1}^{(T_2)}$  reaches at most  $O(p_{\eta_E/\eta})$  and holds throughout  $t \in [T_1, T_2]$ . In this phase, the learning  $E_{2,1}^{(t)}$  of is much faster than the growth of the

first feature  $v_1$  such that  $T_2 - T_1 = o(T_1/d)$ , which is due to the acceleration effects brought by  $B_{1,1}^{(t)} = \Omega(1)$  during this phase.

### **C.1** Induction in Phase II

We will be based on the following induction hypothesis during phase II.

**Inductions C.1** (Phase II). When  $t \in [T_1, T_2]$ , we hypothesize the followings would hold

$$\begin{array}{l} (a) \ \ B_{1,1}^{(t)} = \Theta(1), \ B_{j,\ell}^{(t)} = B_{j,\ell}^{(T_1)}(1 \pm o(1)) = \widetilde{\Theta}(\frac{1}{\sqrt{d}}) \\ for \ \mbox{($j$,")} \ 6 = \mbox{(1,1)} \ and \\ \sin(B_{j,\ell}^{(t)}) = \sin(B_{j,\ell}^{(T_1)}), \\ (b) \ \ |R_{1,2}^{(t)}| = \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})\alpha_1^{O(1)}[R_1^{(t)}]^{1/2}[R_2^{(t)}]^{1/2} \\ (c) \ \ R_1^{(t)} \in [\Omega(\frac{1}{d^{3/4}\alpha_1^2}), O(1)], \ R_2^{(t)} \in [\Omega(\frac{1}{\log d}\sqrt{\eta_E/\eta}), O(1)], \\ \end{array}$$

$$(d) \ E_{1,2}^{(t)} \leq \widetilde{O}(\varrho + \tfrac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2} \ \text{and} \ E_{2,1}^{(t)} \leq O(\sqrt{\eta_E/\eta}).$$

Under Induction C.1, we have some results as direct corollary.

**Claim C.2.** At each iteration  $t \in [T_1, T_2]$ , if Induction B.3 holds, then

(a) 
$$\mathcal{E}_{j}^{(t)} = \Theta(C_{2}[R_{j}^{(t)}]^{3}),$$
  
(b)  $\mathcal{E}_{j,3-j}^{(t)} = \mathcal{E}_{j}^{(t)} \pm \widetilde{O}(E_{j,3-j}^{(t)}(\varrho + \frac{1}{\sqrt{d}})[R_{1}^{(t)}]^{3/2}[R_{2}^{(t)}]^{3/2}) + O((E_{j,3-j}^{(t)})^{2}[R_{3-j}^{(t)}]^{3})$  for each  $j \in [2]$ :

*Proof.* It is trivial to derive (a) from the expression of  $\mathcal{E}_{j}^{(t)}$  and our assumption of  $\xi_{p}$ . For (b) it suffices to directly calculate the expression of  $\mathcal{E}_{j,3-j}^{(t)}$  along with Induction C.1b.  $\square$ 

**Lemma C.3** (variables control in phase II). *In Phase II*  $(t \in [T_1, T_2])$ , *if Induction C.1 holds, then* 

$$(a) \ \Phi_{1}^{(t)} = \widetilde{\Theta}(\frac{1}{\alpha_{1}^{12}}), \ \Phi_{2}^{(t)} = \Theta((C_{2}[R_{2}^{(t)}]^{3} + C_{1}\alpha_{1}^{6}(E_{2,1}^{(t)})^{2})^{-2}),$$

$$(b) \ K_{1,\ell}^{(t)} = \widetilde{O}(\alpha_{\ell}^{6}/d^{3/2}), \ K_{2,\ell}^{(t)} = \widetilde{O}(E_{2,1}^{(t)}\alpha_{\ell}^{6}/d^{3/2} + \alpha_{\ell}^{6}/d^{3})$$

$$(c) \ H_{1,1}^{(t)} = \Theta(C_{1}\alpha_{1}^{6}), \ H_{1,2}^{(t)} = \widetilde{O}([R_{1}^{(t)}]^{3}), \ H_{2,2}^{(t)} = \Theta(C_{2}[R_{2}^{(t)}]^{3}), \ H_{2,1}^{(t)} = \Theta(C_{2}[R_{2}^{(t)}]^{3} + C_{1}\alpha_{1}^{6}(E_{2,1}^{(t)})^{2})$$

*Proof.* The proof of (a) directly follows from Induction C.1a,c and Claim C.2. The proof of (b) follows directly from the expression of  $K_{i}$  and Induction C.1a,d. The proof of (c) is also similar.  $\square$ 

#### C.2 Gradient Lemmas for Phase II

**Lemma C.4** (learning prediction head  $E_{1,2}$ ,  $E_{2,1}$  in phase II). *If Induction C.1 holds at iteration*  $t \in [T_1, T_2]$ , then we have

(a) 
$$-\nabla_{E_{1,2}}L(W^{(t)}, E^{(t)}) = (1 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{3/2}}))\Sigma_{1,1}^{(t)}(-2E_{1,2}^{(t)}[R_2^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2})$$

$$\pm \Sigma_{1,1}^{(t)}\widetilde{O}(\frac{\eta_E/\eta}{\sqrt{d}})\max\{[R_1^{(t)}]^3, \frac{\alpha_1^{O(1)}}{d^{5/2}}\},$$

$$(b) \quad -\nabla_{E_{2,1}}L(W^{(t)}, E^{(t)}) = (1 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{3/2}}))C_0\Phi_2^{(t)}\alpha_1^6(B_{2,1}^{(t)})^3(B_{1,1}^{(t)})^3H_{2,2}^{(t)}$$

$$\pm O(\Sigma_{2,1}^{(t)})(|E_{2,1}^{(t)}|[R_1^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2})$$

*Proof.* We first write down the gradient for  $E_{j,3-j}^{(t)}$ : (ignoring the time superscript  $^{(t)}$ )

$$-\nabla E_{j,3-j}L(W,E) = X C_0 \Phi_{j}\alpha' G_{j,3-j} (B_{33-j,'}H_{j,3-'} - B_{33-j,3-'}K_{j,3-'}) - X \Sigma_{j,'}\nabla E_{j,3-j}E_{j,3-j}$$

$$\stackrel{\cdot}{}\in [2] \qquad \stackrel{\cdot}{}\in [2] \qquad \text{where}$$

$$\nabla_{E_{j,3-j}}\mathcal{E}_{j,3-j}=\mathbb{E}\left[2\langle w_j,\xi_p\rangle^3\langle w_{3-j},\xi_p\rangle^3+2E_{j,3-j}\langle w_{3-j},\xi_p\rangle^6\right]\!.$$
 Thus we have

$$\nabla_{E_{j,3-j}}\mathcal{E}_{j,3-j}^{(t)} = 2E_{j,3-j}^{(t)}[R_{3-j}^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}$$

and by Claim A.1 and Induction C.1a, if (j, ) 6= (1,1)

$$\Sigma_{j,\ell}^{(t)} = O(\Sigma_{1,1}^{(t)}) \frac{(B_{j,\ell}^{(t)})^6 + E_{j,3-j}^{(t)}(B_{3-j,\ell}^{(t)})^3(B_{j,\ell}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}} \le o(\frac{1}{d^{3/2}}) \Sigma_{1,1}^{(t)} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}}$$

Therefore for i = 1:

$$\sum_{\ell \in [2]} \Sigma_{1,\ell}^{(t)} \nabla_{E_{1,2}} \mathcal{E}_{1,2}^{(t)} = (1 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{3/2}})) \Sigma_{1,1}^{(t)} \nabla_{E_{1,2}} \mathcal{E}_{1,2}^{(t)}$$

Now by Induction C.1a,c and Lemma C.3b,c we have  $(B_{1,\ell}^{(t)})^3 H_{1,3-\ell}^{(t)} \leq \max\{\Theta(C_2[R_1^{(t)}]^3), \widetilde{O}(\frac{\alpha_1^6}{d^{3/2}})\}$ , which leads to the bounds

$$|(B_{1,\ell}^{(t)})^3(B_{2,\ell}^{(t)})^3H_{1,3-\ell}^{(t)}| \leq \widetilde{O}(\frac{1}{d^{3/2}})\max\{[R_1^{(t)}]^3,\frac{\alpha_1^6}{d^3}\} \qquad |(B_{1,\ell}^{(t)})^3(B_{2,3-\ell}^{(t)})^3K_{1,3-\ell}^{(t)}| \leq \widetilde{O}(\frac{1}{d^3})$$

which implies

$$\left|\sum_{\ell \in [2]} C_0 \Phi_1^{(t)} \alpha_\ell^6 (B_{1,\ell}^{(t)})^3 ((B_{2,\ell}^{(t)})^3 H_{1,3-\ell}^{(t)} - (B_{2,3-\ell}^{(t)})^3 K_{1,3-\ell}^{(t)})\right| \lesssim \widetilde{O}(\frac{\eta_E/\eta}{\sqrt{d}}) \Sigma_{1,1}^{(t)} \max\{[R_1^{(t)}]^3, \frac{\alpha_1^{O(1)}}{d^{5/2}}\}$$

Combining above together, we have

$$-\nabla_{E_{1,2}}L(W^{(t)},E^{(t)})$$

$$=(1+o(\frac{1}{d^{3/2}}))\Sigma_{1,1}^{(t)}(-2E_{1,2}^{(t)}[R_2^{(t)}]^3\pm O(\overline{R}_{1,2}^{(t)}+\varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}\pm \widetilde{O}(\frac{\eta_E/\eta}{\sqrt{d}})\max\{[R_1^{(t)}]^3,\frac{\alpha_1^{O(1)}}{d^{5/2}}\})$$

For  $-\nabla_{E2,1}L(W^{(t)},E^{(t)})$ , the expression is slightly different, we first observe that by Induction C.1a

$$\Delta_{2,2}^{(t)} \le \widetilde{O}(\frac{1}{d^{3/2}})\Delta_{2,1}^{(t)}$$

Meanwhile, by Induction C.1a and Lemma C.3b,c, we have

$$\Xi_2^{(t)} \le \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^3})C_0C_2\Phi_2^{(t)}[R_2^{(t)}]^3$$

Moreover, we can also calculate  $\Sigma_{2,1}^{(t)}=C_0C_2\alpha_1^6E_{2,1}^{(t)}\Phi_2^{(t)}(B_{2,1}^{(t)})^3)=\widetilde{O}(\frac{\alpha_1^6}{d^{3/2}})\Phi_2^{(t)},\ \Sigma_{2,2}^{(t)}=\widetilde{O}(\frac{\alpha_2^6}{d^3})\Phi_2^{(t)}$ , which gives

$$\sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \nabla_{E_{2,1}} \mathcal{E}_{2,1}^{(t)} = \Sigma_{2,1}^{(t)} (-\Theta(E_{2,1}^{(t)}) [R_1^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2})$$

Now we combine the above results and get

$$-\nabla_{E_{2,1}}L(W^{(t)}, E^{(t)}) = (1 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{3/2}}))C_0\Phi_2^{(t)}\alpha_1^6(B_{2,1}^{(t)})^3(B_{1,1}^{(t)})^3\mathcal{E}_{2,1}^{(t)}$$
$$\pm O(\Sigma_{2,1}^{(t)})(|E_{2,1}^{(t)}|[R_1^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2})$$

**Lemma C.5** (reducing noise in phase II). Suppose Induction C.1 holds at  $t \in [T_1, T_2]$ , then

$$(a) \langle -\nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^{\perp}} w_1^{(t)} \rangle = \Sigma_{1,1}^{(t)} \Theta(-[R_1^{(t)}]^3 \pm \widetilde{O}(|E_{1,2}^{(t)}| + \frac{|E_{2,1}^{(t)}|^2}{d^{3/2}}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2}),$$

$$(b) \langle -\nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^{\perp}} w_2^{(t)} \rangle = \Sigma_{1,1}^{(t)} ((-\Theta(\overline{R}_{1,2}^{(t)}) + O(\varrho)) [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} + \widetilde{O}(|E_{1,2}^{(t)}| + \frac{|E_{2,1}^{(t)}|^2}{d^{3/2}}) R_1^{(t)} [R_2^{(t)}]^2)$$

And furthermore

$$\begin{split} (c) \quad \langle -\nabla_{w_2}L(W^{(t)},E^{(t)}),\Pi_{V^\perp}w_2^{(t)}\rangle &= -\Theta([R_2^{(t)}]^3) \Big( \Sigma_{1,1}^{(t)}\Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big) \\ & \quad \pm O\Big( \sum_{j,\ell} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big); \\ (d) \quad \langle -\nabla_{w_2}L(W^{(t)},E^{(t)}),\Pi_{V^\perp}w_1^{(t)}\rangle &= \Big( \Sigma_{1,1}^{(t)}\Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big) (-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_2^{(t)}]^{5/2} [R_1^{(t)}]^{1/2} \\ &\quad + O\Big( \sum_{j,\ell} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} R_2^{(t)} [R_1^{(t)}]^2 \Big) \end{split}$$

*Proof.* The proof can be obtained directly from some calculation using Claim A.1 as follows: **Proof of (a):** From (A.2), we can obtain that

$$\mathbf{h} - \nabla_{w_1} L(W(t), E(t)), \Pi_{V \perp W1}(t) \mathbf{i} = -\mathbf{X} \Sigma_{(j, t)} \mathbf{h} \nabla_{w_1} \mathbf{E}_{j, (t3) - j, W1}(t) \mathbf{i}$$

$$j, \quad$$

Now from Claim A.1a and Induction C.1a, we know  $(B_{j,\ell}^{(t)})^3 \leq \widetilde{O}(\frac{1}{d^{3/2}})$  and the following

$$\Sigma_{j,\ell}^{(t)} = O(\Sigma_{1,1}^{(t)}) \frac{(B_{j,\ell}^{(t)})^6 + E_{j,3-j}^{(t)}(B_{3-j,\ell}^{(t)})^3(B_{j,\ell}^{(t)})^3}{(B_{1,1}^{(t)})^6} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}} \leq \widetilde{O}(\frac{E_{j,3-j}^{(t)}}{d^{3/2}}) \Sigma_{1,1}^{(t)} \frac{\Phi_j^{(t)}}{\Phi_1^{(t)}} \qquad \text{for any (j, ') 6= (1,1)}$$

From Induction C.1a,c, we know  $((B_{2,\ell}^{(t)})^3 + E_{2,1}^{(t)}(B_{1,\ell}^{(t)})^3)^2 \leq \widetilde{O}(\frac{1}{d^{3/2}})E_{2,1}^{(t)}$  and  $R_2^{(t)} = \Theta(1)$ , which by Claim C.2a,b and Lemma C.3a gives  $\Phi_2^{(t)}/\Phi_1^{(t)} \leq \widetilde{O}(\alpha_1^{O(1)})$ . Combine the bounds above, we can obtain  $\Sigma_{j,\ell}^{(t)} = \widetilde{O}(E_{j,3-j}^{(t)}/d^{3/2})\Sigma_{1,1}^{(t)}$ . We can then directly apply Claim A.1 to prove Lemma C.5a as follows  $\langle -\nabla_{w_1}L(W^{(t)},E^{(t)}),\Pi_{V^\perp}w_1^{(t)}\rangle$ 

**Proof of (b):** For Lemma C.5b, we can use the same analysis for  $\Sigma_{1,1}$  above and Claim A.1(d,e) to get (again we have used  $\Sigma_{j,\ell}^{(t)} = \widetilde{O}(E_{j,3-j}^{(t)})\Sigma_{1,1}^{(t)} = o(\Sigma_{1,1}^{(t)})$ )

$$\begin{split} & \langle -\nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} w_2^{(t)} \rangle \\ &= (1 \pm \widetilde{O}(E_{1,2}^{(t)})) \Sigma_{1,1}^{(t)} \Big( (-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} + E_{1,2}^{(t)} R_1^{(t)} [R_2^{(t)}]^2 \Big) \\ &+ \widetilde{O}(E_{2,1}^{(t)}/d^{3/2}) \Sigma_{1,1}^{(t)} \Big( (-\Theta(\overline{R}_{1,2}^{(t)}) + O(\varrho)) (E_{2,1}^{(t)})^2 [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} + E_{2,1}^{(t)} R_1^{(t)} [R_2^{(t)}]^2 \Big) \\ &= \Sigma_{1,1}^{(t)} ((-\Theta(\overline{R}_{1,2}^{(t)}) + O(\varrho)) [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} + \widetilde{O}(|E_{1,2}^{(t)}| + \frac{|E_{2,1}^{(t)}|^2}{d^{3/2}}) R_1^{(t)} [R_2^{(t)}]^2 ) \end{split}$$

Proof of (c): Similarly to the proof of (a), we can also expand as follows

$$\begin{split} &\langle -\nabla_{w_2} L(W^{(t)}, E^{(t)}), \Pi_{V^{\perp}} w_2^{(t)} \rangle \\ &= (1 \pm O(E_{1,2}^{(t)})) \Sigma_{1,1}^{(t)} \Big( - [R_2^{(t)}]^3 \Theta((E_{1,2}^{(t)})^2) \pm O(E_{1,2}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big) \\ &- \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big( [R_2^{(t)}]^3 \pm O(E_{2,1}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big) \\ &= - [R_2^{(t)}]^3 \Big( \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big) \pm O\Big( \sum_{j,\ell} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big) \end{split}$$

**Proof of (d):** Similarly, we can calculate (again by  $\Sigma_{j,\ell}^{(t)} = \widetilde{O}(E_{j,3-j}^{(t)})\Sigma_{1,1}^{(t)} = o(\Sigma_{1,1}^{(t)})$ )

$$\begin{split} &\langle -\nabla_{w_2}L(W^{(t)},E^{(t)}),\Pi_{V^\perp}w_1^{(t)}\rangle \\ &= \sum_{\ell\in[2]}\Sigma_{1,\ell}^{(t)}\Big((-\Theta(\overline{R}_{1,2}^{(t)})\pm O(\varrho))(E_{1,2}^{(t)})^2[R_2^{(t)}]^{5/2}[R_1^{(t)}]^{1/2} + E_{1,2}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big) \\ &+ \sum_{\ell\in[2]}\Sigma_{2,\ell}^{(t)}\Big((-\Theta(\overline{R}_{1,2}^{(t)})\pm O(\varrho))[R_2^{(t)}]^{5/2}[R_1^{(t)}]^{1/2} + E_{2,1}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big) \\ &= (1\pm\widetilde{O}(E_{1,2}^{(t)}))\Sigma_{1,1}^{(t)}\Big((-\Theta(\overline{R}_{1,2}^{(t)})\pm O(\varrho))(E_{1,2}^{(t)})^2[R_2^{(t)}]^{5/2}[R_1^{(t)}]^{1/2} + E_{1,2}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big) \\ &+ \sum_{\ell\in[2]}\Sigma_{2,\ell}^{(t)}\Big((-\Theta(\overline{R}_{1,2}^{(t)})\pm O(\varrho))[R_2^{(t)}]^{5/2}[R_1^{(t)}]^{1/2} + E_{1,2}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big) \\ &= \Big(\Sigma_{1,1}^{(t)}\Theta((E_{1,2}^{(t)})^2) + \sum_{\ell\in[2]}\Sigma_{2,\ell}^{(t)}\Big)(-\Theta(\overline{R}_{1,2}^{(t)})\pm O(\varrho))[R_2^{(t)}]^{5/2}[R_1^{(t)}]^{5/2}[R_1^{(t)}]^{1/2} + O\Big(\sum_{i,\ell}\Sigma_{j,\ell}^{(t)}E_{j,3-j}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big) \end{split}$$

which completes the proof.

**Lemma C.6** (learning feature  $v_2$  in phase II). For each  $t \in [T_1, T_2]$ , if Induction C.1 holds at iteration t, then we have for each  $j \in [2]$ :

$$|\langle -\nabla_{w_j} L(W^{(t)}, E^{(t)}), v_2 \rangle| \leq \widetilde{O}(\frac{\alpha_2^6 \alpha_1^6}{d^{5/2}}) \left( \Phi_j^{(t)} (|E_{j,3-j}^{(t)}| + [R_j^{(t)}]^3) + \Phi_{3-j}^{(t)} (|E_{3-j,j}^{(t)}| |[R_{3-j}^{(t)}]^3 + \frac{|E_{3-j,j}^{(t)}|^2}{d^{3/2}}) \right)$$

*Proof.* Again as in the proof of Lemma B.9, we expand the notations: (ignoring the superscript <sup>(t)</sup> for the RHS)

$$\langle -\nabla_{w_j} L(W^{(t)}, E^{(t)}), v_2 \rangle = \Lambda_{j,2}^{(t)} + \Gamma_{j,2}^{(t)} - \Upsilon_{j,2}^{(t)}$$
(C.2)

where

$$\begin{split} &\Lambda_{j,2}^{(t)} = C_0 \alpha_2^6 \Phi_j^{(t)} H_{j,1}^{(t)} (B_{j,2}^{(t)})^5 \\ &\Gamma_{j,2}^{(t)} = C_0 \alpha_2^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^2 H_{3-j,1}^{(t)} \\ &\Upsilon_{j,2}^{(t)} = C_0 \alpha_1^6 \left( \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{j,2}^{(t)} + \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{3-j,2}^{(t)} \right) \end{split}$$

Now we further write  $\mathbf{\gamma}_{2}^{(t)}=\mathbf{\Upsilon}_{j,2,1}^{(t)}+\mathbf{\Upsilon}_{j,2,2_{\mathbf{j}}}^{(t)}$  , where

$$\Upsilon_{j,2,1}^{(t)} = C_0 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{j,2}^{(t)} \qquad \Upsilon_{j,2,2}^{(t)} = \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{3-j,2}^{(t)})^2 K_{3-j,2}^{(t)}$$

According to (C.2), we can first compute

$$\begin{split} \Lambda_{j,2}^{(t)} - \Upsilon_{j,2,1}^{(t)} &= C_0 \alpha_2^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^5 H_{j,1}^{(t)} - C_0 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{j,2}^{(t)} \\ &= C_0 \alpha_2^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^5 \left( C_1 \alpha_1^6 ((B_{j,1}^{(t)})^3 + E_{j,3-j}^{(t)} (B_{3-j,1}^{(t)})^3)^2 + C_2 \mathcal{E}_{j,3-j}^{(t)} \right) \\ &- C_0 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 C_1 \alpha_2^6 ((B_{j,2}^{(t)})^3 + E_{j,3-j}^{(t)} (B_{3-j,2}^{(t)})^3) ((B_{j,1}^{(t)})^3 + E_{j,3-j}^{(t)} (B_{3-j,1}^{(t)})^3) \\ &= C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^5 \left( E_{j,3-j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,1}^{(t)})^3 + (E_{j,3-j}^{(t)})^2 (B_{3-j,1}^{(t)})^6 \right) \\ &- C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^2 (B_{3-j,2}^{(t)})^3 E_{j,3-j}^{(t)} \left( (B_{j,1}^{(t)})^6 + E_{j,3-j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,1}^{(t)})^3 \right) \\ &+ C_0 \alpha_2^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^5 C_2 \mathcal{E}_{j,3-j}^{(t)} \end{split}$$

Then we can apply Induction C.1a,c,d, Claim C.2a,b and Lemma C.3a,c to get

$$|\Lambda_{j,2}^{(t)} - \Gamma_{j,2,1}^{(t)}| \le \widetilde{O}(\frac{\alpha_2^6}{\alpha_1^6 d^{5/2}}) \Phi_j^{(t)}(|E_{j,3-j}^{(t)}| + [R_j^{(t)}]^3)$$

where the last inequality is due to Lemma C.3a,c. Similarly, we can also compute for  $\Gamma_{j,2}^{(t)} = \Upsilon_{j,2,2}^{(t)}$ 

$$\begin{split} |\Gamma_{j,2}^{(t)} - \Upsilon_{j,2,2}^{(t)}| &\leq \left| C_0 \alpha_2^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^2 H_{3-j,1}^{(t)} \right| \\ &+ \left| C_0 \alpha_1^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{3-j,2}^{(t)} \right| \\ &\leq \widetilde{O}(\frac{\alpha_1^6 \alpha_2^6}{d^{5/2}}) \Phi_{3-j}^{(t)} |E_{3-j,j}^{(t)}| ([R_{3-j}^{(t)}]^3 + \frac{|E_{3-j,j}^{(t)}|}{d^{3/2}}) \end{split}$$

This completes the proof  $\Box$  **Lemma C.7** (learning feature  $v_1$  in Phase II). For each  $t \in [T_1, T_2]$ , if Induction C.1 holds at iteration t, then we have:

(a) 
$$\langle -\nabla_{w_1} L(W^{(t)}, E^{(t)}), v_1 \rangle = \Theta(\Sigma_{1,1}^{(t)}) [R_1^{(t)}]^3 + \Gamma_{1,1}^{(t)} \pm \widetilde{O}(\alpha_1^{O(1)}/d^{5/2}),$$
  
(b)  $\langle -\nabla_{w_2} L(W^{(t)}, E^{(t)}), v_1 \rangle = \widetilde{O}(\alpha_1^{O(1)}/d^{5/2}) + \widetilde{O}(\frac{\alpha_1^6}{d}) E_{1,2}^{(t)} \Phi_1^{(t)} [R_1^{(t)}]^3$ 

*Proof.* As in the proof of Lemma C.6, we expand the gradient terms:

$$\langle -\nabla_{w_j} L(W^{(t)}, E^{(t)}), v_1 \rangle = \Lambda_{j,2}^{(t)} + \Gamma_{j,2}^{(t)} - \Upsilon_{j,2}^{(t)}$$
(C.3)

where

$$\begin{split} & \Lambda_{j,1}^{(t)} = C_0 \alpha_1^6 \Phi_j^{(t)} H_{j,2}^{(t)} (B_{j,1}^{(t)})^5 \\ & \Gamma_{j,1}^{(t)} = C_0 \alpha_1^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,1}^{(t)})^2 H_{3-j,2}^{(t)} \\ & \Upsilon_{j,1}^{(t)} = C_0 \alpha_1^6 \left( \Phi_j^{(t)} (B_{j,2}^{(t)})^3 (B_{j,1}^{(t)})^2 K_{j,1}^{(t)} + \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,1}^{(t)})^2 K_{3-j,1}^{(t)} \right) \end{split}$$

Indeed, when j = 1, by Induction C.1a and Lemma C.3a,c, we can compute

$$\Lambda_{1,1}^{(t)} = C_0 \alpha_1^6 \Phi_1^{(t)} (B_{1,1}^{(t)})^5 H_{1,2}^{(t)} = \Theta(\Sigma_{1,1}^{(t)}) [R_1^{(t)}]^3$$

and with additionally Lemma C.3b, we also have

$$|\Upsilon_{1,1}^{(t)}| = \left| C_0 \alpha_1^6 \left( \Phi_1^{(t)} (B_{1,2}^{(t)})^3 (B_{1,1}^{(t)})^2 K_{1,1}^{(t)} + \Phi_2^{(t)} E_{2,j}^{(t)} (B_{2,2}^{(t)})^3 (B_{1,1}^{(t)})^2 K_{2,1}^{(t)} \right) \right| \le \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{5/2}})$$

which gives the proof of (a). For (b), we can also apply Induction C.1a and Lemma C.3a,c to get

$$\begin{split} &\Lambda_{2,1}^{(t)} = C_0 \alpha_1^6 \Phi_2^{(t)} H_{2,2}^{(t)} (B_{2,1}^{(t)})^5 \leq \widetilde{O}(\alpha_1^{O(1)}/d^{5/2}) \\ &\Gamma_{2,1}^{(t)} = C_0 \alpha_1^6 \Phi_1^{(t)} E_{1,2}^{(t)} (B_{1,1}^{(t)})^3 (B_{2,1}^{(t)})^2 H_{1,2}^{(t)} \leq \widetilde{O}(\frac{1}{d}) E_{1,2}^{(t)} \Phi_1^{(t)} \frac{[R_1^{(t)}]^3}{\alpha_1^6} \\ &\Upsilon_{2,1}^{(t)} = C_0 \alpha_1^6 \left(\Phi_2^{(t)} (B_{2,2}^{(t)})^3 (B_{2,1}^{(t)})^2 K_{2,1}^{(t)} + \Phi_1^{(t)} E_{1,2}^{(t)} (B_{1,2}^{(t)})^3 (B_{2,1}^{(t)})^2 K_{1,1}^{(t)}\right) \leq \widetilde{O}(\frac{\alpha_1^6}{d^4}) \end{split}$$

this finishes the proof.

### C.3 At the End of Phase II

Now we shall present the main theorem of this section, which gives the result of prediction head  $E_{2,1}^{(t)}$  growth after the feature  $v_1$  is learned in the first stage.

**Lemma C.8** (Phase II). Suppose  $\eta = \frac{1}{\text{poly}(d)}$  is sufficiently small, then Induction C.1 holds for all iteration  $t \in [T_1, T_2]$ , and at iteration  $t = T_2$ , the followings holds:

(a) 
$$B_{1,1}^{(T_2)} = \Theta(1)$$
,  $B_{j,\ell}^{(T_2)} = B_{j,\ell}^{(T_1)}(1 \pm o(1)) = \widetilde{\Theta}(\frac{1}{\sqrt{d}})$  for (j, ') 6= (1,1)

(b) 
$$R_1^{(T_2)} \leq \widetilde{O}(\frac{1}{d^{3/4}}), \ R_2^{(T_2)} = \Theta(\sqrt{\eta_E/\eta}), \ \text{and} \ \overline{R}_{1,2}^{(T_2)} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}), \ \overline{R}_{1,2}^{(T_2)} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

$$(c) \ |E_{1,2}^{(T_2)}| = \widetilde{O}(\varrho + \tfrac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2} \\ \mathit{and} |E_{2,1}^{(T_2)}| = \Theta(\sqrt{\eta_E/\eta})$$

Where the part of learning  $E_{2,1}^{(t)}$  is what we called substitution effect. One can easily verify that  $|E_{2,1}^{(t)}f_1(X^{(1)})| \gg |f_2(X^{(1)})|$  when X is equipped with feature  $v_1$ , as stated in Lemma 5.2.

*Proof.* We first will prove Induction C.1 holds for all iteration  $t \in [T_1, T_2]$ . We shall first prove that if Induction C.1 continues to hold when  $R_2^{(t)} \ge |E_{2,1}^{(t)}|$ , we shall have  $[R_1^{(t)}]$  decreasing at an exponential rate.

**Proof of the decrease of** $R_1^{(t)}$ : Firstly, we write down the update of  $R_1^{(t)}$  using Lemma C.5a:

$$R_1^{(t+1)} = R_1^{(t)} + \eta \Sigma_{1,1}^{(t)} \Theta(-[R_1^{(t)}]^3 \pm O(|E_{1,2}^{(t)}| + \frac{|E_{2,1}^{(t)}|^2}{d^{3/2}}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2})$$

(t) from the expression of  $\Sigma_{1,1}$  in (A.2), and by Induction C.1a and Lemma

C.3a,c, we can compute

$$\Sigma_{1,1}^{(t)} = \Theta(C_0 C_2 \Phi_1^{(t)}) = \Theta(\frac{C_0 C_2}{\alpha_1^{12}})$$

Moreover, from Induction C.1c we know that

$$\begin{split} (|E_{1,2}^{(t)}| + \frac{|E_{2,1}^{(t)}|^2}{d^{3/2}})[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2} &\leq (\widetilde{\Theta}(\frac{1}{d^{3/2}}) + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2})[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2} \\ &\leq (\widetilde{\Theta}(\frac{1}{d^{3/2}}) + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2})[R_1^{(t)}]^{3/2} \end{split}$$

Therefore whenever  $R_1^{(t)} \ge \frac{\alpha_1^{18}}{d^{3/4}}$  (which  $t \le T_2$  suffices), we shall have always have  $(\overline{R}_{1,2}^{(t)} + \varrho)(\widetilde{\Theta}(\frac{1}{d^{3/2}}) + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2})[R_1^{(t)}]^{3/2} \le o([R_1^{(t)}]^3)$ 

which implies, if we set  $T_2' := \min\{t : R_1^{(t)} \ge \frac{1}{d^{3/4}\alpha_1^2}\}$ , then for all  $t \in [T_1, T_2']$ , we will have

$$\begin{split} R_1^{(t+1)} &= R_1^{(t)} + \eta \Sigma_{1,1}^{(t)} \Theta(-[R_1^{(t)}]^3 \pm O(|E_{1,2}^{(t)}| + \frac{|E_{2,1}^{(t)}|^2}{d^{3/2}}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2}) \\ &= R_1^{(t)} - \Theta(\eta \Sigma_{1,1}^{(t)}) [R_1^{(t)}]^3 \end{split} \tag{C.4}$$

$$\leq R_1^{(t)} (1 - \Theta(\frac{\eta C_0^1 C_2}{\alpha^{12}}) \frac{1}{d^{3/2} \alpha^2} \qquad \qquad R_1^{(t)} \geq \frac{1}{d^{3/4}})$$

(since

From the last inequality we know that after  $T_2=T_1+\widetilde{\Theta}(rac{d^{1.5}}{\etalpha_1^{\Omega(1)}})$ , we shall have  $R_1^{(t)}\leq O(rac{lpha_1^{O(1)}}{d^{3/4}})$ .

Moreover, suppose  $T_2' < T_2$ , (which just mean  $R_1^{(s)} \le O(\frac{1}{d^{3/4}\alpha_1^2})$  for some iteration  $s \in [T_1, T_2]$ ) we also have

$$\begin{split} R_1^{(t+1)} &= R_1^{(t)} - \Theta(\eta \Sigma_{1,1}^{(t)}) [R_1^{(t)}]^3 \\ &\geq R_1^{(t)} (1 - \Theta(\frac{\eta C_0 C_2}{\alpha_1^{14}}) \frac{1}{d^{3/2}}) \end{split}$$

So when  $T_2 \leq T_1 + \widetilde{O}(\frac{d^{1.5}\alpha_1^{12}}{\eta})$  iterations, we will have  $R_1^{(t)} \geq R_1^{(s)}(1 - \Theta(\frac{\eta C_0 C_2}{d^{3/2}\alpha_1^{14}}))^{T_2 - T_1} \geq \Omega(R_1^{(t)})$  for all  $t \in [s, T_2]$ , which means we have a lower bound  $R_1^{(t)} \geq \frac{1}{d^{3/4}\alpha_1^2}$  throughout  $t \in [T_1, T_2]$ . This proves Lemma C.8a and also our induction on  $R_1^{(t)}$ .

**Proof of induction for**  $E_{1,2}^{(t)}$ : By Lemma C.4a, we can write

$$\begin{split} -\nabla_{E_{1,2}}L(W^{(t)},E^{(t)}) &= (1+\widetilde{O}(\frac{\alpha_{1}^{O(1)}}{d^{3/2}}))\Sigma_{1,1}^{(t)}(-2E_{1,2}^{(t)}[R_{2}^{(t)}]^{3} \pm O(\overline{R}_{1,2}^{(t)}+\varrho)[R_{1}^{(t)}]^{3/2}[R_{2}^{(t)}]^{3/2}) \\ &\pm \Sigma_{1,1}^{(t)}\widetilde{O}(\frac{\eta_{E}/\eta}{\sqrt{d}})\max\{[R_{1}^{(t)}]^{3},\frac{\alpha_{1}^{O(1)}}{d^{5/2}}\} \\ &= -\Theta(\Sigma_{1,1}^{(t)}[R_{2}^{(t)}]^{3})E_{1,2}^{(t)} \pm O(\Sigma_{1,1}^{(t)})\Big((\overline{R}_{1,2}^{(t)}+\varrho)[R_{1}^{(t)}]^{3/2}[R_{2}^{(t)}]^{3/2} + \widetilde{O}(\frac{\eta_{E}/\eta}{\sqrt{d}})[R_{1}^{(t)}]^{3}\Big) \end{split}$$

Since again from Induction C.1b,c that

$$R_{1,2}^{(t)} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}), R_1^{(t)} = O(1), R_2^{(t)} \in [\sqrt{\eta_E/\eta}, O_{\text{(1)}}],$$
 we

can obtain the update of  $E_{1,2}^{(t)}$  as

$$\begin{split} E_{1,2}^{(t+1)} &= E_{1,2}^{(t)} (1 - \Theta(\eta_E \Sigma_{1,1}^{(t)}[R_2^{(t)}]^3)) \pm \widetilde{O}(\eta_E \Sigma_{1,1}^{(t)}) \Big( (\varrho + \frac{1}{\sqrt{d}}) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} + \widetilde{O}(\frac{\eta_E/\eta}{\sqrt{d}}) [R_1^{(t)}]^3 \Big) \\ &= E_{1,2}^{(t)} (1 - \Theta(\eta_E \Sigma_{1,1}^{(t)}[R_2^{(t)}]^3)) \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \eta_E \Sigma_{1,1}^{(t)} [R_1^{(t)}]^{3/2} \\ &= E_{1,2}^{(t)} (1 - \Theta(\eta_E \Sigma_{1,1}^{(t)}[R_2^{(t)}]^3)) \pm \eta_E \Sigma_{1,1}^{(t)} J_{1,2}^{(t)} \end{split}$$

where  $J_{1,2}^{(t)}=\widetilde{C}(\varrho+\frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2}>_0$  and  $C_{\mathbf{e}}=\Theta(1)_{\mathbf{e}}$  is larger than the hidden constant (including the polylog(d) factors) of  $E_{2,1}^{(T_1)}\leq \widetilde{C}(\varrho+\frac{1}{\sqrt{d}})$  in Lemma B.13d. And then we can compute

$$\begin{split} J_{1,2}^{(t+1)} &= \widetilde{C}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t+1)}]^{3/2} \\ &= \widetilde{C}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2}(1 - \Theta(\eta \Sigma_{1,1}^{(t)})[R_1^{(t)}]^2)^{3/2} \\ &= J_{1,2}^{(t)}(1 - \Theta(\eta^{3/2}(\Sigma_{1,1}^{(t)})^{3/2})[R_1^{(t)}]^3) \end{split} \qquad \text{(due to calculations in (C.4))} \\ &= J_{1,2}^{(t)}(1 - \Theta(\eta^{3/2}(\Sigma_{1,1}^{(t)})^{3/2})[R_1^{(t)}]^3) \end{split} \qquad \text{(because } \eta \Sigma_{1,1}^{(t)} = \frac{\alpha_1^{O(1)}}{\text{poly}(d) \text{ is very small}} \end{split}$$

Now by Lemma B.13d, we know  $|E_{1,2}^{(T_1)}| \leq J_{1,2}^{(T_1)}$ ; then we begin our induction that  $|E_{1,2}^{(t)}| < (\log\log d)J_{1,2}^{(t)}$  at for all iterations  $t\in [T_1,T_2]$ . Now assume we have  $|E_{1,2}^{(t)}|=\frac{1}{2}(\log\log d)J_{1,2}^{(t)4}$ , from above calculations it holds that  $|E_{1,2}^{(t+1)}|=|E_{1,2}^{(t)}|(1-\Theta(\eta\Sigma_{1,1}^{(t)}[R_1^{(t)}]^3))$ . Then we would have

$$\frac{J_{1,2}^{(t+1)}}{J_{1,2}^{(t)}} \ge \left(1 - \Theta(\eta^{3/2}(\Sigma_{1,1}^{(t)})^{3/2})[R_1^{(t)}]^3\right) \ge \left(1 - \Theta(\eta_E \Sigma_{1,1}^{(t)}[R_2^{(t)}]^3)\right) \ge \frac{|E_{1,2}^{(t+1)}|}{|E_{1,2}^{(t)}|}$$

(because of the range of  $R_1^{(t)}$  and  $R_2^{(t)})$ 

This proved that  $|E_{1,2}^{(t+1)}|\lesssim \log\log d\cdot J_{1,2}^{(t+1)}\leq \widetilde{O}(\varrho+\frac{1}{\sqrt{d}})[R_1^{(t+1)}]^{3/2}$  and also the induction can go

**Proof of the growth of**  $E_{2,1}^{(t)}$  and  $T_2 \leq T_1 + O(\frac{d^{1.5}}{\eta \alpha_1^4})$ : According to Lemma C.4b, we can write down the update of  $E_{2,1}^{(t)}$  as

<sup>4</sup>If we want  $|E_{1,2}^{(t)}| > (\log \log d)J_{1,2}^{(t)}$  then as long as  $\eta = \frac{1}{\text{poly}(d)}$  is small enough, we can always assume to have found some iteration  $t^0 \in (T_1,t]$  such that  $|E_{1,2}^{(t')}| = \frac{1}{2}(\log \log d)J_{1,2}^{(t)}$ , and we set  $t=t^0$  and start our argument from that iteration. on until  $t=T_2$ .

$$\begin{split} -\nabla_{E_{2,1}}L(W^{(t)},E^{(t)}) &= (1\pm O(\frac{\alpha_1^{O(1)}}{d^{3/2}}))\Delta_{2,1}^{(t)} \\ &\pm O(\Sigma_{2,1}^{(t)})(|E_{2,1}^{(t)}|[R_1^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)}+\varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}) \end{split}$$

Then, from Lemma C.3a,c and Induction C.1, we have

$$O(\Sigma_{2,1}^{(t)})(|E_{2,1}^{(t)}|[R_1^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}) \leq O \\ \text{polylog } (\underbrace{ \begin{pmatrix} (d) \\ / \end{pmatrix} \Phi_2^{(t)} \leq O( \begin{pmatrix} 1 \\ / \end{pmatrix} \Phi_2^{(t)} \geq O( \begin{pmatrix} 1 \\ / \end{pmatrix} \Phi_2$$

and also

C.2)

$$\left| (1 \pm \widetilde{O}(\frac{\alpha_1^6}{d^{0.3}})) C_0 \Phi_2^{(t)} \alpha_1^6 (B_{2,1}^{(t)})^3 (B_{1,1}^{(t)})^3 H_{2,2}^{(t)} \right| \ge \widetilde{\Theta}(\frac{\alpha_1^6}{d^{3/2}}) \Phi_2^{(t)}$$

Now by Lemma C.3a and Induction C.1a, it allow us to simplify the update to

$$\begin{split} E_{2,1}^{(t+1)} &= E_{2,1}^{(t)} - \eta_E \nabla_{E_{2,1}} L(W^{(t)}, E^{(t)}) \\ &= E_{2,1}^{(t)} + (1 \pm \frac{1}{\alpha_1^{\Omega(1)}}) \eta_E C_0 C_2 \alpha_1^6 \Phi_2^{(t)} (B_{2,1}^{(t)})^3 (B_{1,1}^{(t)})^3 \mathcal{E}_{2,1}^{(t)} \\ &\geq E_{2,1}^{(t)} + \eta_E \widetilde{\Theta}(\frac{1}{d^{3/2} \alpha_1^6}) \mathrm{sign}(B_{1,1}^{(t)}) \mathrm{sign}(B_{2,1}^{(t)}) \end{split} \tag{by Induction C.1 and Claim}$$

Now since  $\text{sign}(B_{j,1}^{(t)}) = \text{sign}(B_{j,1}^{(T_1)})$ , we know there is an iteration  $T'_{2,1} \leq T_1 + O(\frac{d^{1/2}\alpha_1^{O(1)}}{\eta})$  such that for all  $t \in [T'_{2,1}, T_2]$ , it holds

$$|E_{2,1}^{(t)}| = \left| E_{2,1}^{(T_1)} + \sum_{t \in [T_1, T'_{2,1}]} \Theta(\eta_E C_0 C_2 \alpha_1^6) \Phi_2^{(t)} (B_{2,1}^{(t)})^3 (B_{1,1}^{(t)})^3 [R_2^{(t)}]^3 \right|$$

$$= \left| |E_{2,1}^{(T_1)}| \pm \sum_{s \in [T_1, T'_{2,1}]} \eta_E \widetilde{\Theta}(\frac{1}{d^{3/2} \alpha_1^{O(1)}}) \right|$$

$$\in \left[ 2|E_{2,1}^{(T_1)}|, \widetilde{O}(\frac{\alpha_1^{O(1)}}{d}) \right]$$

and thus  $\operatorname{sign}(E_{2,1}^{(t)}) = \prod_{j \in [2]} \operatorname{sign}(B_{j,1}^{(t)})$  and  $|E_{2,1}^{(t)}|$  will be increasing during  $t \in [T_{2,1}', T_2]$ . Thus as long as  $R_2^{(t)} \geq |E_{2,1}^{(t)}|$  continues to hold, after at most  $\widetilde{\Theta}(\frac{d^{1.5}}{\eta \alpha_1^6})$  iterations starting from  $T_1$ , we shall have  $|E_{2,1}^{(t)}| \geq \Omega(\sqrt{\eta_E/\eta})$ .

However, in order to actually prove  $|E_{2,1}^{(T_2)}| = \Theta(\sqrt{\eta_E/\eta})$ , we will need to ensure that (1) there exist some constant  $C = \Omega(p_{\eta_E/\eta})$  such that  $|E_{2,1}^{(t)}| > C$  while  $R_2^{(s)} \ge \frac{1}{\log d} |E_{2,1}^{(t)}|$  for all  $s \in [T_1, t]$ ; (2) we shall have a upper bound  $|E_{2,1}^{(t)}| < O(\sqrt{\eta_E/\eta})$ . They will be done below.

**Proof of**  $E_{2,1}^{(T_2)} = \Theta(\sqrt{\eta_E/\eta})$  and  $T_2 = T_1 + \widetilde{O}(\frac{d^{3/2}\alpha_1^{O(1)}}{\eta})$ : In fact, Induction C.1c are already proved since we have already calculated the dynamics of  $R_1^{(t)}$  and its upper bound and lower bound. In this part we are going to  $\operatorname{prove}^{T_2 = T_1 + \widetilde{\Theta}(\frac{d^{1.5}\alpha_1^{12}}{\eta})}$  (which means  $\operatorname{that} R_2^{(t)} \leq |E_{2,1}|$  can be achieved in  $\widetilde{O}(\frac{d^{3/2}\alpha_1^{12}}{\eta})$  many iterations). From Lemma C.5c, we can write down the update for  $R_2^{(t)}$ 

$$\begin{split} R_2^{(t+1)} &= R_2^{(t)} - 2\eta \langle \nabla_{w_2} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} w_2^{(t)} \rangle + \eta^2 \|\Pi_{V^\perp} \nabla_{w_2} L(W^{(t)}, E^{(t)})\|_2^2 \\ &= R_2^{(t)} - \eta \Theta([R_2^{(t)}]^3) \Big( \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big) \end{split}$$

$$\pm \, \eta O\Big(\sum_{j,\ell} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)}(R_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}\Big) + \frac{\eta}{\mathsf{poly}(d)} \underline{\qquad} \text{the fact that } 1 + \eta = 0$$

 $\|\Pi_{V^{\perp}}\nabla_{w_2}L(W^{(t)},E^{(t)})\|_2^2 \leq \widetilde{\widetilde{O}}(d^2)$  from our assumption on the

(t) noise  $\xi_p$  and a simple bound for  $\Sigma_{i,j}$  as we have done before. Next

we can resort to Induction C.1d

$$\begin{aligned} \mathsf{that}^{|E_{1,2}^{(t)}|} & \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2} \, \mathsf{to} \, \mathsf{derive} \\ & \sum_{s \in [T_1,t]} \eta \Sigma_{1,1}^{(s)} \Theta((E_{1,2}^{(s)})^2) \leq \sum_{s \in [T_1,t]} \widetilde{O}(\varrho^2 + \frac{1}{d}) \eta \Sigma_{1,1}^{(s)}[R_1^{(s)}]^3 \end{aligned}$$

$$\leq \widetilde{O}(\varrho^2 + \frac{1}{d}) = o(1)$$

which is because  $\sum_{t \in [T_1, T_2]} \Theta(\eta \Sigma_{1,1}^{(t)}) [R_1^{(t)}]^3 \leq O_{(1)}$  and  $\Sigma = 0$  as we have calculated in the proof of Induction C.1a above. Similarly, we can also bound

$$\sum_{s \in [T_1,t]} \Sigma_{1,\ell}^{(s)} |E_{1,2}^{(s)}| (|\overline{R}_{1,2}^{(s)}| + \varrho) [R_1^{(s)}]^{3/2} [R_2^{(s)}]^{3/2} \leq \sum_{s \in [T_1,t]} \widetilde{O}(\varrho^2 + \frac{1}{d}) \eta \Sigma_{1,\ell}^{(s)} [R_1^{(s)}]^3 \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) = o(1)$$

Moreover, because  $T_2 \leq T_1 + \widetilde{O}(\frac{d^{3/2}\alpha_1^{12}}{\eta})$  and  $|E_{2,1}^{(t)}| \leq O(1), \ \Phi_2^{(t)} \leq \alpha_1^{O(1)}$  from Induction C.1, we

$$\begin{split} \sum_{s \in [T_1, t]} \eta \Sigma_{2, \ell}^{(s)} |E_{2, 1}^{(s)}| (|\overline{R}_{1, 2}^{(s)}| + \varrho) [R_1^{(s)}]^{3/2} [R_2^{(s)}]^{3/2} &\leq \widetilde{O}(\frac{|E_{2, 1}^{(s)}|^2}{d^{3/2}}) \sum_{s \in [T_1, t]} \eta \Phi_2^{(s)} \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \\ &\leq \widetilde{O}(\frac{\eta}{d^{3/2}}) \cdot \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \cdot \widetilde{O}(\frac{d^{3/2} \alpha_1^{12}}{\eta}) \\ &\leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \alpha_1^{O(1)} = o(1) \end{split}$$

Thus combining all the bounds above, we have proved that for each  $t \in [T_1, T_2]$ , it holds

$$R2(t) = R2(T_1) - X \Theta(\eta \Sigma(2t,1))[R2(t)] 3 \pm o(1)$$

$$s \in [T_1,t]$$

$$= R_2^{(T_1)} - \sum_{s \in [T_1, t]} \Theta(\eta C_0 C_2) E_{2,1}^{(t)} \alpha_1^6 \Phi_2^{(t)} (B_{2,1}^{(t)})^3 (B_{1,1}^{(t)})^3 [R_2^{(t)}]^3 \pm o(1)$$
(C.5)

$$=R_2^{(T_1)} - \sum_{s \in [T_1, t]} \eta E_{2,1}^{(t)} \widetilde{\Theta}(\frac{1}{d^{3/2}}) \Phi_2^{(t)} [R_2^{(t)}]^3 \cdot \operatorname{sign}(E_{2,1}^{(t)}) \cdot \operatorname{sign}(B_{2,1}^{(T_1)}) \cdot \operatorname{sign}(B_{1,1}^{(T_1)}) \pm o(1) \quad (C.6)$$

where the last equality is because  $\operatorname{sign}(B_{j,\ell}^{(t)}) \equiv \operatorname{sign}(B_{j,\ell}^{(T_1)})$  by Induction C.1a. Now from what we have proved above on the growth of  $E_{2,1}^{(t)}$  that  $\operatorname{sign}(E_{2,1}^{(t)}) = \operatorname{sign}(B_{1,1}^{(t)}B_{2,1}^{(t)}) \equiv \operatorname{sign}(B_{1,1}^{(T_1)}B_{2,1}^{(T_1)})$  throughout the rest of phase II (which is just  $t \in [T_{2,1}^{t}, T_{2}]$ ). Recall that

$$R_2^{(T_{2,1}')} = R_2^{(T_1)} \pm o(1) \qquad \qquad \text{(t)} \qquad \qquad \text{(T20,1)} \, \mathbf{X} \qquad \qquad \text{(s)} \qquad \text{(s)} \quad \mathbf{3} \qquad \mathbf{3$$

The above arguments implie for  $t \in [T'_{2,1}, T_2]$ ::

$$R_{2}^{(t+1)} = R_{2}^{(T_{1})} - \sum_{s \in [T'_{2,1},t]} \Theta(\eta C_{0}C_{2}) E_{2,1}^{(s)} \Phi_{2}^{(s)} (B_{2,1}^{(s)})^{3} (B_{1,1}^{(s)})^{3} [R_{2}^{(t)}]^{3} \pm o(1)$$

$$= R_{2}^{(T_{1})} - \Theta(\frac{\eta}{\eta_{E}} |E_{2,1}^{(t)}|^{2}) - o(1)$$

Now we can confirm

(1) there exist a constant  $C = \Theta(p_{\eta_E/\eta})$  such that  $E_{2,1}^{(t)} = C$  if  $R_2^{(t)}$  falls below  $\frac{1}{\log d}|E_{2,1}^{(t)}|$ ;

$$(2) \ \ T_2 = T_1 + \widetilde{\Theta}\big(\frac{d^{3/2}\alpha_1^{12}}{\eta}\big) \ \text{due to the growth} \ |E_{2,1}^{(t+1)}| = |E_{2,1}^{(t)}| + \eta_E \widetilde{\Theta}\big(\frac{1}{d^{3/2}\alpha_1^{12}\sqrt{\eta_E/\eta}}\big) \ \text{for} \ t \in [T_{2,1}', T_2].$$
 which are the desired results.

**Proof of Induction C.1a:** We first obtain from Lemma C.7a that the update of  $B_{1,1}^{(t)}$  can be written as

$$B_{1,1}^{(t+1)} = B_{1,1}^{(t)} + \eta \left( \Theta(\Sigma_{1,1}^{(t)}) \operatorname{sign}(B_{1,1}^{(t)}) [R_1^{(t)}]^3 + \Gamma_{1,1}^{(t)} \pm \widetilde{O}(\alpha_1^{O(1)}/d^{5/2}) \right)$$

Now by what we have calculated above in (C.4), the total decrease of  $R_1^{(t)}$  is (since  $R_1^{(t)}$  is monotone in this phase)

$$X \Theta(\eta \Sigma_{(1t,1)})[R_{1(t)}] \le O(R_{1(T^{1})} - R_{1(T^{2})}) \le O(1)$$

$$t \in [T_{1}, T_{2}]$$

And also since  $T_2 \leq T_1 + \widetilde{\Theta}(rac{d^{3/2}lpha_1^{12}}{\eta})$  , we can bound

$$\sum_{t \in [T_1, T_2]} \widetilde{O}(\alpha_1^6/d^{5/2}) \leq \widetilde{O}(\alpha_1^{O(1)}/d^{5/2}) \cdot \widetilde{O}(\frac{d^{3/2}}{\eta \alpha_1^6}) \leq \widetilde{O}(\alpha_1^{O(1)}/d)$$

(

Now we consider how the  $\Gamma_{1,1}$  term accumulates

$$\sum_{t \in [T_{1}, T_{2}]} \eta \Gamma_{1,1}^{(t)} = \left( \sum_{t \in [T_{1}, T_{2,1}']} + \sum_{t \in [T_{2,1}', T_{2}]} \right) \eta C_{0} \alpha_{1}^{6} E_{2,1}^{(t)} \Phi_{2}^{(t)} (B_{2,1}^{(t)})^{3} (B_{1,1}^{(t)})^{2} H_{2,2}^{(t)}$$

$$\stackrel{@}{=} \widetilde{O}(\frac{\alpha_{1}^{12}}{d}) + \sum_{t \in [T_{2,1}', T_{2}]} O\left( \eta C_{0} \alpha_{1}^{6} \Phi_{2}^{(t)} |B_{2,1}^{(t)}|^{3} |B_{1,1}^{(t)}|^{3} H_{2,2}^{(t)} \right) \operatorname{sign}(B_{1,1}^{(t)})$$

$$= \pm o(1) + O(1) \operatorname{sign}(B_{1,1}^{(t)})$$

where in ¬ we have used  $|E_{2,1}^{(t)}| \leq O(1) \leq O(B_{1,1}^{(t)})$  and  $\operatorname{sign}(E_{2,1}^{(t)}) = \prod_{j \in [2]} \operatorname{sign}(B_{j,1}^{(t)})$  when  $t \in [T_{2,1}, T_{2}]$ . These calculations tell us  $B_{1,1}^{(t)} = B_{1,1}^{(T_{1})} + O(1)\operatorname{sign}(B_{1,1}^{(T_{1})}) \pm O(\frac{1}{\alpha_{1}}) = \Theta(1)$  for all iterations  $t \in [T_{1}, T_{2}]$ . Similarly from Lemma C.7b, for  $B_{2,1}^{(t)}$  we can also write

$$B_{2,1}^{(T+1)} = B_{2,1}^{(t)} + \eta \widetilde{O}(\alpha_1^{O(1)}/d^{5/2}) + \widetilde{O}(\frac{\alpha_1^6}{d}) E_{2,1}^{(t)} \Phi_1^{(t)} [R_1^{(t)}]^3$$

From similar calculations, it holds  $B_{2,1}^{(t)}=B_{2,1}^{(T_1)}\pm\widetilde{O}(\alpha_1^{O(1)}/d)$ , which proves that  $B_{2,1}^{(t)}=B_{2,1}^{(T_1)}(1\pm o(1))$  when  $t\in [T_1,T_2]$ . Now we turn to feature  $v_2$ . By Lemma C.6 we have for  $j\in [2]$ :

$$\begin{split} |\langle -\nabla_{w_{j}}L(W^{(t)}, E^{(t)}), v_{2}\rangle| &\leq \widetilde{O}(\frac{\alpha_{2}^{6}\alpha_{1}^{6}}{d^{5/2}}) \Big(\Phi_{j}^{(t)}(|E_{j,3-j}^{(t)}| + [R_{j}^{(t)}]^{3}) + \Phi_{3-j}^{(t)}(|E_{3-j,j}^{(t)}|[R_{3-j}^{(t)}]^{3} + \frac{|E_{3-j,j}^{(t)}|^{2}}{d^{3/2}})\Big) \\ &\leq \widetilde{O}(\frac{\alpha_{2}^{6}\alpha_{1}^{6}}{d^{5/2}}) \end{split}$$

where the last inequality is from Lemma C.3a and Induction C.1c,d. Thus when  $t \le T_2 = T_1 + \widetilde{O}(\frac{d^{3/2}\alpha_1^{12}}{\eta})$  we would have

$$B_{j,2}^{(t)} = B_{j,2}^{(T_1)} \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d}) = B_{j,2}^{(T_1)}(1 \pm o_{\text{(1)}}) \qquad \qquad \text{since } B_{j,2}^{(T_1)} = \widetilde{\Theta}(\frac{1}{\sqrt{d}}) \text{ by Lemma B.13c}$$

Together they proved Induction C.1a and Lemma C.8a. Moreover, we have also

$$\begin{split} \textbf{Proof of Induction C.1b:} & \text{Firstly, we write down the update of } \frac{R_{1,2}^{(t)}}{R_{1,2}^{(t)}} \text{ using Lemma C.5b,d as follows:} \\ & R_{1,2}^{(t+1)} = R_{1,2}^{(t)} - \eta \langle \nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} w_2^{(t)} \rangle - \eta \langle \nabla_{w_2} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} w_1^{(t)} \rangle \\ & \qquad + \eta^2 \langle \Pi_{V^\perp} \nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} \nabla_{w_2} L(W^{(t)}, E^{(t)}) \rangle \\ & = R_{1,2}^{(t)} + \eta \Sigma_{1,1}^{(t)} ((-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} + \widetilde{O}(|E_{1,2}^{(t)}| + \frac{|E_{2,1}^{(t)}|^2}{d^{3/2}}) R_1^{(t)} [R_2^{(t)}]^2) \\ & \qquad + \eta \Big( \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [n]} \Sigma_{2,\ell}^{(t)} \Big) (-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_2^{(t)}]^{5/2} [R_1^{(t)}]^{1/2} \end{split}$$

where in 
$$+O\Big(\sum_{j,\ell}\eta\Sigma_{j,\ell}^{(t)}E_{j,3-j}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big)+\frac{\eta}{\mathsf{poly}(d)}\frac{\mathsf{the \ last \ inequality \ we \ have \ used}}{|\mathsf{h}\Pi v_{\perp}\nabla_{\mathsf{W}^1}L\big(\mathsf{W}(\mathsf{t}),\!E(\mathsf{t})\big),\!\Pi v_{\perp}\nabla_{\mathsf{W}^2}L\big(\mathsf{W}(\mathsf{t}),\!E(\mathsf{t})\big)\mathsf{i}|}$$

 $\leq k \prod V_{\perp} \nabla_{W^{1}} L(W(t), E(t)) k_{2} k \prod V_{\perp} \nabla_{W^{2}} L(W(t), E(t)) k_{2} \leq Oe(d)$ 

Now from Induction C.1c,d that  $R_2^{(t)} = \Theta(1)$  and  $|E_{1,2}^{(t)}| \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2}, \ |E_{2,1}^{(t)}| \leq O(\sqrt{\eta_E/\eta}), \text{ we can further obtain} |\Sigma_{2,2}^{(t)}| = \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{3/2}})|\Sigma_{2,1}^{(t)}|, \text{ and the bound}$ 

$$\begin{split} R_{1,2}^{(t+1)} &= R_{1,2}^{(t)} \Big( 1 - \Theta(\eta \Sigma_{1,1}^{(t)}) [R_1^{(t)}]^2 - \Theta(\eta (\Sigma_{1,1}^{(t)} (E_{1,2}^{(t)})^2 + \Sigma_{2,1}^{(t)})) [R_2^{(t)}]^2 \Big) \\ & \pm \eta O(\varrho) [R_2^{(t)}]^{1/2} [R_1^{(t)}]^{1/2} \left( O(\Sigma_{1,1}^{(t)}) [R_1^{(t)}]^2 + \left( \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \Sigma_{2,1}^{(t)} \right) [R_2^{(t)}]^2 \right) \end{split}$$

Notice here that there exist a constant  $\mathit{C} = \Theta(1)$ , whenever  $|R_{1,2}^{(t)}| \geq C(\varrho + \frac{1}{\sqrt{d}})[R_2^{(t)}]^{1/2}[R_1^{(t)}]^{1/2}$ , it will holds

$$\begin{split} R_{1,2}^{(t+1)} &= R_{1,2}^{(t)} \Big( 1 - \Theta(\eta \Sigma_{1,1}^{(t)}[R_1^{(t)}]^2) - \Theta(\eta (\Sigma_{1,1}^{(t)}(E_{1,2}^{(t)})^2 + \Sigma_{2,1}^{(t)}))[R_2^{(t)}]^2 \Big) \\ &= R_{1,2}^{(t)} \Big( 1 - \Theta(\eta \Sigma_{1,1}^{(t)}[R_1^{(t)}]^2) - \Theta(\eta (\Sigma_{1,1}^{(t)}(E_{1,2}^{(t)})^2 + \frac{\alpha_1^6}{d^{3/2}} \Sigma_{2,1}^{(t)}))[R_2^{(t)}]^2 \Big) \end{split}$$

Thus we can go through the same analysis as in the proof of induction for  $E_{1,2}^{(t)}$  to derive that

$$|R_{1,2}^{(t)}| \le \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_2^{(t)}]^{1/2}[R_1^{(t)}]^{1/2}$$

which is the desired result. Note that at the end of phase II

Induction C.1a  $\Longrightarrow$  Lemma C.8a Induction C.1b,c  $\Longrightarrow$  Lemma C.8b

Induction C.1d =⇒ Lemma C.8c

We now complete the proof of Lemma C.8.

## D Phase III: The Acceleration Effect of Prediction Head

We shall prove in this section that the growth of  $E_{2,1}^{(t)}$  in the previous phase creates an acceleration effect to the growth of  $B_{2,2}^{(t)}$ , which will finally outrun the growth of  $B_{2,1}^{(t)}$  to win the lottery. We define

$$T_3 := \min \left\{ t : |B_{2,2}^{(t)}| \ge \frac{1}{2} \min\{|B_{1,1}^{(t)}|, \sqrt{\frac{\eta}{\eta_E}}|E_{2,1}^{(t)}|\} \right\}$$
(D.1)

and we call iterations  $t \in [T_2, T_3]$  as the phase III of training and  $t \ge T_3$  as the end phase of training.

### D.1 Induction in Phase III

**Inductions D.1** (Phase III). During  $t \in [T_2, T_3]$ , we hypothesize the following conditions holds.

$$\begin{array}{l} (a) \ |B_{1,1}^{(t)}| = \Theta(1), \ B_{2,1}^{(t)} = B_{2,1}^{(T_2)}(1 \pm o(1)), \ B_{1,2}^{(t)} = B_{1,2}^{(T_2)}(1 \pm o(1)), \ |B_{2,2}^{(t)}| \in [|B_{2,2}^{(T_2)}|, O(1)], \\ (b) \ |E_{2,1}^{(t)}| = \Theta(\sqrt{\eta_E/\eta}), \ \mathrm{sign}(E_{2,1}^{(t)}) = \mathrm{sign}(E_{2,1}^{(T_2)}) \\ and \ |E_{1,2}^{(t)}| \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}, \\ (c) \ R_1^{(t)} \in [\Omega(\frac{1}{d}), O(\frac{d^{o(1)}}{d^{3/4}})], \ [R_2^{(t)}] \in [\frac{1}{\sqrt{d}}, O(\frac{1}{\log d}\sqrt{\eta_E/\eta})] \end{array}$$

As usual, before we prove the induction, we need to derive some useful claims. But firstly we shall give a much cleaner form of  $\nabla_{E_{j,3-j}}L(W^{(t)},E^{(t)})$  to help us understand the learning process of phase III and the end phase.

### Fact D.2. Let us write

$$\Xi_{j}^{(t)} = C_{0}C_{1}\alpha_{1}^{6}\alpha_{2}^{6}\Phi_{j}^{(t)}\left((B_{1,1}^{(t)})^{6}(B_{2,2}^{(t)})^{6} + (B_{2,1}^{(t)})^{6}(B_{1,2}^{(t)})^{6}\right)$$
$$\Delta_{j,\ell}^{(t)} = C_{0}\Phi_{j}^{(t)}\alpha_{\ell}^{6}(B_{j,\ell}^{(t)})^{3}(B_{3-j,\ell}^{(t)})^{3}C_{2}\mathcal{E}_{j,3-j}^{(t)}$$

Then the gradient of  $E^{(t)}_{j,3-j}$  can be written as

$$-\nabla E_{j,3-j}L(W(t),E(t)) = -\Xi(jt)E_{j,}(t3)-j + X\Delta(j,t) - X\Sigma(j,t)\nabla E_{j,3-j}E_{j,}(t3)-j$$

$$`\in [2] \qquad `\in [2]$$

*Proof.* By expanding the gradients of  $E_{j,3-j}^{(t)}$ , we can verify by checking each monomial of polynomials of  $B_{i}$  to obtain the first term, and leave the  $\mathcal{E}_{j,3-j}^{(t)}$  part for the second term.  $\square$ 

**Lemma D.3** (variables control at phase III). For  $t \in [T_2, T_3]$ , if Induction D.1 holds at iteration t, then we have

$$(a) \ \ \Phi_1^{(t)} = \widetilde{\Theta}(\tfrac{1}{\alpha_1^{12}}), \ [Q_2^{(t)}]^{-2} = \Theta(C_2[R_2^{(t)}]^3 + C_1\alpha_2^6(B_{2,2}^{(t)})^6), \ U_2^{(t)} = \Theta(C_1(\alpha_1^6(E_{2,1}^{(t)})^2 + \alpha_2^6(B_{2,2}^{(t)})^6)).$$

(b) 
$$H_{1,1}^{(t)} = \Theta(C_1 \alpha_1^6), \ H_{1,2}^{(t)} \le O(C_2[R_1^{(t)}]^3) + \widetilde{O}(\frac{\alpha_2^6}{d^3}).$$

(c) 
$$H_{2,1}^{(t)} = \Theta(C_1 \alpha_1^6(E_{2,1}^{(t)})^2), H_{2,2}^{(t)} = \Theta(C_2[R_2^{(t)}]^3)$$

(d) 
$$\Sigma_{1,2}^{(t)} \leq \widetilde{O}(\frac{|E_{1,2}^{(t)}|}{d^{3/2}})\Sigma_{1,1}^{(t)}$$

(e) 
$$\mathcal{E}_{j,3-j}^{(t)} = (1 \pm o(1)) \hat{\mathcal{E}}_{j}^{(t)} = O(C_2[R_j^{(t)}]^3)$$

*Proof.* Assuming Induction D.1 holds at  $t \in [T_2, T_3]$ , we can recall the expression of these variables and prove their bounds directly. The bounds for  $\Phi_1$  and  $H_{1,1}$  comes from  $|B_{1,1}^{(t)}| = \Theta(1)$  and

 $|B_{1,2}^{(t)}|, |E_{1,2}^{(t)}| = o$  (1). The bounds for  $Q_2, U_2$  comes from our definition of  $T_3$  in (D.1). The rest of the claims can be derived by similar arguments using Induction D.1.

### D.2 Gradient Lemmas for Phase III

In this subsection, we would give some gradient lemmas concerning the dynamics of our network in Phase III.

**Lemma D.4** (learning feature  $v_2$  in phase III). For each  $t \in [T_2, T_3]$ , if Induction D.1 holds at iteration t, then we have:

$$(a) \ \langle -\nabla_{w_1} L(W^{(t)}, E^{(t)}), v_2 \rangle = \Theta(\frac{(B_{1,2}^{(t)})^2}{(B_{2,2}^{(t)})^2}) E_{2,1}^{(t)} \Lambda_{2,2}^{(t)} \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4}) |E_{2,1}^{(t)}|^2 \Phi_2^{(t)} \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{5/2}}).$$

(b) 
$$\langle -\nabla_{w_2} L(W^{(t)}, E^{(t)}), v_2 \rangle = (1 \pm \widetilde{O}(\frac{1}{d})) \Lambda_{2,2}^{(t)}$$

*Proof.* Since  $\langle -\nabla_{w_j} L(W^{(t)}, E^{(t)}), v_2 \rangle = \Lambda_{j,2}^{(t)} + \Gamma_{j,2}^{(t)} - \Upsilon_{j,2}^{(t)}$ , let us write down the definition of  $\Lambda_{j,2}^{(t)}, \Gamma_{j,2}^{(t)}, \Upsilon_{j,2}^{(t)}$  respectively:

$$\begin{split} &\Lambda_{j,2}^{(t)} = C_0 \alpha_2^6 \Phi_j^{(t)} H_{j,1}^{(t)} (B_{j,2}^{(t)})^5 \\ &\Gamma_{j,2}^{(t)} = C_0 \alpha_2^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^2 H_{3-j,1}^{(t)} \\ &\Upsilon_{j,2}^{(t)} = C_0 \alpha_1^6 \left( \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{j,2}^{(t)} + \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{3-j,2}^{(t)} \right) \end{split}$$

Again we decompose  $\Upsilon_{j,2}^{(t)}=\Upsilon_{j,2,1}^{(t)}+\Upsilon_{j,2,2}^{(t)}$  as in the proof of Lemma C.6, where

$$\Upsilon_{j,2,1}^{(t)} = C_0 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{j,2}^{(t)} \qquad \Upsilon_{j,2,2}^{(t)} = \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{3-j,2}^{(t)})^2 K_{3-j,2}^{(t)}$$

This gives

$$\begin{split} \Lambda_{j,2}^{(t)} - \Upsilon_{j,2,1}^{(t)} &= C_0 \alpha_2^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^5 H_{j,1}^{(t)} - C_0 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{j,2}^{(t)} \\ &= C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^5 \left( E_{j,3-j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,1}^{(t)})^3 + (E_{j,3-j}^{(t)})^2 (B_{3-j,1}^{(t)})^6 \right) \\ &- C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^2 (B_{3-j,2}^{(t)})^3 E_{j,3-j}^{(t)} \left( (B_{j,1}^{(t)})^6 + E_{j,3-j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,1}^{(t)})^3 \right) \\ &+ C_0 \alpha_2^6 \Phi_j^{(t)} (B_{j,2}^{(t)})^5 C_2 \mathcal{E}_{j,3-j}^{(t)} \end{split}$$

When j = 1, from Induction D.1 and Lemma D.3a (which gives  $\Phi_1^{(t)} \leq \alpha_1^{O(1)} \Phi_2^{(t)}$ ), we can crudely obtain  $\left| C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_1^{(t)} (B_{1,2}^{(t)})^5 \left( E_{1,2}^{(t)} (B_{2,1}^{(t)})^3 (B_{1,1}^{(t)})^3 + (E_{1,2}^{(t)})^2 (B_{2,1}^{(t)})^6 \right) \right| \leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4}) \Phi_1^{(t)} |E_{1,2}^{(t)}|$   $\left| C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_1^{(t)} (B_{1,2}^{(t)})^2 (B_{2,2}^{(t)})^3 E_{1,2}^{(t)} \left( (B_{1,1}^{(t)})^6 + E_{1,2}^{(t)} (B_{2,1}^{(t)})^3 (B_{1,1}^{(t)})^3 \right) \right| \leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d}) \Lambda_{2,2}^{(t)} |E_{1,2}^{(t)}|$   $\left| C_0 \alpha_2^6 \Phi_1^{(t)} (B_{1,2}^{(t)})^5 C_2 \mathcal{E}_{1,2}^{(t)} \right| = \widetilde{O}(\frac{\alpha_1^6}{d^{5/2}}) \Sigma_{1,1}^{(t)} [R_1^{(t)}]^3$ 

So we have

$$\Lambda_{1,2}^{(t)} - \Upsilon_{1,2,1}^{(t)} = \widetilde{O}(\frac{\alpha_1^6}{d^{5/2}}) \Sigma_{1,1}^{(t)} [R_1^{(t)}]^3 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d}) \Lambda_{2,2}^{(t)} |E_{1,2}^{(t)}|$$

When j = 2, we can also derive using Lemma D.3 about  $H_{2,1}^{(t)}$  and Induction D.1 about  $B_{2,1}^{(t)}$  and some rearrangement to obtain

$$C_{0}\alpha_{2}^{6}\Phi_{2}^{(t)}(B_{2,2}^{(t)})^{5}\left[C_{1}\alpha_{1}^{6}\left(E_{2,1}^{(t)}(B_{1,1}^{(t)})^{3}(B_{2,1}^{(t)})^{3}+(E_{2,1}^{(t)})^{2}(B_{1,1}^{(t)})^{6}\right)+C_{2}\mathcal{E}_{2,1}^{(t)}\right]=(1\pm\widetilde{O}(\frac{1}{d}))\Lambda_{2,2}^{(t)}$$

$$\left|C_{0}\alpha_{2}^{6}C_{1}\alpha_{1}^{6}\Phi_{2}^{(t)}(B_{2,2}^{(t)})^{2}(B_{1,2}^{(t)})^{3}E_{2,1}^{(t)}\left((B_{2,1}^{(t)})^{6}+E_{2,1}^{(t)}(B_{1,1}^{(t)})^{3}(B_{2,1}^{(t)})^{3}\right)\right|\leq\widetilde{O}(\frac{\alpha_{1}^{O(1)}}{d^{3}})|E_{2,1}^{(t)}|\Phi_{2}^{(t)}$$

which leads to the approximation

$$\Lambda_{2,2}^{(t)} - \Upsilon_{1,2,2}^{(t)} = (1 \pm \widetilde{O}(\frac{1}{d}))\Lambda_{2,2}^{(t)} \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^3})|E_{2,1}^{(t)}|\Phi_2^{(t)}$$

Similarly, we can also calculate

$$\begin{split} \Gamma_{j,2}^{(t)} - \Upsilon_{j,2,2}^{(t)} &= C_0 \alpha_2^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^2 H_{3-j,1}^{(t)} - C_0 \alpha_1^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,2}^{(t)})^2 K_{3-j,2}^{(t)} \\ &= C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_{3-j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^2 E_{3-j,j}^{(t)} \left( E_{3-j,j}^{(t)} (B_{j,1}^{(t)})^3 (B_{3-j,1}^{(t)})^3 + (E_{3-j,j}^{(t)})^2 (B_{j,1}^{(t)})^6 \right) \\ &- C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_{3-j}^{(t)} (B_{j,2}^{(t)})^5 (E_{3-j,j}^{(t)})^2 \left( (B_{3-j,1}^{(t)})^6 + E_{3-j,j}^{(t)} (B_{j,1}^{(t)})^3 (B_{3-j,1}^{(t)})^3 \right) \\ &+ C_0 \alpha_2^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^2 C_2 \mathcal{E}_{3-j,j}^{(t)} \end{split}$$

When j = 1, following similar procedure as above, we can apply Induction D.1 and Lemma D.3 to give

$$\Gamma_{1,2}^{(t)} - \Upsilon_{1,2,2}^{(t)} = \Theta(\frac{(B_{1,2}^{(t)})^2}{(B_{2,2}^{(t)})^2}) E_{2,1}^{(t)} \Lambda_{2,2}^{(t)} \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4}) |E_{2,1}^{(t)}|^2 \Phi_2^{(t)}$$

Note that the first term on the RHS dominates the term  $\pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d})\Lambda_{2,2}^{(t)}|E_{1,2}^{(t)}|$  in the approximation for  $\Lambda_{1,2}^{(t)}-\Upsilon_{1,2,1}^{(t)}$  due to Induction D.1a,b. When j = 2, since  $\Phi^1$   $\leq \widetilde{\Theta}(\frac{1}{\alpha_1^{12}}) \leq \alpha_1^{O(1)}\Phi_2^{(t)}H_{2,1}^{(t)}$  in this phase and  $|B_{1,1}^{(t)}|=O$ (1), we can derive

$$|\Gamma_{2,2}^{(t)} - \Upsilon_{2,2,2}^{(t)}| \le \widetilde{\Theta}(\frac{\alpha_1^{O(1)}}{d^3})(E_{1,2}^{(t)})^2 \Phi_1^{(t)} + \alpha_1^{O(1)}(E_{1,2}^{(t)})^2 \Lambda_{2,2}^{(t)}$$

It can be seen that  $(E_{1,2}^{(t)})^2\Phi_1^{(t)} \leq (E_{2,1}^{(t)})^2\Phi_2^{(t)}$  by Induction D.1 and Lemma D.3. And by similar arguments we can have  $(1^{\pm \widetilde{O}(\frac{1}{d})})\Lambda_{2,2}^{(t)} \geq \frac{1}{d^{\Omega(1)}}\widetilde{O}(\frac{\alpha_1^{O(1)}}{d^3})|E_{2,1}^{(t)}|\Phi_2^{(t)}$ . Combining all the results above, we can finish the proof.

**Lemma D.5** (learning feature  $v_1$  in Phase III). For each  $t \in [T_2, T_3]$ , if Induction D.1 holds at iteration t, then we have: (recall that  $\Delta$ -notation is from Fact D.2)

$$(a) \langle -\nabla_{w_1} L(W^{(t)}, E^{(t)}), v_1 \rangle = \Theta(\Sigma_{1,1}^{(t)}[R_1^{(t)}]^3) \pm O(\frac{(B_{1,2}^{(t)})^3}{(B_{2,2}^{(t)})^3} + \frac{1}{\sqrt{d}}) \alpha_1^{O(1)} \Lambda_{2,2}^{(t)} + \frac{E_{2,1}^{(t)}}{B_{1,1}^{(t)}} \Delta_{2,1}^{(t)} - \frac{B_{2,2}^{(t)}}{B_{1,1}^{(t)}} \Lambda_{2,2}^{(t)}$$

$$(b) \langle -\nabla_{w_2} L(W^{(t)}, E^{(t)}), v_1 \rangle = \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{5/2}}) \Phi_2^{(t)} [R_2^{(t)}]^3 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d}) \Lambda_{2,2}^{(t)} \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^3})$$

Proof. Recall that  $\langle -\nabla_{w_j} L(W^{(t)}, E^{(t)}), v_1 \rangle = \Lambda_{j,1}^{(t)} + \Gamma_{j,1}^{(t)} - \Upsilon_{j,1}^{(t)}$ . Similar to the proof of Lemma D.4, we can decompose  $\Upsilon_{j,-1}^{(t)} = \Upsilon_{j,1,1}^{(t)} + \Upsilon_{j,1,2}^{(t)}$  and do similar calculations:  $\Lambda_{j,1}^{(t)} - \Upsilon_{j,1,1}^{(t)} = C_0 C_1 \alpha_1^6 \alpha_2^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^5 \left( E_{j,3-j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^3 + (E_{j,3-j}^{(t)})^2 (B_{3-j,2}^{(t)})^6 \right)$ 

$$\begin{split} \Lambda_{j,1}^{(t)} - \Upsilon_{j,1,1}^{(t)} &= C_0 C_1 \alpha_1^6 \alpha_2^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^5 \left( E_{j,3-j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^3 + (E_{j,3-j}^{(t)})^2 (B_{3-j,2}^{(t)})^6 \right) \\ &\quad - C_0 C_1 \alpha_1^6 \alpha_2^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^2 (B_{3-j,1}^{(t)})^3 E_{j,3-j}^{(t)} \left( (B_{j,2}^{(t)})^6 + E_{j,3-j}^{(t)} (B_{3-j,2}^{(t)})^3 (B_{j,2}^{(t)})^3 \right) \\ &\quad + C_0 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^5 C_2 \mathcal{E}_{j,3-j}^{(t)} \end{split}$$

When j=1, from Induction D.1 and Lemma D.3a we know  $\Phi_1^{(t)} \leq \alpha_1^{(O(1))}$  during  $t \in [T_2, T_3]$ , which allow us to derive

$$\begin{split} &C_{0}C_{1}\alpha_{1}^{6}\alpha_{2}^{6}\Phi_{1}^{(t)}(B_{1,1}^{(t)})^{5}\left(E_{1,2}^{(t)}(B_{2,2}^{(t)})^{3}(B_{1,2}^{(t)})^{3}+(E_{1,2}^{(t)})^{2}(B_{2,2}^{(t)})^{6}\right)\\ &\leq \widetilde{O}(\Sigma_{1,1}^{(t)}(E_{1,2}^{(t)})^{2})+C_{0}C_{1}\alpha_{1}^{6}\alpha_{2}^{6}(B_{1,1}^{(t)})^{5}E_{1,2}^{(t)}(B_{2,2}^{(t)})^{3}(B_{1,2}^{(t)})^{3}\\ &\leq O(\frac{(B_{1,2}^{(t)})^{3}}{(B_{2,2}^{(t)})^{3}})\alpha_{1}^{O(1)}\Lambda_{2,2}^{(t)}|E_{1,2}^{(t)}|+\Theta(\Sigma_{1,1}^{(t)}[R_{1}^{(t)}]^{3}) \end{split}$$

And

$$\left| C_0 C_1 \alpha_1^6 \alpha_2^6 \Phi_1^{(t)}(B_{1,1}^{(t)})^2 (B_{2,1}^{(t)})^3 E_{1,2}^{(t)} \left( (B_{1,2}^{(t)})^6 + E_{1,2}^{(t)} (B_{2,2}^{(t)})^3 (B_{1,2}^{(t)})^3 \right) \right| \leq \widetilde{O}(\frac{1}{d^{3/2}}) |E_{1,2}^{(t)}| \Lambda_{2,2}^{(t)} |$$

which can be summarized as

$$\Lambda_{1,1}^{(t)} - \Upsilon_{1,1,1}^{(t)} = \Theta(\Sigma_{1,1}^{(t)}[R_1^{(t)}]^3) \pm O(\frac{(B_{1,2}^{(t)})^3}{(B_{2,2}^{(t)})^3} + (B_{2,1}^{(t)})^3 + \frac{1}{\sqrt{d}})|E_{1,2}^{(t)}|\alpha_1^{O(1)}\Lambda_{2,2}^{(t)}|$$

A similar calculation also gives

$$\Lambda_{2,1}^{(t)} - \Upsilon_{2,1,1}^{(t)} = \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{5/2}})\Phi_2^{(t)}[R_2^{(t)}]^3 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4})\Phi_2^{(t)}|E_{2,1}^{(t)}| \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d})\Lambda_{2,2}^{(t)}B_{2,2}^{(t)}$$

Now we turn to the other terms in the gradient, from similar calculations in the proof of Lemma C.6, we have

$$\begin{split} \Gamma_{j,1}^{(t)} - \Upsilon_{j,1,2}^{(t)} &= C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_{3-j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,1}^{(t)})^2 E_{3-j,j}^{(t)} \left( E_{3-j,j}^{(t)} (B_{j,2}^{(t)})^3 (B_{3-j,2}^{(t)})^3 + (E_{3-j,j}^{(t)})^2 (B_{j,2}^{(t)})^6 \right) \\ &- C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_{3-j}^{(t)} (B_{j,1}^{(t)})^5 (E_{3-j,j}^{(t)})^2 \left( (B_{3-j,2}^{(t)})^6 + E_{3-j,j}^{(t)} (B_{j,2}^{(t)})^3 (B_{3-j,2}^{(t)})^3 \right) \\ &+ C_0 \alpha_2^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,1}^{(t)})^3 (B_{j,1}^{(t)})^2 C_2 \mathcal{E}_{3-j,j}^{(t)} \end{split}$$

which also similarly gives

$$\Gamma_{1,1}^{(t)} - \Upsilon_{1,1,2}^{(t)} = \frac{E_{2,1}^{(t)}}{B_{1,1}^{(t)}} \Delta_{2,1}^{(t)} - \frac{B_{2,2}^{(t)}}{B_{1,1}^{(t)}} \Lambda_{2,2}^{(t)} \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{3/2}}) \Lambda_{2,2}^{(t)}$$

and

$$|\Gamma_{2,1}^{(t)} - \Upsilon_{2,1,2}^{(t)}| \leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d})\Phi_1^{(t)}((E_{1,2}^{(t)})^2 + |E_{1,2}^{(t)}|[R_1^{(t)}]^3) \leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^3})$$

which finishes the proof.

**Lemma D.6** (reducing noise in phase III). *Suppose Induction D.1 holds at*  $t \in [T_2, T_3]$ , then we have

$$\begin{split} \langle -\nabla_{w_1}L(W^{(t)},E^{(t)}),\Pi_{V^{\perp}}w_1^{(t)}\rangle &= -\Theta([R_1^{(t)}]^3) \Big( \Sigma_{1,1}^{(t)} + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)}(E_{2,1}^{(t)})^2 \Big) \\ &\pm O\Big( \sum_{j,\ell} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)}(\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big); \\ \langle b\rangle & \langle -\nabla_{w_1}L(W^{(t)},E^{(t)}),\Pi_{V^{\perp}}w_2^{(t)}\rangle &= \Big( \Sigma_{1,1}^{(t)} + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)}(E_{2,1}^{(t)})^2 \Big) (-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} \\ &+ O\Big( \sum_{(j,\ell) \neq (1,2)} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} R_1^{(t)} [R_2^{(t)}]^2 \Big) \\ \langle c\rangle & \langle -\nabla_{w_2}L(W^{(t)},E^{(t)}),\Pi_{V^{\perp}}w_2^{(t)}\rangle &= -\Theta([R_2^{(t)}]^3) \Big( \sum_{\ell \in [2]} \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big) \\ &\pm O\Big( \sum_{j,\ell} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big); \\ \langle d\rangle & \langle -\nabla_{w_2}L(W^{(t)},E^{(t)}),\Pi_{V^{\perp}}w_1^{(t)}\rangle &= \Big( \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big) (-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_2^{(t)}]^{5/2} [R_1^{(t)}]^{1/2} \\ &+ O\Big( \sum_{(j,\ell) \neq (1,2)} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} R_2^{(t)} [R_1^{(t)}]^2 \Big) \end{split}$$

*Proof.* The proof of Lemma D.6 is very similar to Lemma C.5, but we write it down to stress some minor differences. As in (A.2), we first write down

$$h - \nabla_{w_1} L(W(t), E(t)), \Pi V_{\perp} w_1(t) \mathbf{i} = -X \Sigma(j, t) h \nabla_{w_1} E_{j, t}(t_3) - j, w_1(t) \mathbf{i}$$

**Proof of (a):** Combine the bounds above, we can obtain for each  $j \in [2]$ :  $\Sigma_{1,2}^{(t)} = \widetilde{O}(E_{1,2}^{(t)}/d^{3/2})\Sigma_{1,1}^{(t)}$ . We can then directly apply Claim A.1 to prove Lemma D.6a as follows

**Proof of (b):** For Lemma C.5b, we can use the same analysis for  $\Sigma_{1,1}$  above and Claim A.1d,e to get (again we have used  $\Sigma_{1,2}^{(t)} = \widetilde{O}(E_{1,2}^{(t)}/d^{3/2})\Sigma_{1,1}^{(t)}$ )

$$\begin{split} & \langle -\nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^{\perp}} w_2^{(t)} \rangle \\ &= (1 \pm \widetilde{O}(E_{1,2}^{(t)}/d^{3/2})) \Sigma_{1,1}^{(t)} \Big( (-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} + E_{1,2}^{(t)} R_1^{(t)} [R_2^{(t)}]^2 \Big) \\ &+ \Theta(\Sigma_{2,1}^{(t)} + \Sigma_{2,2}^{(t)}) \Big( (-\Theta(\overline{R}_{1,2}^{(t)}) + O(\varrho)) (E_{2,1}^{(t)})^2 [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} + E_{2,1}^{(t)} R_1^{(t)} [R_2^{(t)}]^2 \Big) \\ &= \Big( \Sigma_{1,1}^{(t)} + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} (E_{2,1}^{(t)})^2 \Big) ((-\Theta(\overline{R}_{1,2}^{(t)}) + O(\varrho)) [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} ) \\ &+ O\Big( \sum_{(j,\ell) \neq (2,1)} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} R_1^{(t)} [R_2^{(t)}]^2 \Big) \end{split}$$

**Proof of (c):** Similarly to the proof of (a), we can also expand as follows

$$\begin{split} & \langle -\nabla_{w_2} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} w_2^{(t)} \rangle \\ &= (1 \pm \widetilde{O}(E_{1,2}^{(t)}/d^{3/2})) \Sigma_{1,1}^{(t)} \Big( - [R_2^{(t)}]^3 \Theta((E_{1,2}^{(t)})^2) \pm O(E_{1,2}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big) \\ &- \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big( [R_2^{(t)}]^3 \pm O(E_{2,1}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big) \\ &= -\Theta([R_2^{(t)}]^3) \Big( \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big) \pm O\Big( \sum_{j,\ell} \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big) \end{split}$$

Proof of (d): Similarly, we can calculate

$$\begin{split} &\langle -\nabla_{w_2}L(W^{(t)},E^{(t)}),\Pi_{V^\perp}w_1^{(t)}\rangle\\ &= (1\pm\widetilde{O}(E_{1,2}^{(t)}/d^{3/2}))\Sigma_{1,1}^{(t)}\Big((-\Theta(\overline{R}_{1,2}^{(t)})\pm O(\varrho))(E_{1,2}^{(t)})^2[R_2^{(t)}]^{5/2}[R_1^{(t)}]^{1/2} + E_{1,2}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big)\\ &+ \sum_{\ell\in[2]}\Sigma_{2,\ell}^{(t)}\Big((-\Theta(\overline{R}_{1,2}^{(t)})\pm O(\varrho))[R_2^{(t)}]^{5/2}[R_1^{(t)}]^{1/2} + E_{1,2}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big)\\ &= \Big(\Sigma_{1,1}^{(t)}\Theta((E_{1,2}^{(t)})^2) + \sum_{\ell\in[2]}\Sigma_{2,\ell}^{(t)}\Big)(-\Theta(\overline{R}_{1,2}^{(t)})\pm O(\varrho))[R_2^{(t)}]^{5/2}[R_1^{(t)}]^{1/2}\\ &+ O\Big(\sum_{(j,\ell)\neq(2,1)}\Sigma_{j,\ell}^{(t)}E_{j,3-j}^{(t)}R_2^{(t)}[R_1^{(t)}]^2\Big) \end{split}$$

which completes the proof.

**Lemma D.7** (learning the prediction head in phase III). *If Induction D.1 holds at iteration*  $t \in [T_2, T_3]$ , *then using the notations from Fact D.2, we have* 

$$\begin{split} -\nabla_{E_{j,3-j}}L(W^{(t)},E^{(t)}) &= \Theta(\sum_{\ell \in [2]} \Sigma_{j,\ell}^{(t)}) (-E_{j,3-j}^{(t)}[R_{3-j}^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}) \\ &- \Xi_j^{(t)} E_{j,3-j}^{(t)} + \sum_{\ell \in [2]} \Delta_{j,\ell}^{(t)} \end{split}$$

*Proof.* By Fact D.2, we only need to bound the last term  $\sum_{\ell \in [2]} \sum_{j,\ell}^{(t)} \nabla_{E_{1,2}} \mathcal{E}_{j,3-j}^{(t)}$ , which can be directly obtained from applying Claim A.1.

### D.3 At the End of Phase III

In order to argue that  $B_{2,2}^{(T_2)} = \Omega(1)$  at the end of phase III, we need to define some auxiliary notions. Recall that  $T_3$  is defined in (D.1), and now we further define

$$T_{3,1} := \min\{t : C_1\alpha_2^6(B_{2,2}^{(t)})^6 \ge C_2[R_2^{(t)}]^3\}, \qquad T_{3,2}^{(t)} = \min\{t : |B_{2,2}^{(t)}| \ge \frac{1}{3}\min\{|E_{2,1}^{(t)}|, |B_{1,1}^{(t)}|\}\} \quad (D.2)$$

It can be observed that if Induction D.1 holds for  $t \in [T_2, T_3]$  and our learning rate  $\eta$  is small enough, we shall have  $T_2 < T_{3,1} \le T_{3,2} < T_3$ . Now we are ready to present the main lemma we want to prove in this phase.

**Lemma D.8** (Phase III). Let  $T_3$  be defined as in (D.1). Suppose  $\eta = \frac{1}{\mathsf{poly}(d)}$  is sufficiently small, then Induction D.1 holds for all iteration  $t \in [T_2, T_3]$ , and at iteration  $t = T_3$ , the followings holds:

$$(a) \ |B_{1,1}^{(T_3)}| = \Theta(1), \ |B_{2,2}^{(T_3)}| = \Theta(1), \ B_{j,\ell}^{(T_3)} = B_{j,\ell}^{(T_2)}(1 \pm o(1)) \text{ for $j$ 6= $$};$$

(b) 
$$R_1^{(T_3)} = \widetilde{O}(\frac{1}{d^{3/4}}), \ R_2^{(T_3)} \in [\widetilde{O}(\frac{1}{d^{1/2}}), \widetilde{O}(\frac{1}{d^{1/4}})], \ and \ \overline{R}_{1,2}^{(T_3)} \le \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

$$(c) |E_{2,1}^{(T_2)}| = \Theta(\sqrt{\eta_E/\eta})_{\textit{and}} |E_{1,2}^{(T_2)}| = \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} = \widetilde{O}(\frac{1}{d})$$

Moreover,  $|B_{2,2}^{(t)}|$  is increasing and  $R_2^{(t)}$  is decreasing. The part of learning  $|B_{2,2}^{(t)}|$  till  $\Omega(1)$  and keeping  $B_{2,1}^{(t)}$  close to its initialization is what's been accelerated by the prediction head  $E_{2,1}^{(t)}$ .

The proof of Lemma D.8 will be proven after we have proven Induction D.1, which will again be proven after some intermediate results are proven.

**Lemma D.9** (The growth of  $B_{2,2}^{(t)}$  before  $T_{3,1}$ ). Let  $T_{3,1}$  be defined as in (D.2). If Induction D.1 holds for  $t \in [T_2, T_{3,1}]$ , then we have  $R_2^{(T_{3,1})} \leq \frac{\alpha_1^{12}}{d^{1/4}} and B_{2,2}^{(T_{3,1})} \in [\frac{1}{d^{1/4}}, O(\frac{\alpha_1^{O(1)}}{d^{1/4}})]$  and  $T_{3,1} \leq T_2 + \widetilde{O}(\frac{d^{1.625}\alpha_1}{\eta})^{(1)}$ .

*Proof.* Firstly by Lemma D.6b , we can write down the update of  $R_2^{(t)}$ : (as in Lemma C.8)

$$\pm \, O\Big( \sum_{j,\ell} \eta \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)}(R_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \Big) \pm \frac{\eta}{\operatorname{poly}(d)} \underline{\hspace{2cm}}$$

$$R_2^{(t+1)} = R_2^{(t)} - \eta \Theta([R_2^{(t)}]^3) \Big( \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big)$$

Next, by Claim A.1 and Lemma D.3a combined with Induction D.1a,b, we have  $\widetilde{O}(\frac{|E_{2,1}^{(t)}|}{d^{3/2}})\Sigma_{1,1}^{(t)}\frac{\Phi_1^{(t)}}{\Phi_2^{(t)}}\leq \widetilde{O}(\Sigma_{2,1}^{(t)})$ , which leads to the bound

$$\eta \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) \leq \widetilde{O}(\varrho^2 + \frac{1}{d}) \alpha_1^{O(1)} \eta \Sigma_{1,1}^{(t)} [R_1^{(t)}]^3 [R_2^{(t)}]^3 \leq O(\frac{1}{d^{9/4}}) \eta \Sigma_{1,1}^{(t)} [R_2^{(t)}]^3 \leq O(\frac{\alpha_1^{O(1)}}{d^{3/4}}) \eta \Sigma_{2,1}^{(t)} [R_2^{(t)}]^3$$

Similarly, we can bound the following term

$$\begin{split} \sum_{\ell \in [2]} \eta \Sigma_{1,\ell}^{(t)} |E_{1,2}^{(t)}| (|\overline{R}_{1,2}^{(t)}| + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} &\leq \widetilde{O}(\varrho^2 + \frac{1}{d}) \alpha_1^{O(1)} \sum_{\ell \in [2]} \eta \Sigma_{1,\ell}^{(t)} [R_1^{(t)}]^3 [R_2^{(t)}]^3 \\ &\leq \widetilde{O}(\varrho^2 + \frac{1}{d}) \alpha_1^{O(1)} \frac{1}{d^{9/4}} \sum_{\ell \in [2]} \eta \Sigma_{1,\ell}^{(t)} [R_2^{(t)}]^3 \\ &\leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{3/4}}) \eta \Sigma_{2,1}^{(t)} [R_2^{(t)}]^3 \end{split}$$

Moreover, from Induction D.1c that  $R_2^{(t)} \geq R_1^{(t)}$ , we can also calculate for each  $t \in [T_2, T_{3,1}]$ :

$$\eta \Sigma_{2,\ell}^{(s)} |E_{2,1}^{(t)}| (|\overline{R}_{1,2}^{(t)}| + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \le \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \alpha_1^{O(1)} \eta \Sigma_{2,\ell}^{(t)} [R_2^{(t)}]^3$$

Thus by combining the results above, we have the update of 
$$R_2^{(t)}$$
 at  $t \in [T_2, T_3]$  as follows: 
$$R_2^{(t+1)} = R_2^{(t)} - \eta \Theta([R_2^{(t)}]^3) \Big( \Sigma_{1,1}^{(t)} \Theta((E_{1,2}^{(t)})^2) + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big)$$
 
$$= R_2^{(t)} - \eta (\Sigma_{2,1}^{(t)} + \Sigma_{2,2}^{(t)}) [R_2^{(t)}]^3 \tag{D.3}$$

which implies that  $R_2^{(t)}$  is decreasing throughout phase III. From Lemma D.3a and Induction D.1b, we know that for  $t \in [T_2, T_{3.1}]$ :

$$\Phi_2^{(t)} = Q_2^{(t)} / [U_2^{(t)}]^{3/2} = \Theta(\frac{1}{\sqrt{C_2[R_2^{(t)}]^3} (C_1 \alpha_1^6 (E_{2,1}^{(t)})^2)^{3/2}})$$

which implies (also using a bit of Claim A.1 and Induction D.1a)

$$\begin{split} \Sigma_{2,1}^{(t)}[R_2^{(t)}]^3 &= (1\pm \widetilde{O}(\frac{1}{d^{3/2}}))E_{2,1}^{(t)}\Delta_{2,1}^{(t)} \\ &= (1\pm \widetilde{O}(\frac{1}{d^{3/2}}))(1\pm \widetilde{O}(\frac{1}{d^{3/2}}))C_0C_2\alpha_1^6\Phi_2^{(t)}E_{2,1}^{(t)}(B_{1,1}^{(t)})^3(B_{2,1}^{(t)})^3[R_2^{(t)}]^3 \\ &= \Theta(\frac{C_2^{1/2}[R_2^{(t)}]^{3/2}}{(U_2^{(t)})^{3/2}})C_0\alpha_1^6E_{2,1}^{(t)}(B_{1,1}^{(t)})^3(B_{2,1}^{(t)})^3 \\ &= \Theta(\frac{C_0C_2^{1/2}|B_{2,1}^{(T_2)}|^3}{C_1^{3/2}\alpha_1^3|E_{2,1}^{(T_2)}|})[R_2^{(t)}]^{3/2} \\ &= \Theta(\frac{C_0C_2^{1/2}|B_{2,1}^{(T_2)}|^3}{C_1^{3/2}\alpha_1^3|E_{2,1}^{(T_2)}|})[R_2^{(t)}]^{3/2} \\ \text{(because } B_{2,1}^{(t)} = B_{2,1}^{(T_2)}(1\pm o(1)), \ B_{1,1}^{(t)} = \Theta(B_{1,1}^{(T_2)}) \text{ and } E_{2,1}^{(t)} = \Theta(E_{2,1}^{(T_2)}) \text{sign}(B_{1,1}^{(T_2)}B_{2,1}^{(T_2)})) \end{split}$$

And for  $\Sigma_{2,2}$ , from some simple calcualtions (using Claim A.1), we have

- when  $|B_{2,2}^{(t)}| \leq \frac{\alpha_1}{\alpha_2} \sqrt{|B_{2,1}^{(T_2)}|}$ , we would have  $\Sigma_{2,2}^{(t)} \leq O(\Sigma_{2,1}^{(t)})$ ;
- otherwise, we have  $\Sigma_{2,1}^{(t)} + \Sigma_{2,2}^{(t)} = \Theta(\Sigma_{2,2}^{(t)}).$

So by (D.3), we know  $R_2$  is decreasing for  $t \in [T_2, T_{3,1}]$  by at least

$$\begin{split} R_2^{(t+1)} & \leq R_2^{(t)} - \eta \Theta(\frac{C_0 C_2^{1/2} |B_{2,1}^{(T_2)}|^3}{C_1^{3/2} \alpha_1^3 |E_{2,1}^{(T_2)}|}) [R_2^{(t)}]^{3/2} \leq R_2^{(t)} (1 - \eta \zeta [R_2^{(t)}]^{1/2} \\ & \qquad \qquad \qquad \qquad \\ & \qquad \qquad \zeta := \Theta(\frac{C_0 C_2^{1/2} |B_{2,1}^{(T_2)}|^3}{C_1^{3/2} \alpha_1^3 |E_{2,1}^{(T_2)}|}) = \widetilde{\Theta}(\frac{\sqrt{\eta/\eta_E}}{d^{3/2} \alpha_1^3}). \text{ By this update, we can} \end{split}$$

where

 $\text{prove}^{T_{3,1}} \le T_2 + O(\frac{d^{3/2+1/8}\alpha_1^{O(1)}}{\eta}).$ 

In order to do that, we can first see that for some  $t_{3,1}' \in [T_2 + \widetilde{\Theta}(\frac{d^{3/2}\alpha_1^2\sqrt{\eta_E/\eta}}{\eta}), T_2 + \widetilde{\Theta}(\frac{d^{3/2}\alpha_1^4\sqrt{\eta_E/\eta}}{\eta})],$  we shall have  $R_2^{(t_{3,1}')} \leq d^{-1/4}$ . Indeed, suppose otherwise  $R_2^{(t_{3,1}'-1)} \geq d^{-1/4}$ , then (D.4) implies

$$\begin{split} R_2^{(t_{3,1}')} &\leq R_2^{(t_{3,1}'-1)} (1 - \eta \zeta [R_2^{(t_{3,1}'-1)}]^{1/2}) \leq R_2^{(t_{3,1}'-1)} (1 - \eta \zeta \frac{1}{d^{1/8}}) \\ &\leq R_2^{(T_2)} \left( 1 - \Theta(\frac{C_0 C_2^{1/2} \sqrt{\eta/\eta_E}}{C_1^{3/2} d^{3/2} \alpha_1^3}) \frac{\eta}{d^{1/8}} \right)^{t_{3,1}'-T_2-1} \\ &\leq O(\sqrt{\eta_E/\eta}) \left( 1 - \Theta(\frac{C_0 C_2^{1/2} \sqrt{\eta/\eta_E}}{C_1^{3/2} d^{3/2} \alpha_1^3}) \frac{\eta}{d^{1/8}} \right)^{t_{3,1}'-T_2-1} \end{split}$$

which means there must exist an iteration  $t_{3,1}' \in [T_2 + \widetilde{\Theta}(\frac{d^{3/2}\alpha_1^2\sqrt{\eta_E/\eta}}{\eta}), T_2 + \widetilde{\Theta}(\frac{d^{3/2}\alpha_1^4\sqrt{\eta_E/\eta}}{\eta})] \text{ such that } R_2^{(t_{3,1}'-1)} \geq d^{-1/4} \text{ (so the above update bound is still valid when the RHS is for } t \leq t_{3,1}' - 1)$ 

and  $R_2^{(t'_{3,1})} < d^{-1/4}$ . Next we need to prove that at  $t = t'_{3,1}$ , it holds  $C_1 \alpha_2^6 (B_{2,2}^{(t)})^6 \ge C_2 [R_2^{(t)}]^3$ . Let us discuss several possible cases:

- 1. Suppose  $|B_{2,2}^{(t'_{3,1})}| \geq \frac{\alpha_1}{\alpha_2} |B_{2,1}^{(T_1)}|^{1/2} \geq \Theta(\frac{1}{d^{1/4}})$  (by Induction D.1a and Lemma D.8), then we already have  $C_1 \alpha_2^6 (B_{2,2}^{(t'_{3,1})})^6 \geq C_2 [R_2^{(t'_{3,1})}]^3$  and  $T_{3,1} \leq t'_{3,1}$ ;
- $\begin{aligned} \text{2. Suppose otherwise} &|B_{2,2}^{(t'_{3,1})}| \leq \frac{\alpha_1}{\alpha_2} |B_{2,1}^{(T_1)}|^{1/2} \text{, then we shall have } \Sigma^{(t)}_{2,2} \leq O(\Sigma^{(t)}_{2,1}) \text{. So the update of } R^{(t)}_2 \\ &\text{during } t \in [T_2,T_{3,1}] \text{ can be written as } \\ &R^{(t+1)}_2 = R^{(t)}_2 \Theta(\eta \Sigma^{(t)}_{2,1})[R^{(t)}_2]^3 = R^{(t)}_2(1 \Theta(\eta \zeta)[R^{(t)}_2]^{1/2}) \end{aligned}$

Let  $t_{3,2}'=\min\{t:R_2^{(t)}\leq 2d^{-1/4}\}$  be an iteration between  $T_2$  and  $t_{3,1}'$ , we shall have

$$\sum_{t \in [t_{3,2}',t_{3,1}']} \eta \zeta[R_2^{(t)}]^{3/2} = \Theta(R_2^{(t_{3,2}')} - R_2^{(t_{3,1}')}) = \Theta(\frac{1}{d^{1/4}}) \qquad \text{and} \quad R_2^{(t)} \in [0.99 \frac{1}{d^{1/4}}, 2.01 \frac{1}{d^{1/4}}]$$

which also implies  $t_{3,1}'-t_{3,2}'=\Theta(\frac{d^{1/8}}{\eta\zeta})=\widetilde{\Theta}(\frac{d^{3/2+1/8}\alpha_1^3\sqrt{\eta_E/\eta}}{\eta})$ . In this case, let us look at the update of  $B_{2,2}^{(t)}$  at  $t\in[T_2,T_3]$ . By Lemma D.42, we have

$$B_{2,2}^{(t+1)} = B_{2,2}^{(t)} + \eta(1 \pm \widetilde{O}(\frac{1}{d}))\Lambda_{2,2}^{(t)}$$

It is not hard to see  $|B_{2,2}^{(t)}|$  is monotonically increasing. Also by Induction D.1a and Lemma D.3a, if we sum together the update between  $t_{3,2}'$  and  $t_{3,1}'$  as follows: (suppose the sign of  $B_{2,2}^{(t_{3,2}')}$  is positive for now, the negative case can be similarly dealt with)

$$\begin{split} B_{2,2}^{(t'_{3,2})} &\text{ is positive for now, the negative case can be similarly dealt with)} \\ B_{2,2}^{(t'_{3,2})} + \sum_{t \in [t'_{3,2},t'_{3,1}]} \eta(1 \pm \widetilde{O}(\frac{1}{d})) \Lambda_{2,2}^{(t)} = \sum_{t \in [t'_{3,2},t'_{3,1}]} \Theta(\frac{\eta C_0 C_1 \alpha_1^6 \alpha_2^6 (E_{2,1}^{(T_2)})^2}{\sqrt{C_2 [R_2^{(t)}]^3} (C_1 \alpha_1^6 (E_{2,1}^{(T_2)})^2)^{3/2}}) (B_{2,2}^{(t)})^5 \\ &\geq B_{2,2}^{(t'_{3,2})} + (B_{2,2}^{(T_2)})^4 \sum_{t \in [t'_{3,2},t'_{3,1}]} \Theta(\frac{\eta C_0 \alpha_1^3 \alpha_2^6 B_{2,2}^{(t)}}{C_1^{1/2} C_2^{1/2} [R_2^{(t)}]^{3/2} |E_{2,1}^{(T_2)}|}) \\ &\geq B_{2,2}^{(t'_{3,2})} \prod_{t = t'_{3,2}} \left(1 + \eta \widetilde{\Theta}(\frac{\alpha_1^3 \alpha_2^6}{d^{3/2 + 1/8} \sqrt{\eta_E/\eta}})\right) \\ &\geq \widetilde{\Theta}(\frac{1}{\sqrt{d}}) \left(1 + \eta \widetilde{\Theta}(\frac{\alpha_1^3 \alpha_2^6}{d^{3/2 + 1/8} \sqrt{\eta_E/\eta}})\right)^{\widetilde{\Theta}(d^{3/2} \alpha_1^3 \sqrt{\eta_E/\eta/\eta})} \\ &\geq \Omega(e^{\alpha_1}) \end{split}$$

which is a contradiction to our assumption  $|B_{2,2}^{(t'_{3,1})}| \leq \frac{\alpha_1}{\alpha_2}|B_{2,1}^{(T_1)}|^{1/2}$ . Since  $|B_{2,2}^{(t)}|$  is monotonically increasing, we know there must exist some iteration  $t \leq t'_{3,1}$  such that  $|B_{2,2}^{(t)}| \geq \frac{\alpha_1}{\alpha_2}|B_{2,1}^{(T_1)}|^{1/2}$ . Which means  $T_{3,1} \leq t'_{3,1}$ .

Thus we proved the bound of  $T_{3,1} \leq T_2 + \widetilde{\Theta}(\frac{d^{3/2}\alpha_1^{O(1)}}{\eta})$ .

Using similar arguments, we can prove that  $R_2^{(T_{3,1})} \leq \frac{\alpha_1^{O(1)}}{d^{1/4}}$ . Indeed, we can set  $T_{3,3} \coloneqq \min\{t: |B_{2,2}^{(t_{3,1})}| \geq \frac{\alpha_1}{\alpha_2}|B_{2,1}^{(T_1)}|^{1/2}\}$ . From our arguments in this proof, we know  $\Sigma_{2,2}^{(t_{3,2})} \leq O(\Sigma_{2,2}^{(t)})$  for  $t \leq T_{3,3}$ . Now we can further choose  $t_{3,3}' = \min\{t: R_2^{(t)} \leq a\}$  for some  $a = \frac{\alpha_1^{12}}{d^{1/4}}$  to be some iteration with  $R_2^{(t)} \geq a_{\text{for}} t \in [T_2, t_{3,3}']$  and  $t_{3,3}' - T_2 = \Theta(\frac{\sqrt{a \log d}}{\eta \zeta})$ . Now we can work out the update of  $B_{2,2}^{(t)}$  during  $t \in [T_2, t_{3,3}']$  again to see that  $B_{2,2}^{(t_{3,3})} \leq B_{2,2}^{(T_2)} \left(1 + \eta \widetilde{\Theta}(\frac{\alpha_1^3 \alpha_2^6}{d^2 a^{3/2} \sqrt{\eta_{E/\eta}}})\right)^{\frac{\sqrt{a}}{\eta \zeta}} \leq \widetilde{O}(\frac{1}{\sqrt{d}})$ . This would prove that  $t_{3,3}' \leq T_{3,3}$  and  $t_{3,3}' \geq T_{3,4}$  by our arguments above and the fact that  $t_{3,3}' \geq t_{3,4}' \leq t_$ 

Now we proceed to characterize the learning of  $B_{2,2}^{(t)}$  during  $t \in [T_{3,1}, T_{3,2}]$ .

**Lemma D.10** (The growth of  $B_{2,2}^{(t)}$  until  $T_3$ ). Let  $T_{3,1}$ ,  $T_{3,2}$  be defined as in (D.2). If Induction D.1 holds true for all  $t \in [T_2, T_3]$ , then we have  $T_{3,2} = T_{3,1} + \widetilde{O}(\frac{d^{1/4}\alpha_1^{O(1)}}{\eta})$  and  $T_3 \leq T_{3,2} + \widetilde{O}(\frac{\alpha_1^{O(1)}}{\eta})$ .

*Proof.* We first calculate the bound for  $T_{3,2}$ . After  $T_{3,1}$ , since  $|B_{2,2}^{(t)}|$  is increasing while  $R_2^{(t)}$  is decreasing by Induction D.1. So by Lemma D.3a, we have

$$[Q_2^{(t)}]^{-2} = \Theta(C_1\alpha_2^6(B_{2,2}^{(t)})^6), \quad \Phi_2^{(t)} = Q_2^{(t)}/[U_2^{(t)}]^{3/2} = \Theta((C_1^{3/2}\alpha_2^3\alpha_1^9|B_{2,2}^{(t)})|^3|E_{2,1}^{(t)}|^3)^{-1})$$

So according to Lemma D.4, we would have for all  $t \in [T_{3,1}, T_{3,2})$ :

$$\langle -\nabla_{w_2} L(W^{(t)}, E^{(t)}), v_2 \rangle = (1 \pm o(1)) \Lambda_{2,2}^{(t)} = \Theta(\frac{1}{C_1^{3/2} \alpha_1^9 |E_{2,1}^{(T_2)}|^3}) (B_{2,2}^{(t)})^2 \operatorname{sign}(B_{2,2}^{(t)})$$

where we have used  $(E_{2,1}^{(t)})^3 = \Theta((E_{2,1}^{(T_2)})^3)$  from Induction D.1a. So when  $t \in [T_{3,1}, T_{3,2}]$ , we can (t) write down the explicit form of

 $\Lambda_{2,2}$  and use Lemma D.3d to derive

$$|B_{2,2}^{(t+1)}| = |B_{2,2}^{(t)}| + \eta \Theta\left(\frac{C_1 \alpha_1^6 |E_{2,1}^{(T_2)}|^2}{C_1^{3/2} \alpha_1^9 |E_{2,1}^{(T_2)}|^3}\right) (B_{2,2}^{(t)})^2$$

$$\geq |B_{2,2}^{(t)}| \left(1 + \Theta\left(\frac{1}{C_1 \alpha_1^{O(1)}}\right) |B_{2,2}^{(T_{3,1})}|\right)$$

$$\geq |B_{2,2}^{(t)}| \left(1 + \Theta\left(\frac{1}{C_1 \alpha_1^{O(1)}}\right) \frac{1}{d^{1/4}}\right)$$

Thus after  $\widetilde{O}(\frac{d^{1/4}\alpha^{O(1)}}{\eta})$  many iterations, we would have  $|B_{2,2}^{(t)}| \geq \frac{1}{3}\min\{|E_{2,1}^{(t)}|,|B_{1,1}^{(t)}|\}$ . Now let us deal with the growth of  $|B_{2,2}^{(t)}|$  at  $t \in [T_{3,2},T_{3,3}]$ . During this stage, since  $B_{2,2}^{(t)}$  is still increasing and  $|E_{2,1}^{(t)}|=|E_{2,1}^{(T_2)}|$  by Induction D.1, we have from Lemma D.3a that

$$\Phi_2^{(t)} = Q_2^{(t)}/[U_2^{(t)}]^{3/2} = \Theta(\frac{1}{C_1^2\alpha_2^{12}(B_{2,2}^{(t)})^{12}}) \geq \Theta(\frac{1}{C_1^2\alpha_1^{O(1)}})$$

And we can redo the calcualtions as above to  $\operatorname{get}^{T_3} \leq T_{3,2} + \widetilde{O}(\frac{\alpha_1^{O(1)}}{\eta}) \operatorname{since}^{\sqrt{\eta/\eta_E}}|E_{2,1}^{(t)}| \operatorname{and} |B_1^{(t_1)}|$ are both  $\Theta(1)$  according to Induction D.1a,b

Proving The Main Lemma. Now we finally begin to prove Lemma D.8.

*Proof of Lemma D.8.* We start with proving Induction D.1.

**Proof of Induction D.1a:** From Lemma D.5, we know the update of  $B_{1,1}^{(t)}$  can be written as

$$B_{1,1}^{(t+1)} = B_{1,1}^{(t)} + \Theta(\eta \Sigma_{1,1}^{(t)}[R_1^{(t)}]^3) \pm \eta O(\frac{(B_{1,2}^{(t)})^3}{(B_{2,2}^{(t)})^3} + \frac{1}{\sqrt{d}})\alpha_1^{O(1)}\Lambda_{2,2}^{(t)} + \frac{E_{2,1}^{(t)}}{B_{1,1}^{(t)}}\eta \Delta_{2,1}^{(t)} - \frac{B_{2,2}^{(t)}}{B_{1,1}^{(t)}}\eta \Lambda_{2,2}^{(t)}$$

Since from Lemma D.9 and Lemma D.10, we know  $T_3 \leq \widetilde{O}(\frac{d^{1.625}\alpha_1^{O(1)}}{\eta})$  and from Claim A.1 and

Since from Lemma D.9 and Lemma D.10, we know 
$$S = O(\eta)$$
 and from Cla Induction D.1a,c we have  $\Sigma_{1,1}^{(t)}[R_1^{(t)}]^3 \leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{2.25}})$ , we shall have 
$$\sum_{s \in [T_2,t)} \Theta(\eta \Sigma_{1,1}^{(s)}[R_1^{(s)}]^3) \leq \widetilde{O}(\frac{d^{1.625}\alpha_1^{O(1)}}{\eta}) \widetilde{O}(\frac{\eta \alpha_1^{O(1)}}{d^{2.25}}) \leq \frac{1}{\sqrt{d}} = o(1)$$

Further more, by applying Lemma G.3 to  $x_t = B_{2,2}^{(t)}$  with  $q^0 = q - 2$ , and notice that sign(  $B_{j,2}^{(t)}$ ) =  $sign(B_{j,2}^{(T_2)})$  for all  $t \in [T_2, T_3]$ , we also have

$$\left| \sum_{s \in [T_2,t)} O(\frac{(B_{1,2}^{(s)})^3}{(B_{2,2}^{(s)})^3}) \alpha_1^{O(1)} \eta \Lambda_{2,2}^{(s)} \right| \leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{\sqrt{d}})$$

Now we turn to the last two terms. We first see that from the expression (D.3) of  $R_2^{(t)}$ ,'s update, we have that (note that  $\operatorname{sign}(E_{2,1}^{(t)}\Delta_{2,1}^{(t)})=1)$ 

$$\sum_{s \in [T_2,t)} \frac{E_{2,1}^{(s)}}{|B_{1,1}^{(s)}|} \eta \Delta_{2,1}^{(s)} = \sum_{s \in [T_2,t)} \frac{1}{|B_{1,1}^{(s)}|} \Theta(\eta \Sigma_{2,1}^{(s)}[R_2^{(s)}]^3) = \Theta(\frac{\sqrt{\eta_E/\eta}}{|B_{1,1}^{(T_2)}|}) = \Theta(\sqrt{\eta_E/\eta})$$

where we have used the fact that  $\Sigma_{2,1}^{(t)}[R_2^{(t)}]^3 = (1 \pm O(\frac{1}{d}))E_{2,1}^{(t)}\Delta_{2,1}^{(t)} \\ ext{and} \\ \sum_{s \in [T_2,t)} \eta \Sigma_{2,1}^{(s)}[R_2^{(s)}]^3 \lesssim R_2^{(T_2)}$ from (D.3) (which holds for all  $t \in [T_2, T_3]$ ). And also, the analysis above shows that

$$|B_{1,1}^{(t)}| = |B_{1,1}^{(T_2)}| + O(\sqrt{\eta_E/\eta}) - \sum_{s \in [T_2, t]} \frac{B_{2,2}^{(s)}}{B_{1,1}^{(s)}} \eta \Lambda_{2,2}^{(s)}$$

for all  $t \in [T_2, T_3]$ , which means that either  $\sum_{s \in [T_2, t]} \frac{B_{2,2}^{(s)}}{|B_{1,1}^{(s)}|} \eta \Lambda_{2,2}^{(s)} \leq \sum_{s \in [T_2, t)} \frac{E_{2,1}^{(s)}}{|B_{1,1}^{(s)}|} \eta \Delta_{2,1}^{(s)}$  and we have

$$|B_{1,1}^{(t)}| \geq |B_{1,1}^{(T_2)}|$$
 holds throughout  $t \in [T_2, T_3]$ , or that  $\sum_{s \in [T_2, t]} \frac{B_{2,2}^{(s)}}{|B_{1,1}^{(s)}|} \eta \Lambda_{2,2}^{(s)} \geq \Omega(\sqrt{\eta_E/\eta})$ , in

which case we would have  $|B_{1,1}^{(t)}|$  to be actually decreasing (as  $B_{2,2}^{(t)}$  is increasing). Now that since  $B_{1,1}^{(T_2)}$ =  $\Theta(1)$ , we can easily see by our definition of  $T_3$  and the monotonicity of  $B_{1,1}^{(t)}$  after going below  $B_{1,1}^{(T_2)} - \Omega(\sqrt{\eta_E/\eta})$  that  $B_{1,1}^{(t)} \geq 0.49B_{1,1}^{(T_2)} = \Omega(1)$  for all  $t \in [T_2, T_3]$ .

Next let us look at the change of  $B_{2,1}^{(t)}$ . From Lemma D.5, we can write down the update of  $B_{2,1}^{(t)}$ 

$$B_{2,1}^{(t+1)} = B_{2,1}^{(t)} + \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{5/2}}) \eta \Phi_2^{(t)} [R_2^{(t)}]^3 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{d}) \eta \Lambda_{2,2}^{(t)} \pm \widetilde{O}(\frac{\eta \alpha_1^{O(1)}}{d^3})$$

For the first term, according to Lemma D.9 and Lemma D.10 and  $R_2^{(t)} \leq O(\sqrt{\eta_E/\eta}) = o_{(1)}$  for all  $t \in [T_2, T_3]$  by Induction D.1c, we have  $\Phi_2^{(t)}[R_2^{(t)}]^3 \leq \alpha_1^{O(1)}$  for all  $t \in [T_2, T_3]$  $[T_2,T_3]$  and

$$\sum_{s \in [T_2, t]} \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{5/2}}) \eta \Phi_2^{(s)} [R_2^{(s)}]^3 \le \widetilde{O}(\frac{d^{1.625} \alpha_1^{O(1)}}{\eta}) \eta \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{5/2}}) \le \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{7/8}})$$

And similarly as in the proof of induction for 
$$B_{1,1}^{(t)}$$
, we have 
$$\sum_{s \in [T_2,t]} \widetilde{O}(\frac{\alpha_1^{O(1)}}{d}) \eta \Lambda_{2,2}^{(s)} \leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d}), \quad \sum_{s \in [T_2,t]} \widetilde{O}(\frac{\eta \alpha_1^{O(1)}}{d^3}) \leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d})$$

which proved the induction for  $B_{2,1}^{(t)}$  since  $|B_{2,1}^{(T_2)}| = \widetilde{\Theta}(\frac{1}{\sqrt{d}})$ 

Next we go on for the induction of  $B_{1,2}^{(t)}$ , we write down its update:

$$B_{1,2}^{(t+1)} = B_{1,2}^{(t)} + \Theta(\frac{(B_{1,2}^{(t)})^2}{(B_{2,2}^{(t)})^2}) E_{2,1}^{(t)} \eta \Lambda_{2,2}^{(t)} \pm \eta \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4}) |E_{2,1}^{(t)}|^2 \Phi_2^{(t)} \pm \eta \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^{5/2}})$$

By Lemma D.9 and Lemma D.10, we have for any  $t \in [T_2, T_3]$ 

$$X_{s \in [T_2, t]} \eta Q\left(\frac{\alpha_1^{O(1)}}{d^{5/2}}\right) \le \sqrt[4]{\frac{1}{d}}$$

$$1$$

$$\text{polylog}(d)$$

and also

$$\sum_{s \in [T_2, t]} \eta \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4}) |E_{2,1}^{(t)}|^2 \Phi_2^{(t)} \leq \left( \sum_{s \in [T_2, T_{3,1}]} + \sum_{s \in [T_{3,1}, T_3]} \right) \eta \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4}) |E_{2,1}^{(t)}|^2 \Phi_2^{(t)} \\
\leq \eta \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4}) \cdot (T_{3,1} - T_2) \cdot O(\alpha_1^{O(1)} d^{3/8}) + \eta \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^4}) (T_3 - T_{3,1}) \\
\leq \widetilde{O}(\frac{\alpha_1^{O(1)}}{d^2})$$

Now we consider the term  $\Theta(\frac{(B_{1,2}^{(t)})^2}{(B_{2,2}^{(t)})^2})E_{2,1}^{(t)}\eta\Lambda_{2,2}^{(t)}$  , we have by Induction D.1a that

$$\left| \sum_{s \in [T_2,t]} \Theta(\frac{(B_{1,2}^{(t)})^2}{(B_{2,2}^{(t)})^2}) E_{2,1}^{(t)} \eta \Lambda_{2,2}^{(t)} \right| \leq O(\sqrt{\eta_E/\eta} (B_{1,2}^{(T_2)})^2) \sum_{s \in [T_2,t]} \eta \frac{|\Lambda_{2,2}^{(t)}|}{(B_{2,2}^{(t)})^2}$$

where we have used our induction hypothesis that  $B_{1,2}^{(t)}=B_{1,2}^{(T_2)}(1\pm o_{1})$ . Using Lemma G.3 by setting  $x_t=B_{2,2}^{(t)},\ q'=3$ , and  $A=\Theta(1)\geq d^{\Omega(1)}B_{2,2}^{(T_2)}$ , it holds that

$$\begin{array}{l} X \\ s \in [T_{2},t] \\ \theta(\frac{(1,2)}{(B_{2,2}^{(t)})^2}) E_{2,1}^{(t)} \eta \Lambda_{2,2}^{(t)} \leq O(^n \overline{\eta_E/\eta}) \frac{(1,2)}{|B_{2,2}^{(T_2)}|} \leq O(^n \overline{\eta_E/\eta}) \frac{(1,2)}{|B_{2,2}^{(0)}|} \leq \sqrt{\frac{1}{d}} \\ B & 2 \\ B & B(0) 2 \\ \end{array}$$

where in the second inequality we have used Lemma B.13c, Lemma C.8a and Lemma B.1, and in the last our choice of  $\eta_E/\eta \leq \rho_{\text{polylog}} 1_{(d)}$ . This ensures the induction can go on until  $t = T_3$ . And we finished our proof of Induction D.1a.

$$\begin{split} & \textbf{Proof of Induction D.1b:} \text{ Let us write down the update of}^{E_{1,2}^{(t)}} \text{ using Lemma D.7:} \\ & E_{1,2}^{(t+1)} = E_{1,2}^{(t)} (1 - \eta_E \Xi_1^{(t)}) + \sum_{\ell \in [2]} \Theta(\eta_E \Sigma_{1,\ell}^{(t)}) (-E_{1,2}^{(t)} [R_2^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2}) + \sum_{\ell \in [2]} \eta_E \Delta_{1,\ell}^{(t)} \\ & = E_{1,2}^{(t)} (1 - \eta_E \Xi_1^{(t)} - \sum_{\ell \in [2]} \Theta(\eta_E \Sigma_{1,\ell}^{(t)}) [R_2^{(t)}]^3) + \widetilde{O}(\frac{\eta_E}{d^{3/2}}) \Phi_1^{(t)} [R_1^{(t)}]^3 \\ & \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \sum_{\ell \in [2]} \eta_E \Sigma_{1,\ell}^{(t)} [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \\ & = E_{1,2}^{(t)} (1 - \eta_E \Xi_1^{(t)} - \Theta(\eta_E \Sigma_{1,1}^{(t)}) [R_2^{(t)}]^3) \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \eta_E \Sigma_{1,1}^{(t)} [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \end{split}$$

where in the last inequality we have used  $R_2^{(t)} \geq R_1^{(t)}$  from Induction D.1c and  $\Sigma_{1,1}^{(t)} \geq \Omega(\Phi_1^{(t)})$ ,

 $\Sigma_{2,1}^{(t)} \leq \widetilde{O}(rac{1}{d^{3/2}})\Sigma_{1,1}^{(t)}$  from Claim A.1 and Induction D.1a. Now we can use the same analysis in the proof of Lemma C.8 on $E_{1,2}^{(t)}$  to prove the desired claim, which we do not repeat here.

As for 
$$E_{2,1}^{(t)}$$
, we can obtain similar expressions: 
$$E_{2,1}^{(t+1)} = E_{2,1}^{(t)} (1 - \eta_E \Xi_2^{(t)} - \sum_{\ell \in [2]} \Theta(\eta_E \Sigma_{2,\ell}^{(t)}) [R_1^{(t)}]^3) \\ \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \sum_{\ell \in [2]} \Theta(\eta_E \Sigma_{2,\ell}^{(t)}) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} + \sum_{\ell \in [2]} \eta_E \Delta_{2,\ell}^{(t)}$$

Now we can obtain bounds for each terms as

$$\sum_{s \in [T_2,t]} \sum_{\ell \in [2]} \Theta(\eta_E \Sigma_{2,\ell}^{(s)}) [R_1^{(s)}]^3 \leq \widetilde{O}(\frac{\eta_E \alpha_1^{O(1)}}{d^2}) \cdot \widetilde{O}(\frac{d^{1.625} \alpha_1^{O(1)}}{\eta}) \leq \frac{1}{d^{3/4}}$$

and by (D.3) in Lemma D.9, we also have for any  $t \in [T_2, T_3]$ 

$$\sum_{s \in [T_2, t]} \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \sum_{\ell \in [2]} \Theta(\eta_E \Sigma_{2, \ell}^{(s)}) [R_1^{(s)}]^{3/2} [R_2^{(s)}]^{3/2} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \sum_{s \in [T_2, t]} \sum_{\ell \in [2]} \Theta(\eta_E \Sigma_{2, \ell}^{(s)}) [R_2^{(s)}]^3 \\
\leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) R_2^{(T_2)} \\
\leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

And also by using our induction and by (D.3) in Lemma D.9:

$$\sum_{s \in [T_2, t]} \sum_{\ell \in [2]} \eta_E \Delta_{2, \ell}^{(s)} \le \sum_{s \in [T_2, t]} \frac{\eta_E / \eta}{|E_{2, 1}^{(t)}|} \Theta(\eta \Sigma_{2, 1}^{(s)} + \eta \Sigma_{2, 2}^{(s)}) [R_2^{(s)}]^3 \le \frac{\eta_E / \eta}{|E_{2, 1}^{(T_2)}|} R_2^{(T_2)} \le O(\frac{\eta_E / \eta}{\log d}) = o(\sqrt{\eta_E / \eta})$$

Finally, we can calculate

$$\sum_{s \in [T_2, t]} \eta_E \Xi_2^{(t)} E_{2, 1}^{(t)} = \sum_{s \in [T_2, t]} \frac{\eta_E}{\eta} \frac{B_{2, 2}^{(t)}}{E_{2, 1}^{(t)}} \eta \Lambda_{2, 2}^{(t)}$$

By resorting to the defintion of  $T_3$  and go through similar analysis as for the induction of  $E_{1,1}^{(t)}$ , we can obtain that  $|E_{2,1}^{(t)}|$  is either above  $|E_{2,1}^{(T_2)}|$  (1 + o(1)) or is decreasing and always above  $\frac{1}{2}|E_{2,1}^{(T_2)}|$ . This proves Induction D.1b.

**Proof of Induction D.1c:** The proof of induction of  $R_2^{(t)}$  is half done in Lemma D.9, we only need to complete the part when  $t \in [T_{3,1}, T_3]$ , since by (D.3), we always have  $R_2^{(t)}$  to be decreasing by

And when  $t \in [T_{3,1}, T_3]$ , we have

X (s) 3/8+o(1) 
$$\Theta(\eta \Sigma_{2, \ \ } \leq {}^{O}_{e}(\eta d)$$
 \(\frac{1}{2}\)

So if we suppose  $R_2^{(T_3)} \leq \frac{1}{\sqrt{d}}$ , we shall have for  $T_3 - T_{3,1} = O(d^{1/4+o(1)}/\eta)$  many iterations that

$$R_2^{(t+1)} \ge R_2^{(T_{3,1})} (1 - \frac{\eta}{d^{5/8}})^{T_3 - T_{3,1}} \ge \Omega(R_2^{(T_{3,1})}) \ge \frac{1}{d^{1/4}}$$
 (by Lemma D.9)

So it negates our supposition, which completes the proof of the induction for  $R_2^{(t)}$  in  $t \in [T_2, T_3]$ .

Now we turn to the proof of induction for  $R_1^{(t)}$ , we write down its update: (as in Lemma C.8)  $R_1^{(t+1)} = R_1^{(t)} - \Theta(\eta[R_1^{(t)}]^3) \Big( \Sigma_{1,1}^{(t)} + \sum_{\ell \in \mathbb{N}} \Sigma_{2,\ell}^{(t)} (E_{2,1}^{(t)})^2 \Big)$ 

$$\begin{split} \text{It is straightforward to} & \ \pm O\Big(\sum_{j,\ell} \eta \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)}(R_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}\Big) \pm \frac{\eta}{\mathsf{poly}(d)} \\ & \ \sum_{\ell \in [2]} \Sigma_{1,\ell}^{(t)} |E_{1,2}^{(t)}| (\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})^2 \sum_{\ell \in [2]} \Sigma_{1,\ell}^{(t)} [R_1^{(t)}]^3 [R_2^{(t)}]^3 \end{split} \quad \text{derived}$$

and when  $t \in [T_2, T_{3,1}]$ :

$$\begin{split} \sum_{s \in [T_2, t]} \sum_{\ell \in [2]} \eta \Sigma_{2, \ell}^{(s)} |E_{2, 1}^{(s)}| (\overline{R}_{1, 2}^{(s)} + \varrho) [R_1^{(s)}]^{3/2} [R_2^{(s)}]^{3/2} &\leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \frac{d^{o(1)} d^{3/8}}{d^{9/8}} \sum_{s \in [T_2, t]} \sum_{\ell \in [2]} \eta \Sigma_{2, \ell}^{(s)} [R_2^{(s)}]^3 \\ &\leq o(\frac{d^{o(1)}}{d^{3/4}}) \end{split}$$

and when  $t \in [T_{3,1}, T_3]$ :

$$\sum_{s \in [T_2,t]} \sum_{\ell \in [2]} \eta \Sigma_{2,\ell}^{(s)} |E_{2,1}^{(s)}| (\overline{R}_{1,2}^{(s)} + \varrho) [R_1^{(s)}]^{3/2} [R_2^{(s)}]^{3/2} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \frac{d^{o(1)} d^{3/8}}{d^{9/8}} \eta \widetilde{O}(\frac{d^{1/4 + o(1)}}{\eta}) \leq O(\frac{1}{d})$$

So these combined with Lemma C.8 proved that  $R_1^{(t)} \leq O(\frac{d^{o(1)}}{d^{3/4}})$  for all  $t \in [T_2, T_3]$ . We can go through some similar analysis about  $R_2^{(t)}$  to get that  $R_1^{(t)} \geq \frac{1}{d}$  for all  $t \in [T_2, T_3]$ .

Finally we begin to prove the induction of  $\overline{R}_{1,2}^{(t)}$ . Similarly as in the proof of Lemma C.8, we first write down

$$\begin{split} R_{1,2}^{(t+1)} &= R_{1,2}^{(t)} - \eta \langle \nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} w_2^{(t)} \rangle - \eta \langle \nabla_{w_2} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} w_1^{(t)} \rangle \\ &+ \eta^2 \langle \Pi_{V^\perp} \nabla_{w_1} L(W^{(t)}, E^{(t)}), \Pi_{V^\perp} \nabla_{w_2} L(W^{(t)}, E^{(t)}) \rangle \\ &= R_{1,2}^{(t)} + \eta \Big( \Sigma_{1,1}^{(t)} + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} (E_{2,1}^{(t)})^2 \Big) (-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_1^{(t)}]^{5/2} [R_2^{(t)}]^{1/2} \\ &+ \eta \Big( \Sigma_{1,1}^{(t)} \Theta(E_{1,2}^{(t)})^2 + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} \Big) (-\Theta(\overline{R}_{1,2}^{(t)}) \pm O(\varrho)) [R_2^{(t)}]^{5/2} [R_1^{(t)}]^{1/2} \end{split}$$

Note that since 
$$+O\Big(\sum_{(j,\ell)\neq (1,2)} \eta \Sigma_{j,\ell}^{(t)} E_{j,3-j}^{(t)} (R_1^{(t)}[R_2^{(t)}]^2 + R_2^{(t)}[R_1^{(t)}]^2) \Big) \pm \frac{\eta}{\operatorname{poly}(d)} \\ |E_{1,2}^{(t)}| \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) [R_2^{(t)}]^{3/2} [R_1^{(t)}]^{3/2} \inf_{\mathsf{and}} R_1^{(t)} \leq O(\frac{1}{d^{3/4}}), \text{ it holds}$$

$$\begin{split} \sum_{(j,\ell)\neq(1,2)} \eta \Sigma_{j,\ell}^{(t)} |E_{j,3-j}^{(t)}| R_2^{(t)} [R_1^{(t)}]^2 &\leq \sum_{(j,\ell)\neq(1,2)} \eta \Sigma_{j,\ell}^{(t)} |E_{j,3-j}^{(t)}| R_1^{(t)} [R_2^{(t)}]^2 \\ &\leq o\left(\Sigma_{1,1}^{(t)} [R_1^{(t)}]^2 + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)} [R_2^{(t)}]^2\right) \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) [R_2^{(t)}]^{1/2} [R_1^{(t)}]^{1/2} \end{split}$$

so the update becomes

$$\begin{split} R_{1,2}^{(t+1)} &= R_{1,2}^{(t)} \left(1 - \eta \Theta\Big(\Sigma_{1,1}^{(t)} + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)}(E_{2,1}^{(t)})^2\Big) [R_1^{(t)}]^2 - \eta \Theta\Big(\Sigma_{1,1}^{(t)}(E_{1,2}^{(t)})^2 + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)}\Big) [R_2^{(t)}]^2 \right) \\ &\pm \eta \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) [R_1^{(t)}]^{1/2} [R_2^{(t)}]^{1/2} \Theta\Big(\Sigma_{1,1}^{(t)} + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)}(E_{2,1}^{(t)})^2\Big) [R_1^{(t)}]^2 \\ &\pm \eta \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) [R_1^{(t)}]^{1/2} [R_2^{(t)}]^{1/2} \Theta\Big(\Sigma_{1,1}^{(t)}(E_{1,2}^{(t)})^2 + \sum_{\ell \in [2]} \Sigma_{2,\ell}^{(t)}\Big) [R_2^{(t)}]^2 \end{split}$$

Now we can use the same arguments as in the proof of  $\overline{R}_{1,2}^{(t)}$  in Lemma C.8 to conclude.

**Proof of Lemma D.8a,b,c:** Indeed, at the end of phase III:

 $\begin{array}{lll} \mbox{Induction D.1a} & \Longrightarrow & \mbox{Lemma D.8a} \\ \mbox{Induction D.1b} & \Longrightarrow & \mbox{Lemma D.8c} \\ \mbox{Induction D.1c} & \Longrightarrow & \mbox{Lemma D.8b} \\ \end{array}$ 

Now we have completed the whole proof.

# E The End Phase: Convergence

When we arrive at  $t = T_3$ , we have already obtained the representation we want for the encoder network f(X), where  $v_1$  and  $v_2$  are satisfactorily learned by different neurons. In the last phase, we prove that such features are the solutions that the algorithm are converging to, which gives a stronger guarantee than just accidentally finding the solution at some intermediate steps.

To prove the convergence, we need to ensure all the good properties that we got through the training still holds. Fortunately, mosts of Induction D.1 still hold, as we summarized below:

**Inductions E.1.** At the end phase, i.e. when  $t \in [T_3,T]$ , Induction D.1a continues to hold except that  $|B_{2,2}^{(t)}| = \Theta(1)$ , Induction D.1b will hold except that for  $|E_{2,1}^{(t)}|$  only the upper bound still holds, and the upper bounds in Induction D.1c still hold while the lower bounds for  $R_1^{(t)}, R_2^{(t)}$  is 1/poly(d).

Moreover, there is a constant C = O(1) such that when  $t \geq T_3 + \frac{\alpha_1^C}{\eta}$ , we would have  $|E_{2,1}^{(t)}| \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}$ 

Now we present the main theorem of the paper, which we shall prove in this section.

**Theorem E.2** (End phase: convergence). For some  $T_4 = T_3 + \frac{d^{2+o(1)}}{\eta}$  and  $T = \text{poly}(d)/\eta$ , we have for all  $t \in [T_4, T]$  that Induction E.1 holds true and:

(a) Successful learning of both 
$$v_1, v_2 \colon |B_{1,1}^{(t)}|, |B_{2,2}^{(t)}| = \Theta(1)$$
 while  $|B_{2,1}^{(t)}|, |B_{1,2}^{(t)}| = \widetilde{O}(\frac{1}{\sqrt{d}})$ .

- (b) Successful denoising at the end:  $R_j^{(t)} \leq R_j^{(T_3)} (1 \widetilde{\Theta}(\tfrac{1}{\alpha_j^6})[R_j^{(t)}]^2) \text{ for all } j \in [2].$
- (c) Prediction head is close to identity:  $|E_{j,3-j}^{(t)}| \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2}[R_1^{(t)}]^{3/2}$  for all  $j \in [2]$ ; In fact, (b) and (c) also imply for some sufficiently large  $t = \text{poly}(d)/\eta$ , it holds  $R_j^{(t)} \leq \frac{1}{\text{poly}(d)}$  and  $|E_{j,3-j}^{(t)}| \le \frac{1}{\text{poly}(d)}$  for all  $j \in [2]$ .

And we have a simple corollary for the objective convergence.

Corollary E.3 (objective convergence, with prediction head). Let OPT denote the global minimum of the population objective (A.1). It is easy to derive that  $OPT = 2 - 2\frac{C_0}{C_1} = \Theta(\frac{1}{\log d})$ . We have for some sufficiently large  $t \ge poly(d)/\eta$ :

$$L(W(t),E(t)) \leq \mathsf{OPT} + \frac{1}{\mathsf{poly}(d)}$$

Now we need to establish some auxiliary lemmas:

**Lemma E.4.** For some  $t \in [T_3, poly(d)/\eta]$ , if Induction E.1 holds from  $T_3$  to t, we have Lemma D.6 holds at t.

*Proof.* Simple from similar calculations in the proof of Lemma D.6.

**Lemma E.5.** For some  $t \in [T_3, poly(d)/\eta]$ , if Induction E.1 holds from  $T_3$  to t, we have for each  $j \in [2]$ that  $X X \eta \Sigma(s)[R(s)] 3 \le O(R_j(T_3)), \forall j \in [2]_{j,`}$  $s \in [T_3, t] \in [2]$ 

*Proof.* Notice that when Induction E.1 holds, we always have

$$X(\Sigma(j,t) + \Sigma(3t-)j,(E_3(t-)j,j)2) = (1 \pm o(1)) X \Sigma(j,t)$$
  
` $\in$ [2]

we can use Lemma E.4 to obtain the update of  $R_2^{(t)}$  as in the calculations when we obtained (D.3):

$$R2(t) = R2(T_3) - X X \Theta(\eta \Sigma 2(s, \gamma)) [R2(s)] 3$$

$$s \in [T_3, t) \in [2]$$

which means that  $R_2^{(t)}$  is decreasing from  $T_3$  to t. Summing up the update, the part of  $R_2^{(t)}$  is solved. For the part of  $R_1^{(t)}$ , we separately discuss when  $|E_{2,1}^{(t)}|$  is larger than or smaller than  $\widetilde{O}(\varrho+\frac{1}{\sqrt{d}})[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2}$ . When the former happens, which we know from Induction E.1 that it

cannot last until some  $t_4'=T_3+rac{lpha_1^{O(1)}}{\eta}$  many iterations, we have for  $t\in[T_3,t_4']$ 

$$\sum_{s \in [T_3,t)} \sum_{(j,\ell) \in [2]^2} \eta \Sigma_{j,\ell}^{(s)} |E_{j,3-j}^{(s)}| (\overline{R}_{1,2}^{(s)} + \varrho) [R_1^{(s)}]^{3/2} [R_2^{(s)}]^{3/2} \leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \frac{\alpha_1^{O(1)}}{d} R_1^{(T_3)} \leq \frac{1}{d} R_1^{(T_3)}$$

Now for  $t \ge t_4'$  we can simply go through similar calculations as in the proof of Induction D.1c to obtain

$$\begin{split} \sum_{s \in [t_4', t)} \sum_{(j, \ell) \in [2]^2} \eta \Sigma_{j, \ell}^{(s)} |E_{j, 3-j}^{(s)}| (\overline{R}_{1, 2}^{(s)} + \varrho) [R_1^{(s)}]^{3/2} [R_2^{(s)}]^{3/2} &\leq \sum_{s \in [t_4', t)} \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})^2 \sum_{(j, \ell) \in [2]^2} \eta \Sigma_{j, \ell}^{(s)} [R_1^{(s)}]^3 [R_2^{(s)}]^3 \\ &\leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})^2 R_2^{(T_3)} \max_{s \in [t_4', t)} [R_1^{(s)}]^3 \\ &\leq \frac{1}{d} R_1^{(T_3)} \end{split}$$

So by applying Lemma E.4a and Lemma D.6, we have

$$R1(t) = (1 \pm o(1))R1(T_3) - X \times \Theta(\eta \Sigma(j, s))[R1(s)]3$$

$$s \in [T_3, t) \in [2]$$

which proves the claim.

**Lemma E.6.** For some  $t \in [T_3, \mathsf{poly}(d)/\eta]$ , if Induction E.1 holds from  $T_3$  to t. Then we have  $|E_{j,3-j}^{(t)}|$  is decreasing until  $|E_{j,3-j}^{(t)}| \leq O(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2} + \widetilde{O}(\frac{1}{d^{3/2}})[R_j^{(t)}]^3$ . Moreover, we have for each

$$\left| \sum_{s \in [T_3, t]} \eta_E \Xi_j^{(t)} E_{j, 3-j}^{(s)} \right| \le |E_{j, 3-j}^{(T_3)}| + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \le O(\sqrt{\eta_E/\eta})$$

*Proof.* We can go through the same calculations in the proof of Induction D.1b (using Fact D.2) to obtain

$$\begin{split} E_{j,3-j}^{(t+1)} &= E_{j,3-j}^{(t)} (1 - \eta_E \Xi_j^{(t)}) + \sum_{\ell \in [2]} \eta_E \Delta_{j,\ell}^{(t)} \\ &+ \sum_{\ell \in [2]} \Theta(\eta_E \Sigma_{j,\ell}^{(t)}) (-E_{j,3-j}^{(t)} [R_{3-j}^{(t)}]^3 \pm O(\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2}) \\ &= E_{j,3-j}^{(t)} (1 - \eta_E \Xi_j^{(t)} - \eta_E \Theta(\Sigma_{j,j}^{(t)} [R_{3-j}^{(t)}]^3)) + \widetilde{O}(\frac{1}{d^{3/2}}) \sum_{\ell \in [2]} \eta_E \Sigma_{j,\ell}^{(t)} [R_j^{(t)}]^3 \\ &\pm O(\eta_E \Sigma_{j,j}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \end{split}$$

where we have used in the second equality that  $\sum_{\ell \in [2]} \Delta_{j,\ell}^{(t)} \leq \widetilde{O}(\frac{1}{d^{3/2}}) \sum_{\ell \in [2]} \Sigma_{j,\ell}^{(t)} [R_j^{(t)}]^3$  and also

 $\Sigma_{j,3-j}^{(t)} \leq O(\tfrac{1}{d^{3/2}}) \Sigma_{j,j}^{(t)} \text{ for both } j \in [2] \text{ when Induction E.1 holds. Note that from above calculations, there exist a constant C such that if } -|E_{j,3-j}^{(t)}| \geq C(R_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2} + \sum_{\ell \in [2]} \eta_E \Delta_{j,\ell}^{(t)}$ , we have  $|E_{2,1}^{(t)}|$  to be decreasing. Now it

$$\begin{split} \sum_{s \in [T_3,t]} O(\eta_E \Sigma_{j,j}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} &\leq \sum_{s \in [T_3,t]} O(\eta_E \Sigma_{1,1}^{(t)} + \eta_E \Sigma_{2,2}^{(t)}) (\overline{R}_{1,2}^{(t)} + \varrho) ([R_1^{(t)}]^3 + [R_2^{(t)}]^3) \\ &\leq \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \end{split}$$

Also note that  $\sum_{j,j}^{(t)}[R_{3-j}^{(t)}]^3 \leq$ 

which is from Induction E.1, Induction D.1c and Lemma E.4.  $O(\frac{d^{o(1)}}{d^{3/4}})\Xi_j^{(t)}$  at this stage, we have

$$E_{3-j,j}^{(t)} = E_{j,3-j}^{(T_3)} - \sum_{s \in [T_3,t)} \Xi_j^{(s)} E_{j,3-j}^{(s)} + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

Recalling the expression of  $\Xi_i$  finishes the proof.

**Lemma E.7.** Recall  $T_2$  defined in (C.1) and  $T_3$  defined in (D.1), we have

$$\sqrt{\eta/\eta_E} \max_{t \le T_3} |E_{2,1}^{(t)}| \le \sum_{t \le T_2} \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} + \frac{1}{\alpha_1^{\Omega(1)}}$$

To prove this lemma, we need a simple claim.

**Claim E.8.** If  $\{x_t\}_{t < T}, x_t \ge 0$  is an increasing sequence and  $C = \Theta(1)$  is a constant such that  $\sqrt{\frac{1}{2}}$ 

$$\begin{aligned} x_{t+1} - x_t &\leq O(\eta) \ and \ \Pr_{t < T} x_t (x_{t+1} - x_t) = C, \ then \ for \ each \ \delta \in \binom{t}{d}, 1 \ it \ holds \ |x_T - C| \leq O(\delta^2 + x_0^2 + O(\frac{\log_d}{d})) \end{aligned}$$

*Proof.* Indeed, for every  $g \in 0,1,...$ , we define  $T_g := \min\{t : x_t \ge (1 + \delta)^g x_0\}$ . and define  $b := \min\{g : ((1 + \delta)^g x_0)^2 \ge C - \delta^2\}$ . Now for any g < b, we have

$$\sum_{t \in [\mathcal{T}_{g}, \mathcal{T}_{g+1}]} x_{t}(x_{t+1} - x_{t}) \ge x_{\mathcal{T}_{g}}(x_{\mathcal{T}_{g+1}} - x_{\mathcal{T}_{g}}) \ge (1 + \delta)^{g} \delta(1 + \delta)^{g-1} x_{0}^{2} - \frac{1}{d} = \delta(1 + \delta)^{2g-1} x_{0}^{2} - \frac{1}{d}$$

By our definition of  $T_g$ , we can further get

$$C = \sum_{t < T} x_t (x_{t+1} - x_t) = \sum_{g=1}^b \sum_{t \in [\mathcal{T}_0, \mathcal{T}_{g+1}]} x_t (x_{t+1} - x_t) \ge (1 + \delta)^{2b} x_0^2 - x_0^2 - \frac{b}{d} \ge C - \delta^2 - x_0^2 - \frac{b}{d}$$

And also we have  $C \le (\max_{t \le T} x_t)^{\mathbf{P}_{t < T}} (x_{t+1} - x_t) = x_T^2$ , so we have  $|x_T^2 - C| \le \delta^2 + x_0^2 + \frac{b}{d}$ , where  $b = O(\log(C)/\log(1+\delta)) \le O(\log d)$ , which proves the claim.  $\square$ 

Proof of Lemma E.7. From the proof of Lemma C.8 and Lemma D.8 we know that

$$\max_{t \leq T_3} |E_{2,1}^{(t)}| \leq \sum_{t \leq T_3} (1 \pm \frac{1}{\alpha_1^{\Omega(1)}}) \eta_E |\Delta_{2,1}^{(t)}| + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

And since from the proof of Lemma C.8 we know that

$$\begin{split} R_2^{(T_3)} &= R_2^{(0)} - \sum_{t \le T_3} (1 \pm \widetilde{O}(\frac{1}{d^{3/2}})) \eta \Sigma_{2,1}^{(t)} \mathcal{E}_{2,1}^{(t)} \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \\ &= (1 \pm \widetilde{O}(\frac{1}{d^{3/2}})) \sum_{t \le T_3} E_{2,1}^{(t)} \Delta_{2,1}^{(t)} \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \end{split}$$

We can define some alternative variables  $\widetilde{E}_{2,1}^{(t)}$  updated as  $\widetilde{E}_{2,1}^{(t+1)}=\widetilde{E}_{2,1}^{(t)}+\eta_E\Delta_{2,1}^{(t)}$  and  $\widetilde{R}_2^{(t+1)}=$ 

 $\widetilde{R}_{2}^{(t)} - \widetilde{E}_{2,1}^{(t)} \Delta_{2,1}^{(t)}$ . It is easy to see that  $|E_{2,1}^{(t)} - \widetilde{E}_{2,1}^{(t)}| \le \frac{1}{\alpha_{1}^{\Omega(1)}} \max_{t \le T_{3}} |E_{2,1}^{(t)}|$ . From above calculations, we know  $\frac{\eta}{\eta_E} \sum_{t \in [T_1, T_3]} \widetilde{E}_{2,1}^{(t)} (\widetilde{E}_{2,1}^{(t+1)} - \widetilde{E}_{2,1}^{(t)}) = \widetilde{R}_2^{(T_1)} \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) + O(\frac{1}{d^{1/4}})$ , which by Claim E.8 implies

$$\sqrt{\eta/\eta_E}|\widetilde{E}_{2,1}^{(T_3)}| = \sqrt{\widetilde{R}_2^{(T_1)}} \pm O(\frac{1}{d^{1/4}}) = \sqrt{2} \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \pm O(\frac{1}{d^{1/4}})$$

And when we turn back, we shall have -

 $\sqrt{\eta/\eta_E}\max_{t\leq T_3}|E_{2,1}^{(t)}|\leq \sqrt{2}+rac{1}{lpha_1^{\Omega(1)}}$  . Now we can use similar

techniques on  $B_{1,1}^{(t)}$  and  $B_1^{(t)}$ . Indeed, from (C.4) and similar arguments in phase I, we know for all  $t \in \mathbb{R}$  $[T_1, T_2]$ 

$$\begin{split} R_1^{(t+1)} &= R_1^{(0)} - \sum_{s \leq t} (1 \pm \widetilde{O}(\frac{1}{d^{3/2}})) \eta \Sigma_{1,1}^{(s)} \mathcal{E}_{1,2}^{(s)} \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}} \\ R_1^{(t+1)} &\leq R_1^{(t)} (1 - \widetilde{O}(\frac{\eta}{\alpha_1^6}) [R_1^{(t)}]^2) \end{split}$$
 (E.1)

So one can obtain that at some iteration  $t' = T_1 + O(\frac{d\alpha_1^{O(1)}}{\eta})$ , we shall have  $R_1^{(t)} \leq O(\frac{1}{\sqrt{d}})$  for all  $t \geq t^0$ .

Now let us consider the growth of 
$$B_{1,1}^{(t)}$$
 before  $t^0$ , which clearly constitutes of 
$$B_{1,1}^{(t')} = B_{1,1}^{(T_1)} + \sum_{t \in [T_1,t')} (\Lambda_{1,1}^{(t)} + \Gamma_{1,1}^{(t)} - \Upsilon_{1,1}^{(t)})$$
 
$$= B_{1,1}^{(T_1)} + \sum_{t \in [T_1,t')} \left( \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} \mathrm{sign}(B_{1,1}^{(t)}) + \eta \Gamma_{1,1}^{(t)} - \eta \Upsilon_{1,1}^{(t)} \right)$$
 
$$= B_{1,1}^{(0)} + \sum_{t \in [T_1,t')} \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} \mathrm{sign}(B_{1,1}^{(t)}) + \sum_{t \in [T_1,t')} \eta \left( \Gamma_{1,1}^{(t)} - \Upsilon_{1,1}^{(t)} \right) + \widetilde{O}(\frac{1}{\sqrt{d}})$$

where the last one comes from the proof of Lemma B.13. Moreover by using the same arguments in the proof of Lemma C.8 we can easily prove that

$$\left| \sum_{t \in [T_1, t')} (\Gamma_{1,1}^{(t)} - \Upsilon_{1,1}^{(t)}) \right| \leq \widetilde{O}(\frac{1}{\sqrt{d}}) \quad \Rightarrow \quad \sum_{t < t'} \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} \geq |B_{1,1}^{(t')}| - |B_{1,1}^{(0)}| - \widetilde{O}(\frac{1}{\sqrt{d}})$$

And for  $t \in [t^0, T_2]$ , we also have by (E.1) that

$$\sum_{t \in [t', T_2]} \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} \le \sum_{t \in [t', T_2)} \eta \Sigma_{1,1}^{(t)} \mathcal{E}_{1,2}^{(t)} \le O(\frac{1}{\sqrt{d}})$$

 $\text{Recall } R_1^{(0)} = \textstyle \sum_{t \in [0,t')} (1 \pm \widetilde{O}(\frac{1}{d^{3/2}})) \eta \Sigma_{1,1}^{(t)} \mathcal{E}_{1,2}^{(t)} \pm \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \text{ by (E.1) and } R_1^{(t)} \leq O(\frac{1}{\sqrt{d}}) \text{ for } t \geq t^0.$ 

Now we can finally go through the same analysis using Claim E.8 on  $B_{1,1}^{(t)}$  and  $B_{1}^{(t)}$  during  $t \in [0,t^0]$  as above to obtain that

$$\sum_{t \le T_2} \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} \ge (1 - \widetilde{O}(\frac{1}{d^{3/2}})) \sqrt{R_1^{(0)}} - \widetilde{O}(\frac{1}{\sqrt{d}}) = 1 - \widetilde{O}(\varrho + \frac{1}{\sqrt{d}})$$

Combining the results, we finishes the proof.

Now we are prepared to prove Theorem E.2.

### **E.1** Proof of Convergence

*Proof of Theorem E.2.* First we start with the  $B_{j,\ell}^{(t)}$ s. Indeed, we can go through similar calculations to see that all gradients  $h-\nabla_{w_j}L(W^{(t)},E^{(t)})$ , v i can be decomposed into

$$\langle -\nabla_{w_j} L(W^{(t)}, E^{(t)}), v_\ell \rangle = (\Lambda_{j,\ell}^{(t)} - \Upsilon_{j,\ell,1}^{(t)}) + (\Gamma_{j,\ell}^{(t)} - \Upsilon_{j,\ell,2}^{(t)})$$

$$\begin{split} \text{where } \Lambda_{j,\ell}^{(t)} &- \Upsilon_{j,\ell,1}^{(t)} \text{ and } \Gamma_{j,\ell}^{(t)} - \Upsilon_{j,\ell,2}^{(t)} \text{ can be expressed as} \\ \Lambda_{j,\ell}^{(t)} &- \Upsilon_{j,\ell,1}^{(t)} = C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_j^{(t)} (B_{j,\ell}^{(t)})^5 \left( E_{j,3-j}^{(t)} (B_{3-j,3-\ell}^{(t)})^3 (B_{j,3-\ell}^{(t)})^3 + (E_{j,3-j}^{(t)})^2 (B_{3-j,3-\ell}^{(t)})^6 \right) \\ &- C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_j^{(t)} (B_{j,\ell}^{(t)})^2 (B_{3-j,\ell}^{(t)})^3 E_{j,3-j}^{(t)} \left( (B_{j,3-\ell}^{(t)})^6 + E_{j,3-j}^{(t)} (B_{3-j,3-\ell}^{(t)})^3 (B_{j,3-\ell}^{(t)})^3 \right) \\ &+ C_0 \alpha_2^6 \Phi_j^{(t)} (B_{j,\ell}^{(t)})^5 C_2 \mathcal{E}_{j,3-j}^{(t)} \\ &\Gamma_{j,\ell}^{(t)} - \Upsilon_{j,\ell,2}^{(t)} = C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_{3-j}^{(t)} (B_{3-j,\ell}^{(t)})^3 (B_{j,\ell}^{(t)})^2 E_{3-j,j}^{(t)} \left( E_{3-j,j}^{(t)} (B_{j,3-\ell}^{(t)})^3 (B_{j,3-\ell}^{(t)})^3 (B_{3-j,3-\ell}^{(t)})^3 (B_{j,3-\ell}^{(t)})^3 \right) \\ &- C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_{3-j}^{(t)} (B_{j,\ell}^{(t)})^5 (E_{3-j,j}^{(t)})^2 \left( (B_{3-j,3-\ell}^{(t)})^6 + E_{3-j,j}^{(t)} (B_{j,3-\ell}^{(t)})^3 (B_{3-j,3-\ell}^{(t)})^3 \right) \\ &+ C_0 \alpha_2^6 \Phi_{3-j}^{(t)} E_{3-j,j}^{(t)} (B_{3-j,\ell}^{(t)})^3 (B_{j,\ell}^{(t)})^2 C_2 \mathcal{E}_{3-j,j}^{(t)} \end{split}$$

Firstly, for all the terms that contain factors of  $(B_{j,\ell}^{(t)})^2(B_{3-j,\ell}^{(t)})^2$  (or  $(B_{j,\ell}^{(t)})^2(B_{j,3-\ell}^{(t)})^2$ ), we can apply Lemma E.6, our Induction E.1 assumption and  $|E_{j,3-j}^{(t)}| \leq O(1), \forall t \in [T_3,T]$  to obtain that their (multiplicated by  $\eta$ ) summation over  $t \in [T_3,T]$  is absolutely bounded by  $O_{\mathbf{C}}(d^1)$ . So we can move on to deal with all other terms. When j = 1, Using Lemma E.6, we have

$$\sum_{t \in [T_3, T]} \eta C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_j^{(t)} |B_{j,\ell}^{(t)}|^5 (E_{j,3-j}^{(t)})^2 (B_{3-j,3-\ell}^{(t)})^6 = \sum_{t \in [T_3, T]} \frac{\eta \Xi_j^{(t)}}{|B_{j,\ell}^{(t)}|} (E_{j,3-j}^{(t)})^2 \\ \leq \sqrt{\frac{\eta}{\eta_E}} |E_{j,3-j}^{(T_3)}| + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) = O(1)$$

And the sign of LHS is  $sign(B_{j,\ell}^{(t)})$ . Moreover, for j = 1, from Lemma E.7 and Lemma E.6 we also have

$$\begin{split} \sum_{t \in [T_3, T]} \eta C_0 \alpha_2^6 C_1 \alpha_1^6 \Phi_2^{(t)} |B_{1,1}^{(t)}|^5 (E_{2,1}^{(t)})^2 (B_{2,2}^{(t)})^6 &\leq \sqrt{\frac{\eta}{\eta_E}} \left| \sum_{t \in [T_3, T]} \eta_E \Xi_j^{(t)} E_{j, 3-j}^{(t)} \right| \\ &\leq \sqrt{\frac{\eta}{\eta_E}} |E_{2,1}^{(T_3)}| + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) \\ &\leq \sum_{t \leq T_2} \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} + \frac{1}{\alpha_1^{\Omega(1)}} \end{split}$$

Since we have

$$B_{1,1}^{(T_2)} = \sum_{s < T_2} \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} + \sum_{s < T_2} \frac{\eta \Sigma_{2,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{2,1}^{(t)} - \sum_{t \in [T_3,T]} \frac{\eta \Xi_j^{(t)}}{|B_{j,\ell}^{(t)}|} (E_{j,3-j}^{(t)})^2$$

And since by Induction C.1 we have  $= |B_{1,1}^{(t)}| \quad \Theta(1)$  during  $t \in [T_1, T_2]$  and  $\sum_{t \in [T_1, T_2]} \eta \Sigma_{2,1}^{(t)} \geq \sqrt{1 - t}$ 

 $R^{(T_1)} - o(1) = 2 - o(1)$ . For all the other terms in the gradient, we can apply Lemma E.6, our Induction E.1 assumption and  $|E_{j,3-j}^{(t)}| \leq O(1)$  so we have for  $t \in [T_3,T]$ 

$$\begin{split} |B_{1,1}^{(t)}| &= \sum_{s \leq T_2} \frac{\eta \Sigma_{1,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{1,2}^{(t)} + \sum_{s \leq T_2} \frac{\eta \Sigma_{2,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{2,1}^{(t)} - \sum_{t \in [T_3,T]} \frac{\eta \Xi_j^{(t)}}{|B_{j,\ell}^{(t)}|} (E_{j,3-j}^{(t)})^2 - o(1) \\ &\geq \sqrt{\eta/\eta_E} \max_{t \leq T_3} |E_{2,1}^{(t)}| + \sum_{s \leq T_2} \frac{\eta \Sigma_{2,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{2,1}^{(t)} - \sqrt{\frac{\eta}{\eta_E}} |E_{j,3-j}^{(T_3)}| + \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) - o(1) \\ &\geq \sum_{s \leq T_2} \frac{\eta \Sigma_{2,1}^{(t)}}{|B_{1,1}^{(t)}|} \mathcal{E}_{2,1}^{(t)} - o(1) \geq \Omega(1) \end{split}$$

which also proved  $|B_{1,1}^{(t)}| = O$  (1) since all the terms on the RHS are absolutely O(1) bounded. Since one can see from Lemma E.6 that  $|E_{2,1}^{(t)}|$  is decreasing before it reaches  $d^1$ ). Moreover this proves  $\sqrt{\eta/\eta_E}|E_{2,1}^{(t)}| \leq B_{1,1}^{(t)}$  for all  $t \in [T_3,T]$ , and also the fact that

$$B_{1,1}^{(t)} \ge \Omega(1), \quad \forall t \in [T_3, T]$$

The case of  $B_{2,2}^{(t)}$  is much more simple as  $E_{1,2}^{(t)} \leq \widetilde{O}(\frac{1}{d})$  throughout  $t \in [T_3,T]$  by Lemma E.6 and Lemma D.8c, Now we can go through the similar calculations again to obtain that  $B_{2,2}^{(t)} = \Theta(1)$  for all  $t \in [T_3,T]$ . When j 6= `, all the terms calculated in the expansion of  $\Lambda^{(j,t)} - \Upsilon^{(j,t)} = 1$  and  $\Gamma^{(j,t)} - \Upsilon^{(j,t)} = 1$  contain factors of  $(B_{2,1}^{(t)})^2 = \widetilde{O}(\frac{1}{d})$  or  $(B_{1,2}^{(t)})^2 = \widetilde{O}(\frac{1}{d})$ . So we can similarly use Lemma E.6 as before to derive that  $B_{j,3-j}^{(t)} = B_{j,3-j}^{(T_3)} (1 \pm \widetilde{O}(\frac{\alpha_1^{O(1)}}{\sqrt{d}}))$  for all  $t \in [T_3,T]$  and  $j \in [2]$ .

As for the prediction head, the induction of  $E_{1,2}^{(t)}$  follows from exactly the same proof in Lemma D.8. The part of  $E_{2,1}^{(t)}$  is half done in Lemma E.6. It suffices to notice that  $\Xi^2 = \widetilde{\Theta}(\frac{\alpha_1^6}{\alpha_2^6})$  and if

$$\begin{split} |E_{2,1}^{(t)}| &\geq C(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2} \text{ for some } \textit{C} = \textit{O}(\textbf{1}) \text{, then} \\ E_{2,1}^{(t+1)} &= E_{2,1}^{(t)}(1 - \eta_E\Xi_2^{(t)} - \eta_E\Theta(\Sigma_{2,2}^{(t)}[R_1^{(t)}]^3)) + \widetilde{O}(\frac{1}{d^{3/2}})\sum_{\ell \in [2]} \eta_E\Sigma_{2,\ell}^{(t)}[R_2^{(t)}]^3 \\ & \pm O(\eta_E\Sigma_{2,2}^{(t)})(\overline{R}_{1,2}^{(t)} + \varrho)[R_1^{(t)}]^{3/2}[R_2^{(t)}]^{3/2} \\ &\leq E_{2,1}^{(t)}(1 - \widetilde{\Theta}(\frac{\eta\alpha_1^6}{\alpha_2^6})) \end{split}$$

So after  $\frac{\alpha_1^{O(1)}}{\eta}$  many epochs will we have

$$|E_{2,1}^{(t)}| \le (\log d) |\overline{R}_{1,2}^{(t)} + \varrho|[R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2} \le \widetilde{O}(\varrho + \frac{1}{\sqrt{d}}) [R_1^{(t)}]^{3/2} [R_2^{(t)}]^{3/2}$$

as desired. And the rest of the induction of  $E_{2,1}^{(t)}$  is the same as in the induction arguments of  $E_{1,2}^{(t)}$  in Lemma D.8.

The induction of  $R_1^{(t)}, R_2^{(t)}$  and  $R_{1,2}^{(t)}$  is exactly the same as those in the proof of Lemma D.8 except here we only need  $R_1^{(t)}/R_2^{(t)} \in [\frac{1}{\alpha_1^{O(1)}}, \alpha_1^{O(1)}]$  after  $T_4$ . Indeed, from the update of  $R_j^{(t)}$  (which can be easily worked out), we have

$$R_j^{(t+1)} = R_j^{(t)} (1 - \Theta(\eta \Sigma_{j,j}^{(t)}) [R_j^{(t)}]^2) = R_j^{(t)} (1 - \widetilde{\Theta}(\frac{\eta}{\alpha_j^6}) [R_j^{(t)}]^2)$$

Now after  $\frac{d^2\alpha_1^{O(1)}}{\eta}$  many epochs, we can obtain from similar arguments in Lemma D.8 that  $R_1^{(t)}/R_2^{(t)} \in \left[\frac{1}{\alpha_1^{O(1)}}, \alpha_1^{O(1)}\right]$  and  $R_j^{(t)} \leq \frac{1}{d}$ . The induction can go on untill  $t = \text{poly}(d)/\eta$ .

For the convergence of  $B_{1,1}^{(t)}$  and  $B_{2,2}^{(t)}$  after  $t=T_4$ , notice that their change depends on  $\sum_{t\geq T_4} \frac{E_{j,3-j}^{(t)}}{B_{j,j}^{(t)}} \Xi_j^{(t)}$ , which stays very small after  $T_4$ , we have that  $|B_{j,j}^{(t)}-B_{j,j}^{(T_4)}|\leq o_{\text{(1)}} \text{ for all } j\in \text{[2]}.$  This finishes the whole proof.

## F Learning Without Prediction Head

When we do not use prediction head in the network architecture, the analysis is much simpler. We can reuse most of the gradient calculations in previous sections as long as we set  $E^{(t)}$  to the identity. Note that here we allow  $m \ge 1$  to be any positive integer.

**Theorem F.1** (learning without the prediction head). Let m be any positive integer. If we keep  $E^{(t)} \equiv I_m$  during the whole training process, then for all  $t \in [\widetilde{\Omega}(\frac{d^2}{\eta})]$ , poly $(d)/\eta$ , we shall have  $|B_{j,1}^{(t)}| = \Theta(1), \ |B_{j,2}^{(t)}| = \widetilde{O}(\frac{1}{\sqrt{d}})$  and  $R_j^{(t)} = O(\frac{1}{d^{1-o(1)}})$  for all  $j \in [m]$  with probability 1 - o(1).

Moreover, for a longer training time  $t = \text{poly}(d)/\eta$ , we would have  $R_j^{(t)} \leq \frac{1}{\text{poly}(d)}$  for all  $j \in [m]$ .

Moreover, it is direct to obtain a objective convergence result similar to Corollary E.3.

**Corollary F.2** (objective convergence, without prediction head). Let OPT denote the global minimum of the population objective (A.1). When trained with  $E^{(t)} \equiv I_m$ , we have for some sufficiently large  $t \ge \text{poly}(d)/\eta$ :

$$L(W^{(t)},I_m) \leq \mathsf{OPT} + \frac{1}{\mathsf{poly}(d)}$$

*Proof of Theorem F.1.* The proof is easy to obtain since it is very similar to some proofs in previous sections, and we only sketch it here. Indeed, using the calculations in Lemma D.5 and Lemma D.4

and 
$$\operatorname{set}^{E_{i,j}^{(t)}}, i \neq j \in [m]$$
 to zero. We shall have (note that here  $\mathcal{E}_{j,r}^{(t)} \equiv \mathcal{E}_{j}^{(t)}$  for any  $r$  6=  $j$ )  $\langle -\nabla_{w_{j}}L(W^{(t)}, E^{(t)}), v_{\ell} \rangle = C_{0}C_{2}\alpha_{\ell}^{6}(B_{j,\ell}^{(t)})^{5}\Phi_{j}^{(t)}\mathcal{E}_{j}^{(t)} = \Theta(C_{0}C_{2}\alpha_{\ell}^{6}\Phi_{j}^{(t)}(B_{j,\ell}^{(t)})^{5}[R_{j}^{(t)}]^{3})$ 

Now we can go through the similar induction arguments as in the proof of Lemma B.13 (with TPM lemma to distinguish the learning speed) to obtain that for each  $j \in [m]$ :

$$|B_{j,1}^{(t)}| = \Theta(1), \quad |B_{j,2}^{(t)}| = |B_{j,2}^{(0)}|(1 \pm o(1)), \quad \forall j \in [m]$$
 (when  $t \ge \frac{d^2}{\eta}$ )

When this is proven, we can also reuse the calculations as in the proof of Lemma C.5 to obtain that

$$R_j^{(t+1)} = R_j^{(t)} (1 - \Theta(\eta \Sigma_{j,1}^{(t)}) [R_j^{(t)}]^2) = R_j^{(t)} (1 - \Theta(\eta C_0 C_2 \alpha_1^6 \Phi_j^{(t)} (B_{j,1}^{(t)})^6 [R_j^{(t)}]^2), \quad \forall j \in [m]$$

So again after some  $t=\widetilde{O}(\frac{d^2}{\eta})$ , we shall have  $R_j^{(t)} \leq O(\frac{d^{o(1)}}{d})$ . While the decrease of  $R_j^{(t)}$  is happening, we can make induction that  $|B_{j,2}^{(t)}| = |B_{j,2}^{(0)}| (1 \pm o(1))$ , since if it holds for all previous iterations before t, then

$$X \eta |h - \nabla_{j} L(W(s), E(s)), v i| = X \eta Co\alpha_{26} \Phi_{j(s)} |B_{j,(s2)}| 5C_{2} E_{j(s)} w$$

$$s \le t - 1$$

$$\leq \sup_{polylog(d)} |B_{j,2}^{(0)}|$$

$$1$$

where ¬ is due to Corollary G.2, where 
$$x_t = |B_{j,1}^{(t)}|_{\text{and }} y_t = |B_{j,2}^{(t)}|_{\text{and }} S_t \le \frac{1}{1 \cdot v} \le \frac{1}{1 \cdot v}$$

 $O(\log d)x_0$ . which finishes the proof.

### **G** Tensor Power Method Bounds

In this section, we give two lemmas related to the tensor power method that can help us in previous sections' proofs.

**Lemma G.1** (TPM, adapted from [3]). Consider an increasing sequence  $x_t \ge 0$  defined by  $x_{t+1} = x_t + \eta C_t x^q t$  for some integer  $q \ge 3$  and  $C_t > 0$ , and suippose for some A > 0 there exist  $t^0 \ge 0$  such that  $x_{t^0} \ge A$ . Then for every  $\delta > 0$ , and every  $\eta \in (0,1)$ :

$$\sum_{t \ge 0, x_t \le A} \eta C_t \ge \left( \frac{\delta (1+\delta)^{-1}}{(1+\delta)^{q-1} - 1} \left( 1 - \left( \frac{(1+\delta)x_0}{A} \right)^{q-1} \right) - \frac{O(\eta A^q)}{x_0} \frac{\log(A/x_0)}{\log(1+\delta)} \right) \cdot \frac{1}{x_0^{q-1}}$$

$$\sum_{t \ge 0, x_t \le A} \eta C_t \le \left( \frac{(1+\delta)^{q-1}}{q-1} + \frac{O(\eta A^q)}{x_0} \frac{\log(A/x_0)}{\log(1+\delta)} \right) \cdot \frac{1}{x_0^{q-1}}$$

This lemma has a corollary:

**Corollary G.2** (TPM, from [3]). Let  $q \ge 3$  be a constant and  $x_0, y_0 = o(1)$  and A = O(1). Let  $\{x_t, y_t\}_{t\ge 0}$  be two positive sequences updated as

- $x_{t+1} = x_t + \eta C_t x_t^q$  for some  $C_t > 0$ ;
- $y_{t+1} = y_t + \eta S_t C_t y_t^q$  for some  $S_t > 0$ .

Suppose  $x_0 \ge y_0(\max_{t:x_t \le A} S_t)^{\frac{1}{q-1}}(1 + \sup_{\text{polylog}} 1_{(d)})$ , then  $y_t \le O_{\mathbf{e}}(y_0)$  for all t such that  $x_t \le A$ . Moreover, if  $x_0 \ge y_0(\max_{t:x_t \le A} S_t)^{\frac{1}{q-1}}\log(d)$ , we would have  $|y_t - y_0| \lesssim \sup_{\text{polylog}(d)} |y_0|$ .

Moreover, we prove the following lemma for comparing the updates of different variables.

**Lemma G.3** (TPM of different degrees). Consider an increasing sequences  $x_t \ge 0$  defined by  $x_{t+1} = x_t + 1$  $\eta C_t x^q_t$ , for some integer  $q > q^0 \ge 3$  and  $q^0 \le q - 2$ , and  $C_t > 0$ , and further suppose given A = O(1), there exists  $t^0 \ge 0$ , $x_{t^0} \ge A$ . Then for every  $\delta > 0$  and every  $\eta \in (0,1)$ :

$$\sum_{t \ge 0, x_t \le A} \eta C_t x_t^{q'} \le (1+\delta)^{q'} \left( O(1) + \eta b A^q \right) \frac{1}{x_0^{q-q'-1}}$$

$$\sum_{t \ge 0, x_t \le A} \eta C_t x_t^{q'} \ge (1+\delta)^{-q'} \left( \delta (1+\delta)^{-1} \frac{1 - (1+\delta)^{-b(q-q'-1)}}{1 - (1+\delta)^{-(q-q'-1)}} - \eta b A^q \right) \frac{1}{x_0^{q-q'-1}}$$

where  $b = \Theta(\log(A/x_0)/\log(1+\delta))$ . When  $A = x_0 d^{\Theta(1)}$ ,  $\eta = o(A^{1}_{q\delta})$  and q = O(1), then

$$\sum_{t \ge 0, x_t \le A} \eta C_t x_t^{q'} = \Theta(\frac{1}{x_0^{q-q'-1}})$$

*Proof.* For every  $g \in 0,1,...$ , we define  $T_g := \min\{t : x_t \ge (1+\delta)^g x_0\}$ . and define  $b := \min\{g : (1+\delta)^g \ge A\}$ , we can write down the following two inequalities according to the update of  $x_t$ :

$$X \eta C_t[(1+\delta)gx]_q \le (1+\delta)xT_g - xT_g + \eta A_q \le \delta(1+\delta)gx0 + \eta A_q$$

$$t \in [T_g, T_{g+1}]$$

$$X \eta C_t[(1+\delta)g+1X_0]q \ge (1+\delta)xT_g - xT_g - \eta Aq \ge \delta(1+\delta)gx_0 - \eta Aq$$

$$t \in [T_g, T_{g+1}]$$

where  $g+1 \le b$ . Dividing both sides by  $[(1+\delta)^g x_0]^{q-q_0}$  in the first inequality and  $[(1+\delta)^{g+1} x_0]^{q-q_0}$  in the second, we have

$$\sum_{t \in [\mathcal{T}_g, \mathcal{T}_{g+1}]} \eta C_t [(1+\delta)^g x_0]^{q'} \le \frac{\delta}{(1+\delta)^{g(q-q'-1)}} \frac{1}{x_0^{q-q'-1}} + \frac{\eta A^q}{x_0^{q-q'-1}}$$

$$\sum_{t \in [\mathcal{T}_g, \mathcal{T}_{g+1}]} \eta C_t [(1+\delta)^{g+1} x_0]^{q'} \ge \frac{\delta (1+\delta)^{-1}}{(1+\delta)^{(g+1)(q-q'-1)}} \frac{1}{x_0^{q-q'-1}} - \frac{\eta A^q}{x_0^{q-q'-1}}$$

Therefore if we sum over 
$$g = 0,...,b$$
, then 
$$\sum_{t \geq 0, x_t \leq A} \eta C_t x_t^{q'} \leq \sum_{t \geq 0, x_t \leq A} \eta C_t [(1+\delta)^{g+1} x_0]^{q'}$$
 
$$= (1+\delta)^{q'} \sum_{t \geq 0, x_t \leq A} \eta C_t [(1+\delta)^g x_0]^{q'}$$
 
$$\leq (1+\delta)^{q'} \sum_{0 \leq g \leq b} \left( \frac{\delta}{(1+\delta)^{g(q-q'-1)}} \frac{1}{x_0^{q-q'-1}} + \frac{\eta A^q}{x_0^{q-q'-1}} \right)$$
 
$$= (1+\delta)^{q'} O\left( \frac{\delta}{(1+\delta)^{q-q'-1}} + \eta b A^q \right) \frac{1}{x_0^{q-q'-1}}$$
 
$$\leq (1+\delta)^{q'} O\left( \frac{1}{q-q'-1} + \eta b A^q \right) \frac{1}{x_0^{q-q'-1}}$$

For the lower bound, we also have

$$\sum_{t\geq 0, x_t \leq A} \eta C_t x_t^{q'} \geq (1+\delta)^{-q'} \sum_{t\geq 0, x_t \leq A} \eta C_t [(1+\delta)^{g+1} x_0]^{q'}$$

$$\geq (1+\delta)^{-q'} \sum_{0\leq g \leq b} \left( \frac{\delta(1+\delta)^{-1}}{(1+\delta)^{(g+1)(q-q'-1)}} - \eta A^q \right) \frac{1}{x_0^{q-q'-1}}$$

$$= (1+\delta)^{-q'} \left( \delta(1+\delta)^{-1} \frac{1 - (1+\delta)^{-b(q-q'-1)}}{1 - (1+\delta)^{-(q-q'-1)}} - \eta b A^q \right) \frac{1}{x_0^{q-q'-1}}$$

$$= (1+\delta)^{-q'} \left( \delta(1+\delta)^{-1} \frac{1 - (1+\delta)^{-b(q-q'-1)}}{1 - (1+\delta)^{-(q-q'-1)}} - \eta b A^q \right) \frac{1}{x_0^{q-q'-1}}$$

Inserting  $b = \Theta(\log(A/x_0)/\log(1 + \delta))$  proves the lower bound. For the last one we can choose  $\delta = \frac{1}{\sqrt{\log d}}$  to get:

$$b = \Theta(\text{polylog}^{}(d)), \quad \frac{\delta(1-(1+\delta)^{-b(q-q'-1)})}{1-(1+\delta)^{-(q-q'-1)}} = \Omega(1), \quad (1+\delta)^{-q'} = \Omega(1),$$

which proves the claim.

#### References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels? *Advances in Neural Information Processing Systems*, 32, 2019.
- [2] Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413*, 2020.
- [3] Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- [4] Zeyuan Allen-Zhu and Yuanzhi Li. Forward super-resolution: How can gans learn hierarchical generative models for real-world distributions. *arXiv preprint arXiv:2106.02619*, 2021.
- [5] Zeyuan Allen-Zhu and Yuanzhi Li. Feature purification: How adversarial training performs robust deep learning. In 62nd IEEE Annual Symposium on Foundations of Computer Science, FOCS 2021, Denver, CO, USA, February 7-10, 2022, pages 977–988. IEEE, 2021. URL https://doi.org/10.1109/FOCS52979.2021.00098.
- [6] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *NeurIPS 2019 : Thirty-third Conference on Neural Information Processing Systems*, pages 6158–6169, 2019.
- [7] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *ICML 2019 : Thirty-sixth International Conference on Machine Learning*, pages 242–252, 2019.
- [8] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the convergence rate of training recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 6676–6688, 2019.

- [9] Alexandr Andoni, Rina Panigrahy, Gregory Valiant, and Li Zhang. Learning polynomials with neural networks. In *International conference on machine learning*, pages 1908–1916. PMLR, 2014.
- [10] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [11] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pages 322–332. PMLR, 2019.
- [12] Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A theoretical analysis of contrastive unsupervised representation learning. *arXiv preprint arXiv:1902.09229*, 2019.
- [13] Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Dipendra Misra. Investigating the role of negatives in contrastive representation learning. *arXiv preprint arXiv:2106.09943*, 2021.
- [14] Yamini Bansal, Gal Kaplun, and Boaz Barak. For self-supervised learning, rationality implies generalization, provably. *arXiv preprint arXiv:2010.08508*, 2020.
- [15] Han Bao, Yoshihiro Nagano, and Kento Nozawa. Sharp learning bounds for contrastive unsupervised representation learning. *arXiv preprint arXiv:2110.02501*, 2021.
- [16] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.
- [17] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- [18] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [19] Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian inputs. In *Proceedings of the 34<sup>th</sup> International Conference on Machine LearningVolume* 70, pages 605–614. JMLR. org, 2017.
- [20] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020.
- [21] Mathilde Caron, Hugo Touvron, Ishan Misra, Herv'e J'egou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021.
- [22] Shuxiao Chen, Edgar Dobriban, and Jane H. Lee. A group-theoretic framework for data augmentation. *Journal of Machine Learning Research*, 21(245):1–71, 2020.
- [23] Sitan Chen, Jerry Li, Yuanzhi Li, and Anru R Zhang. Learning polynomial transformations. *arXiv* preprint arXiv:2204.04209, 2022.

- [24] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML 2020: 37th International Conference on Machine Learning*, volume 1, pages 1597–1607, 2020.
- [25] Ting Chen, Calvin Luo, and Lala Li. Intriguing properties of contrastive losses. *Advances in Neural Information Processing Systems*, 34, 2021.
- [26] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 15750–15758, 2021.
- [27] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [28] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9640–9649, 2021.
- [29] Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on Learning Theory*, pages 1305–1338. PMLR, 2020.
- [30] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186, 2019.
- [31] Simon S Du, Jason D Lee, and Yuandong Tian. When is a convolutional filter easy to learn? In *International Conference on Learning Representations*, 2018.
- [32] Simon S. Du, Jason D. Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *36<sup>th</sup> International Conference on Machine Learning, ICML 2019*, pages 1675–1685, 2019.
- [33] Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for selfsupervised representation learning. In *International Conference on Machine Learning*, pages 3015–3024. PMLR, 2021.
- [34] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021.
- [35] Antoine Gautier, Quynh N Nguyen, and Matthias Hein. Globally optimal training of generalized polynomial neural networks with nonlinear spectral methods. *Advances in Neural Information Processing Systems*, 29, 2016.
- [36] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of lazy training of two-layers neural network. In *Advances in Neural Information Processing Systems 32:*Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 9108–9118, 2019.
- [37] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural networks outperform kernel methods? In *Advances in Neural Information Processing Systems 33:*

- Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.
- [38] Sven Gowal, Po-Sen Huang, Aaron van den Oord, Timothy Mann, and Pushmeet Kohli. Self-supervised adversarial robustness for the low-label, high-data regime. In *International Conference on Learning Representations*, 2020.
- [39] Jean-Bastien Grill, Florian Strub, Florent Altch'e, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33:21271–21284, 2020.
- [40] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- [41] Zhaohan Daniel Guo, Bernardo Avila Pires, Bilal Piot, Jean-Bastien Grill, Florent Altch'e, R'emi Munos, and Mohammad Gheshlaghi Azar. Bootstrap latent-predictive representations for multitask reinforcement learning. In *International Conference on Machine Learning*, pages 3875–3886. PMLR, 2020.
- [42] Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for selfsupervised deep learning with spectral contrastive loss. *arXiv preprint arXiv:2106.04156*, 2021.
- [43] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9729–9738, 2020.
- [44] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Doll'ar, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.
- [45] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2019.
- [46] Tianyu Hua, Wenxiao Wang, Zihui Xue, Sucheng Ren, Yue Wang, and Hang Zhao. On feature decorrelation in self-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9598–9608, 2021.
- [47] Weiran Huang, Mingyang Yi, and Xuyang Zhao. Towards the generalization of contrastive self-supervised learning. *arXiv preprint arXiv:2111.00743*, 2021.
- [48] Samy Jelassi and Yuanzhi Li. Towards understanding how momentum improves generalization in deep learning, 2022. URL https://openreview.net/forum?id=lf0W6tcWmh-.
- [49] Samy Jelassi, Arthur Mensch, Gauthier Gidel, and Yuanzhi Li. Adam is no better than normalized SGD: Dissecting how adaptivity improves GAN performance, 2022. URL https://openreview.net/forum?id=D9SuLzhgK9.
- [50] Wenlong Ji, Zhun Deng, Ryumei Nakada, James Zou, and Linjun Zhang. The power of contrast for feature learning: A theoretical analysis. *arXiv preprint arXiv:2110.02473*, 2021.

- [51] Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pages 1772–1798. PMLR, 2019.
- [52] Li Jing, Pascal Vincent, Yann LeCun, and Yuandong Tian. Understanding dimensional collapse in contrastive self-supervised learning. *arXiv preprint arXiv:2110.09348*, 2021.
- [53] Stefani Karp, Ezra Winston, Yuanzhi Li, and Aarti Singh. Local signal adaptivity: Provable feature learning in neural networks beyond kernels. *Advances in Neural Information Processing Systems*, 34, 2021.
- [54] Joe Kileel, Matthew Trager, and Joan Bruna. On the expressive power of deep polynomial neural networks. *Advances in neural information processing systems*, 32, 2019.
- [55] Jason D Lee, Qi Lei, Nikunj Saunshi, and Jiacheng Zhuo. Predicting what you already know helps: Provable self-supervised learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- [56] Yuanzhi Li and Zehao Dou. Making method of moments great again?–how can gans learn distributions. *arXiv preprint arXiv:2003.04033*, 2020.
- [57] Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. In *Advances in neural information processing systems*, pages 597–607, 2017.
- [58] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in overparameterized matrix sensing and neural networks with quadratic activations. In *COLT 2018: 31st Annual Conference on Learning Theory*, pages 2–47, 2018.
- [59] Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial large learning rate in training neural networks. In *NeurIPS 2019 : Thirty-third Conference on Neural Information Processing Systems*, pages 11674–11685, 2019.
- [60] Yuanzhi Li, Tengyu Ma, and Hongyang R. Zhang. Learning over-parametrized two-layer relu neural networks beyond ntk. In *COLT*, pages 2613–2682, 2020.
- [61] Bingbin Liu, Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. Analyzing and improving the optimization landscape of noise-contrastive estimation. *arXiv preprint arXiv:2110.11271*, 2021.
- [62] Bingbin Liu, Daniel Hsu, Pradeep Ravikumar, and Andrej Risteski. Masked prediction tasks: a parameter identifiability view. *arXiv preprint arXiv:2202.09305*, 2022.
- [63] Zeping Luo, Cindy Weng, Shiyou Wu, Mo Zhou, and Rong Ge. One objective for all models self-supervised learning for topic models. *arXiv preprint arXiv:2203.03539*, 2022.
- [64] Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.
- [65] Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, and Kunio Kashino. Byol for audio: Self-supervised learning for general-purpose audio representation. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2021.
- [66] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

- [67] Ashwini Pokle, Jinjin Tian, Yuchen Li, and Andrej Risteski. Contrasting the landscape of contrastive and non-contrastive learning. *arXiv preprint arXiv:2203.15702*, 2022.
- [68] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [69] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [70] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *OpenAI blog*, 2022. URL https://cdn.openai.com/papers/dall-e-2.pdf.
- [71] Noam Razin and Nadav Cohen. Implicit regularization in deep learning may not be explainable by norms. *Advances in neural information processing systems*, 33:21174–21187, 2020.
- [72] Pierre H. Richemond, Jean-Bastien Grill, Florent Altch'e, Corentin Tallec, Florian Strub, Andrew Brock, Samuel Smith, Soham De, Razvan Pascanu, Bilal Piot, and Michal Valko. BYOL works even without batch statistics. *arXiv:2010.10241 [cs, stat]*, October 2020.
- [73] Joshua Robinson, Li Sun, Ke Yu, Kayhan Batmanghelich, Stefanie Jegelka, and Suvrit Sra. Can contrastive learning avoid shortcut solutions? *Advances in Neural Information Processing Systems*, 34, 2021.
- [74] Nikunj Saunshi, Sadhika Malladi, and Sanjeev Arora. A mathematical exploration of why language models help solve downstream tasks. *arXiv preprint arXiv:2010.03648*, 2020.
- [75] Nikunj Saunshi, Jordan Ash, Surbhi Goel, Dipendra Misra, Cyril Zhang, Sanjeev Arora, Sham Kakade, and Akshay Krishnamurthy. Understanding contrastive learning requires incorporating inductive biases. *arXiv* preprint *arXiv*:2202.14037, 2022.
- [76] Mahdi Soltanolkotabi. Learning relus via gradient descent. In *Advances in Neural Information Processing Systems*, volume 30, pages 2007–2017, 2017.
- [77] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- [78] Jiaye Teng, Weiran Huang, and Haowei He. Can pretext-based self-supervised learning be boosted by downstream data? a theoretical analysis. *arXiv* preprint arXiv:2103.03568, 2021.
- [79] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *European conference on computer vision*, pages 776–794. Springer, 2020.
- [80] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning. In *Advances in Neural Information Processing Systems*, volume 33, 2020.

- [81] Yuandong Tian. An analytical formula of population gradient for two-layered relu network and its applications in convergence and critical point analysis. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3404–3413. JMLR. org, 2017.
- [82] Yuandong Tian, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning dynamics without contrastive pairs. In *International Conference on Machine Learning*, pages 10268–10278. PMLR, 2021.
- [83] Christopher Tosh, Akshay Krishnamurthy, and Daniel Hsu. Contrastive estimation reveals topic posterior information to linear models. *arXiv preprint arXiv:2003.02234*, 2020.
- [84] Christopher Tosh, Akshay Krishnamurthy, and Daniel Hsu. Contrastive learning, multi-view redundancy, and linear models. In *Algorithmic Learning Theory*, pages 1179–1206. PMLR, 2021.
- [85] Yao-Hung Hubert Tsai, Yue Wu, Ruslan Salakhutdinov, and Louis-Philippe Morency. Selfsupervised learning from a multi-view perspective. *arXiv* preprint arXiv:2006.05576, 2020.
- [86] Julius Von Ku¨gelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Sch¨olkopf, Michel Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style. *Advances in Neural Information Processing Systems*, 34, 2021.
- [87] Xiang Wang, Xinlei Chen, Simon S Du, and Yuandong Tian. Towards demystifying representation learning with non-contrastive self-supervision. *arXiv preprint arXiv:2110.04947*, 2021.
- [88] Colin Wei, Sang Michael Xie, and Tengyu Ma. Why do pretrained language models help in downstream tasks? an analysis of head and prompt tuning. *Advances in Neural Information Processing Systems*, 34, 2021.
- [89] Zixin Wen and Yuanzhi Li. Toward understanding the feature learning process of selfsupervised contrastive learning. In *International Conference on Machine Learning*, pages 11112–11122. PMLR, 2021.
- [90] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5754–5764, 2019.
- [91] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and St'ephane Deny. Barlow twins: Selfsupervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. PMLR, 2021.
- [92] Chaoning Zhang, Kang Zhang, Chenshuang Zhang, Trung X Pham, Chang D Yoo, and In So Kweon. How does simsiam avoid collapse without negative samples? a unified understanding with self-supervised contrastive learning. *arXiv preprint arXiv:2203.16262*, 2022.
- [93] Kai Zhong, Zhao Song, Prateek Jain, Peter L Bartlett, and Inderjit S Dhillon. Recovery guarantees for one-hidden-layer neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 4140–4149. JMLR. org, 2017.
- [94] Difan Zou, Yuan Cao, Yuanzhi Li, and Quanquan Gu. Understanding the generalization of adam in learning neural networks with proper regularization. *arXiv preprint arXiv:2108.11371*, 2021.