

Why Fine-grained Labels in Pretraining Benefit Generalization?

Guan Zhe Hong
Purdue University

hong288@purdue.edu

Yin Cui*
NVIDIA

yinc@nvidia.com

Ariel Fuxman
Google Research

afuxman@google.com

Stanley H. Chan
Purdue University

stanchan@purdue.edu

Enming Luo
Google Research

enming@google.com

Reviewed on OpenReview: <https://openreview.net/forum?id=FojAV72owK>

Abstract

Recent studies show that pretraining a deep neural network with fine-grained labeled data, followed by fine-tuning on coarse-labeled data for downstream tasks, often yields better generalization than pretraining with coarse-labeled data. While there is ample empirical evidence supporting this, the theoretical justification remains an open problem. This paper addresses this gap by introducing a “hierarchical multi-view” structure to confine the input data distribution. Under this framework, we prove that: 1) coarse-grained pretraining only allows a neural network to learn the common features well, while 2) fine-grained pretraining helps the network learn the rare features in addition to the common ones, leading to improved accuracy on hard downstream test samples.

1 Introduction

We consider the theory of label granularity in deep learning. By label granularity, we mean a hierarchy of training labels specifying how detailed each label subclass needs to be (See Figure 1).

Having access to different granularity of labels offers us the freedom of training a classifier using a different level of precision. For example, instead of differentiating between `dogs` and `cats`, we can train a classifier to differentiate a `Poodle` dog and a `Persian` cat. The latter classification task is undoubtedly harder. However, recent studies found that if one uses fine-grained labels to *pre-train* a backbone, the pre-trained backbone will help the downstream neural networks generalize better (Chen et al., 2018). Vision transformers, for example, are well-known to require pretraining on large datasets with thousands of classes for effective downstream generalization (Dosovitskiy et al., 2021; He et al., 2016; Krizhevsky et al., 2012).

To convince readers who are less familiar with this particular training strategy, we conduct an experiment on ImageNet with details described in Appendix A.2 (we also include experiments on iNaturalist 2021 in Appendix A). Our experiment is limited in scale due to its high demand on the computing resources. Figure 2 shows an experiment of pre-training on ImageNet21k and fine-tuning the pre-trained network using ImageNet1k. The labels used in the ImageNet21k is based on WordNet Hierarchy. The downstream task

*Work done at Google Research.

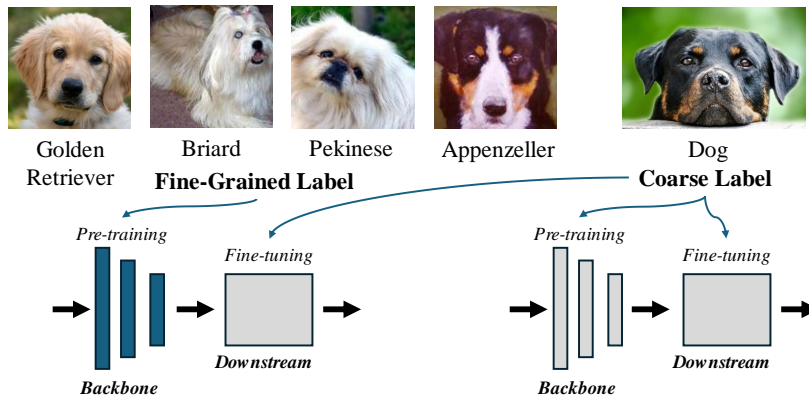


Figure 1: The goal of this paper is to provide a theoretical justification of why fine-grained labels in pre-training benefit generalization.

is ImageNet1k classification. The x -axis of this plot indicates the number of pre-training classes whereas the y -axis shows the validation accuracy for the ImageNet1k classification task. It is evident from the plot that as we increase the number of classes (hence a finer label granularity in pre-training), the downstream classification task’s performance is improved.

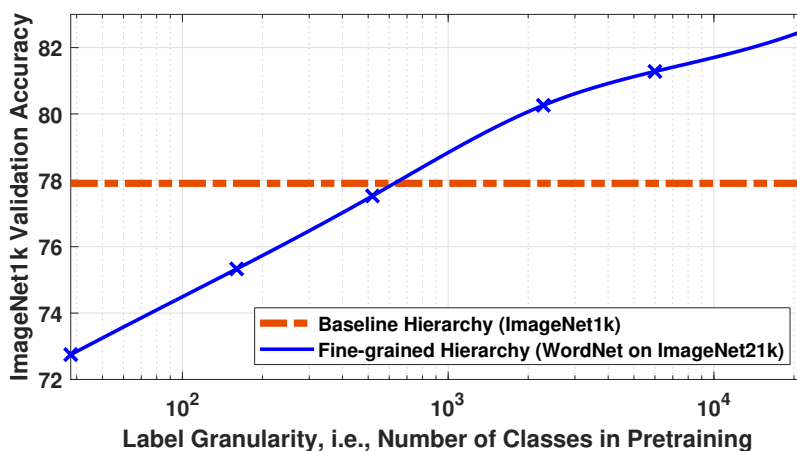


Figure 2: ImageNet21k ImageNet1k transfer using a ViT-B/16 model. [Blue]: pretrained on the WordNet hierarchy of ImageNet21k, *finetuned* on ImageNet1k. [Red]: baseline, trained and evaluated on ImageNet1k.

1.1 Goal of this paper

The above experimental finding may sound familiar to practitioners who frequently train large models. In fact, experimental evidence on this subject is abundant (Mahajan et al., 2018; Singh et al., 2022; Yan et al., 2020; Shnarch et al., 2022; Juan et al., 2020; Yang et al., 2021; Chen et al., 2018; Ridnik et al., 2021; Son et al., 2023; Ngiam et al., 2018; Cui et al., 2018; 2019a). However, the theoretical explanation remains an open problem. Our goal in this paper is to provide a *theoretical* justification. The core question we ask is:

— Theoretical Question —

Why does pretraining at a high label granularity benefit generalization?

Certainly, this grand challenge can be impossible to answer in full because of the uncontrollable complexity of the practical situations. To say something concrete, we focus on a tractable (sub-)problem under a controlled setting:

- *Simple* scheme: We pretrain a backbone on a classification task and then finetune it for a target problem;

- Assume *negligible distribution shift* between the input distributions of the source and target datasets;
- The label functions for both datasets *align* well in terms of the features which they consider discriminative;
- The labels are error-free.

1.2 Main results and theoretical contributions

Our main result is based on analyzing a two-layer convolutional neural network with ReLU activation. We assume that the data distribution satisfies a certain *hierarchical multi-view* condition (to be discussed in Section 4.1). The optimization algorithm is stochastic gradient descent. Such problem settings are consistent with published works on this subject (Allen-Zhu & Li, 2023b; 2022; Shen et al., 2022b; Jelassi & Li, 2022). Our conclusions are as follows.

— Theoretical results —

1. *Coarse-grained* pretraining *only* allows the neural network to learn the common features well. Therefore, when testing, the test error on easy samples is $o(1)$ (i.e., small) whereas the error on hard samples is $\Omega(1)$ (i.e., large).
2. *Fine-grained* pretraining helps the network learn the rare features *in addition* to the common ones, thus improving its test error on *hard* samples. In particular, the test error rate on both easy and hard test samples are $o(1)$ (i.e., small).

To our knowledge, a precise characterization of the test error presented in this paper has never been reported in the literature. The key enablers of our theoretical finding are the concepts of *hierarchical multi-view* and *representation-label correspondence*. We summarize these two concepts below:

1. *Hierarchical multi-view*. To understand the label granularity problem, we argue that it is necessary for coarse and fine-grained classes to be distinguished by their corresponding input features. This is consistent with the *multi-view* data property pioneered by Allen-Zhu & Li (2023b). We call this a *hierarchical multi-view* structure. The hierarchical multi-view structure on the data makes us different from many other deep learning theory works that assume simple or no structure in the input data (Kawaguchi, 2016; Allen-Zhu & Li, 2023a; Ba et al., 2022; 2023; Damian et al., 2022; Kumar et al., 2023; Ju et al., 2021).
2. *Representation-label correspondence*. Representation learning aims to recognize features in the input data. As will be shown later in the paper, under the hierarchical multi-view data assumption, label complexity (i.e., how complex the labels are) during training influences the representation complexity (i.e., how many and what types of features are learnt), which further influences the model’s generalization performance. Studying label granularity through understanding the neural network’s feature-learning process is a departure from the literature which focuses on *feature selection* (Jacot et al., 2018; Ju et al., 2021; 2022; Pezeshki et al., 2021; Arora et al., 2019), i.e., selecting a subset of pre-determined features.

2 Related Work

2.1 Our theoretical setting compared to the literature

The subject of label granularity is immensely related to how to make a deep neural network (DNN) generalize better. In the existing literature, this is mostly explained through the lens of implicit regularization and bias towards simpler solutions to prevent overfitting even when DNNs are highly overparameterized (Lyu et al., 2021; Kalimeris et al., 2019; Ji & Telgarsky, 2019; De Palma et al., 2019; Huh et al., 2017). An alternative approach is the concept of shortcut learning which argues that deep networks can learn overly simple solutions. As such, deep networks achieve high training and testing accuracy on in-distribution data but generalize poorly to challenging downstream tasks (Geirhos et al., 2020; Shah et al., 2020; Pezeshki et al., 2021).

By examining these papers, we believe that Shah et al. (2020); Pezeshki et al. (2021) are the closest to ours because they demonstrate that DNNs perform shortcut learning and respond weakly to features that have a weak presence in the training data. However, our work departs from Shah et al. (2020); Pezeshki et al. (2021) in several key ways.

1. We focus on how the pretraining label space affects classification generalization, while Shah et al. (2020); Pezeshki et al. (2021) primarily focus on demonstrating that simplicity bias can be harmful to generalization.
2. The core theoretical tool used by Pezeshki et al. (2021) is the neural tangent kernel (NTK) model, which is unsuitable for analyzing the label granularity problem because the feature extractor of an NTK model barely changes after pretraining.
3. The theoretical setting in Shah et al. (2020) is limited because they use the hinge loss while we use a more standard exponential-tailed cross-entropy loss.
4. Our data distribution assumptions are more realistic, as they capture feature hierarchies in natural images, which has direct impact on the downstream generalization power of the pretrained model.

2.2 Our analytic tool compared to literature

Our theoretical analysis is inspired by a recent line of work by Allen-Zhu & Li (2022; 2023b); Shen et al. (2022b). These papers analyze the feature learning dynamics of neural networks by tracking how the hidden neurons of shallow nonlinear neural networks evolve to solve dictionary-learning-like problems. We adopt a *multi-view* approach to the data distribution which was first proposed in Allen-Zhu & Li (2023b). However, the learning problems we analyze and the results we aim to show are fundamentally different. As such, we derive the gradient descent dynamics of the neural network from scratch.

2.3 Consistency with existing empirical results

We stress that our theoretical findings are consistent with the reported empirical results in the literature, especially those that aim to improve classification accuracy by manipulating the pre-training label space (Mahajan et al., 2018; Singh et al., 2022; Yan et al., 2020; Shnarch et al., 2022; Juan et al., 2020; Yang et al., 2021; Chen et al., 2018; Ridnik et al., 2021; Son et al., 2023; Ngiam et al., 2018; Cui et al., 2018; 2019a). For example, Mahajan et al. (2018); Singh et al. (2022) use hashtags from Instagram as pretraining labels, Yan et al. (2020); Shnarch et al. (2022) apply clustering on the data first and then treat the cluster IDs as pretraining labels, Juan et al. (2020) use the queries from image search results, Yang et al. (2021) apply image transformations such as rotation to augment the label space, and Chen et al. (2018); Ridnik et al. (2021) include fine-grained manual hierarchies in their pretraining processes. Our results corroborate the utility of pretraining on fine-grained label space.

On the empirical end, there is also work focusing on exploiting the hierarchical structures present in (human-generated) label space to improve classification accuracy (Yan et al., 2015; Zhu & Bain, 2017; Goyal & Ghosh, 2020; Sun et al., 2017; Zelikman et al., 2022; Silla & Freitas, 2011; Shkodrani et al., 2021; Bilal et al., 2017; Goo et al., 2016). For example, Yan et al. (2015) adapt the network architecture to learn super-classes at each hierarchical level, Zhu & Bain (2017) add hierarchical losses in the hierarchical classification task, Goyal & Ghosh (2020) propose a hierarchical curriculum loss for curriculum learning. Our results do not directly validate these practices because we are more interested in understanding the influence of label granularity on model generalization.

3 Notations and Intuitions

3.1 Notations and training schemes

For a DNN-based classifier, given input image \mathbf{X} , we can write its (pre-logit) output for class c as

$$\underbrace{F_c(\mathbf{X})}_{\text{pre-logit output for class } c} = \left\langle \underbrace{\mathbf{a}_c}_{\text{linear classifier}}, \underbrace{\mathbf{h}}_{\text{backbone network}} \left(\underbrace{\Theta}_{\text{network parameter}}; \mathbf{X} \right) \right\rangle, \quad (1)$$

where \mathbf{a}_c is the linear classifier for class c , $\mathbf{h}(\Theta; \cdot)$ is the network backbone with parameter Θ .

Referring to Figure 1, label granularity concerns about two datasets: X^{src} for the source (typically fine-grained) and X^{tgt} for the target (typically coarse-grained). The corresponding labels are Y^{src} and Y^{tgt} , respectively. A dataset can be represented as $D = (X, Y)$. For instance, the source training dataset is $D_{\text{train}}^{\text{src}} = (X_{\text{train}}^{\text{src}}, Y_{\text{train}}^{\text{src}})$.

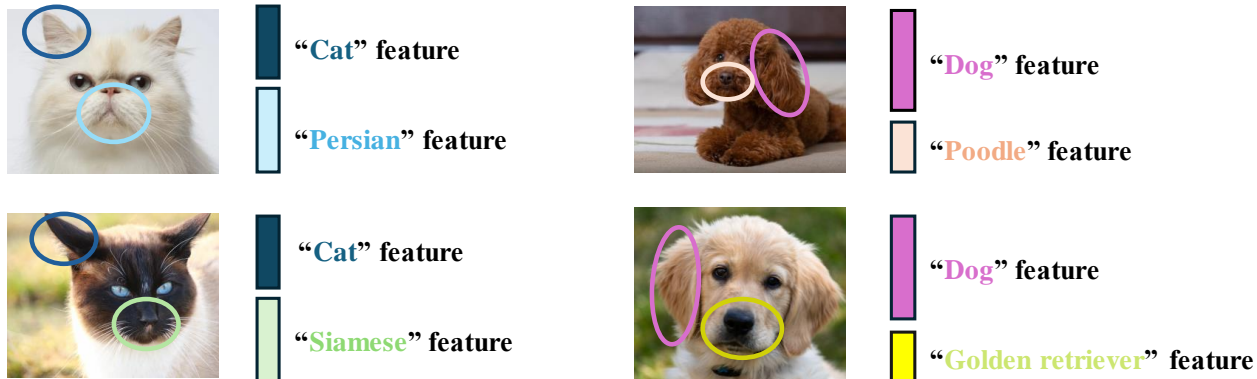


Figure 3: A simplified symbolic representation of the cat versus dog problem.

The relevant training and testing datasets are denoted as $D_{\text{train}}^{\text{src}}, D_{\text{train}}^{\text{tgt}}, D_{\text{test}}^{\text{tgt}}$. Finally, the granularity of a label set is denoted as $G(Y)$, which represents the total number of classes.

The two learning methodologies of interest are as follows.

1. Baseline: Train $F_c(\cdot)$ using $D_{\text{train}}^{\text{tgt}}$. Test $F_c(\cdot)$ using $D_{\text{test}}^{\text{tgt}}$.
2. Fine-to-coarse: Train $F_c(\cdot)$ using $D_{\text{train}}^{\text{src}}$. This gives us the pretrained feature extractor $h(\Theta_{\text{train}}^{\text{src}}; \cdot)$. Then finetune $F_c(\cdot)$ using $D_{\text{train}}^{\text{tgt}}$. Test the resulting $F_c(\cdot)$ using $D_{\text{test}}^{\text{tgt}}$.

3.2 Intuition: why higher granularity improves generalization

Consider the following toy example. There are two classes: **cat** and **dog**. Our goal is to build a binary classifier. Let’s discuss how the two training schemes would work, with an illustration shown in Figure 3.

1. **Baseline.** The baseline method tries to identify the common features that can distinguish most of the cats from dogs, for instance, the shape of the animal’s ear as shown in Figure 3. These features are often the most noticeable ones because they appear the most frequently. Of course, there are hard samples, e.g., a close-up shot of a cat’s fur. They pose limited influence during training because they are relatively rare in natural images.
2. **Fine-to-coarse.** With fine-grained labels, each subclass has its own unique visual features that are only dominant within that subclass. However, fine-grained features are not as common in the dataset, hence making them more difficult to be noticed. Therefore, if we only present the coarse labels in the pre-training stage, the learner is allowed to take shortcuts by learning only the common features to achieve low training loss. One strategy to force the learner to learn the rarer features is to explicitly label the fine-grained classes. This means that within each fine-grained class, the fine-grained features become as easy to notice as the common features. As a result, even if common features are weakly present or missing in a hard test sample, the network can still be reasonably robust to distracting irrelevant patterns due to its ability to recognize (some of) the finer-grained features.

4 Problem Formulation

Our first theoretical contribution is a new data model, the hierarchical multi-view model. This model consists of four definitions. Compared to existing theories studying feature learning of neural networks in the literature (Allen-Zhu & Li, 2023b; 2022; Shen et al., 2022b; Jelassi & Li, 2022), these four definitions are better formulated to the label granularity problem. For the sake of brevity, we present the core concepts of our data model here, and delay its full specification to Appendix B. Following data model specifications, we also discuss characteristics of the learner, a two-layer nonlinear convolutional neural network.

4.1 New data model: hierarchical multi-view

We consider the setting where an input sample $\mathbf{X} \in \mathbb{R}^{dP}$ consists of P patches $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_P$ with $\mathbf{x}_p \in \mathbb{R}^d$, where d is sufficiently large, and all our asymptotic statements are made with respect to d .

For analytic tractability, we consider two levels of label hierarchy. The root of this hierarchy has two superclasses $+1$ and -1 . The superclass $+1$ has k_+ subclasses. We denote these k_+ subclasses as $(+1, c_1), \dots, (+1, c_{k_+})$. We can do the same for the superclass -1 which has k_- subclasses. Each subclass has two types of features: the common features and the fine-grained features. The two types of features are sufficiently different in the sense they have zero correlation and equal magnitude. This leads to the following definition.

Definition 4.1 (Features). We define **features** as elements of a fixed orthonormal dictionary $V = \{\mathbf{v}_i\}_{i=1}^d \in \mathbb{R}^d$. The common and fine-grained features are

- Common feature: $\mathbf{v}_+ \in V$ and $\mathbf{v}_- \in V$
- Fine-grained feature of subclass c : $\mathbf{v}_{+,c} \in V$ and $\mathbf{v}_{-,c} \in V$

The usage of an orthonormal dictionary is again a choice of our model. We choose so because it is more tractable. With features defined, we can now specify patches in an input sample.

Definition 4.2 (Input patches). We define three types of patches for $y \in \{+, -\}$:

- (Common-feature patches) are defined as $\mathbf{x}_p = \alpha_p \mathbf{v}_y + \zeta_p$, where $\alpha_p \in [0, 1]$, and $\zeta_p \sim \mathcal{N}(\mathbf{0}, \sigma_\zeta^2 \mathbf{I}_d)$.
- (Subclass-feature patches) are defined as $\mathbf{x}_p = \alpha_p \mathbf{v}_{y,c} + \zeta_p$, where $\alpha_p \in [0, 1]$, and $\zeta_p \sim \mathcal{N}(\mathbf{0}, \sigma_\zeta^2 \mathbf{I}_d)$.
- (Noise patches) are defined as $\mathbf{x}_p = \zeta_p$.

Within an input sample $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_P)$, there are approximately s^* common-feature patches and s^* subclass-feature patches, the rest are all noise patches. Moreover, within a sample, the choice of y has to be consistent across the feature patches. Lastly, the positions of the features patches are random.

These definitions of the input patches are illustrated in Figure 4.

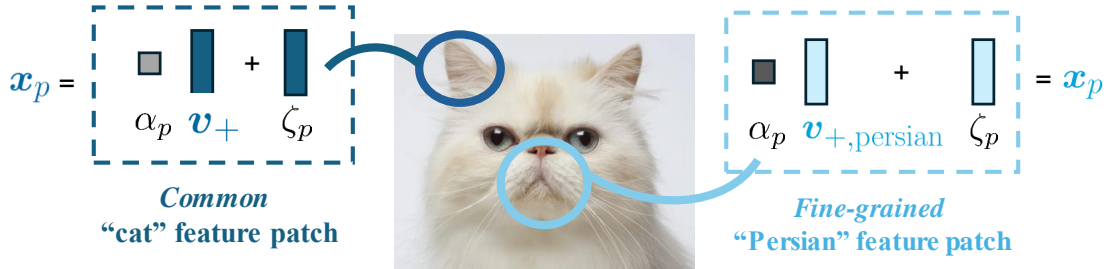


Figure 4: Illustration of features and patches.

Some comments: An **easy sample** is generated according to Definition 4.2. A **hard sample** is generated in the same way as easy samples, except the common-feature patches are replaced by noise patches, and we replace a small number of noise patches by “feature-noise” patches, which are of the form $\mathbf{x}_p = \alpha_p^\dagger \mathbf{v}_- + \zeta_p$, where $\alpha_p^\dagger \in o(1)$, and set one of the noise patches to $\zeta_p^* \sim \mathcal{N}(\mathbf{0}, \sigma_\zeta^2 \mathbf{I}_d)$ with $\sigma_\zeta \ll \sigma_\zeta^*$; these patches serve the role of “distracting patterns”.

Definition 4.3 (Source dataset’s label mapping). We say a sample \mathbf{X} belongs to the $+1$ superclass if any one of its common- or subclass-feature patches contains \mathbf{v}_+ or $\mathbf{v}_{+,c}$ for any $c \in [k_+]$. It belongs to the $(+, c)$ subclass if any one of its subclass-feature patches contains $\mathbf{v}_{+,c}$.

Definition 4.4 (Source training set). We assume the input samples of the source training set as $\mathcal{X}_{\text{train}}^{\text{src}}$ are generated as in Definition 4.2; the corresponding labels are generated following Definition 4.3. Overall, we denote the source training dataset $D_{\text{train}}^{\text{src}}$.

Relation to multi-view. Our data model is inspired by the multi-view concept first proposed in Allen-Zhu & Li (2023b), as we (1) use an orthonormal dictionary to define the features, (2) define an input consisting

of many disjoint high-dimensional patches, and (3) assume the existence of *multiple* discriminative features per class. The reason why the original multi-view property is insufficient for our problem is that it does not consider any label hierarchy nor its link to the input structure. We resolve this issue by following our intuition that classes at different hierarchy levels should be distinguished by their corresponding features: this naturally defines a feature hierarchy, with an exact correspondence with the label hierarchy.

Target dataset. To ensure that baseline and fine-grained training have no unfair advantage over each other, we post a set of new characterizations on the *target* dataset:

1. The input samples in the target dataset is generated according to Definition 4.2.
2. The true label function is identical across the source and target datasets.
3. Since we are studying the “fine-to-coarse” transfer direction, the target problem’s label space is the *root* of the hierarchy, meaning that any element of $\mathcal{Y}_{\text{train}}^{\text{tgt}}$ or $\mathcal{Y}_{\text{test}}^{\text{tgt}}$ must belong to the label space $\{+1, -1\}$.

Therefore, in our setting, only \mathcal{Y}^{src} and \mathcal{Y}^{tgt} can differ (in distribution) due to different choices in the label granularity level. In this idealized setting, we have essentially made baseline training and coarse-grained pretraining the *same* procedure. Therefore, an equally valid way to view our theory’s setting is to consider $D_{\text{train}}^{\text{tgt}}$ the same as $D_{\text{train}}^{\text{src}}$ except with coarse-grained labels. In other words, we pretrain the network on two versions of the source dataset $D_{\text{train}}^{\text{src,coarse}}$ and $D_{\text{train}}^{\text{src,fine}}$, and then compare the two models on $D_{\text{test}}^{\text{tgt}}$ (which has coarse-grained labels).

4.2 Characteristics about the learner

Our model about the learner is consistent with Allen-Zhu & Li (2023b; 2022); Shen et al. (2022b). The learner is a two-layer average-pooling convolutional ReLU network:

$$F_c(\mathbf{X}) = \sum_{r=1}^m a_{c,r} \sum_{p=1}^P \sigma(\mathbf{w}_{c,r} \cdot \mathbf{x}_p + b_{c,r}), \quad (2)$$

where m is a low-degree polynomial in d and denotes the width of the network, $\sigma(\cdot) = \max(0, \cdot)$ is the ReLU nonlinearity, and c denotes the class. We perform a *random initialization* of $\mathbf{w}_{c,r}^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I}_d)$ with $\sigma_0^2 = 1/\text{poly}(d)$; we set $b_{c,r}^{(0)} = -\Theta(\sigma_0 \sqrt{\ln(d)})$ and manually tune it, similar to Allen-Zhu & Li (2022). Cross-entropy is the training loss for both baseline and transfer training. To simplify analysis and to focus solely on the learning of the feature extractor, we freeze $a_{c,r} = 1$ during all baseline and transfer training phases, and we use the fine-grained model for binary classification as follows: $\hat{F}_+(\mathbf{X}) = \max_{c \in [k_+]} F_{+,c}(\mathbf{X})$, $\hat{F}_-(\mathbf{X}) = \max_{c \in [k_-]} F_{-,c}(\mathbf{X})$. See Appendix B.2 and the beginning of Appendix G for details of learner characteristics and training algorithm.

5 Theoretical results and proof strategy

Our second theoretical contribution lies in establishing a correspondence between the *complexity of the labels* and *complexity of the network’s representations*. Under the assumption of the hierarchical multi-view data structure, the following are true:

1. If trained with coarse-grained labels (i.e. overly simple labels), the network only learns the common features well, so its representations of the data is overly simple;
2. In contrast, training with fine-grained labels helps the network learn the fine-grained features well in addition to the common ones, so its representation of the data is more complex.

The difference in representation complexity leads to the difference in the network’s downstream test accuracy.

5.1 Main results

Theorem 5.1 (Coarse-label training: baseline). (*Summary*). *Let the number of subclasses be lower-bounded: $k_y \sim \text{polylog}(d)$. With high probability, with proper choice of step size, there exists a time $T^* \sim \text{poly}(d)$ such*

that for any $T \in [T^*, \text{poly}(d)]$, the training loss is upper bounded according to

$$L(F^{(T)}) = o(1) \quad (3)$$

Moreover, for an **easy** test sample $(\mathbf{X}_{\text{easy}}, y)$, the probability of making a classification mistake is small:

$$\mathbb{P} \left[F_y^{(T)}(\mathbf{X}_{\text{easy}}) \neq F_y^{(T)}(\mathbf{X}_{\text{easy}}) \right] = o(1), \quad \text{for } y' = y. \quad (4)$$

However, for all $t \in [0, \text{poly}(d)]$, given a **hard** test sample $(\mathbf{X}_{\text{hard}}, y)$, the probability of making a classification mistake is large:

$$\mathbb{P} \left[F_y^{(t)}(\mathbf{X}_{\text{hard}}) \neq F_y^{(t)}(\mathbf{X}_{\text{hard}}) \right] = \Omega(1), \quad \text{for } y' = y. \quad (5)$$

This theorem essentially says that, with a mild lower bound on the number of fine-grained classes, if we only train on the *easy* samples with *coarse* labels, it is virtually impossible for the network to learn the fine-grained features even if we give it as much practically reachable amount of time and training samples as possible. Consequently, the network would perform poorly on the hard downstream test samples: if the sample is missing the *common* features, then the network can be easily misled by the noise present in the sample. To see the full setup and statement of this theorem, please see Appendix B and E. Its proof spans Appendix C to E.

Theorem 5.2 (Fine-grained-label training). *(Summary). Assume the same setting as in Theorem 5.1, except let the labels be fine-grained and $k_y = d^{0.4}$ (number of subclasses not pathologically large; see Section 7 for its discussion). Within $\text{poly}(d)$ time, the probability of making a classification mistake is small:*

$$\mathbb{P} \left[\widehat{F}_y^{(T)}(\mathbf{X}) \neq \widehat{F}_y^{(T)}(\mathbf{X}) \right] = o(1) \quad \text{for } y' = y, \quad (6)$$

on the target binary problem on both **easy** and **hard** test samples.

The full version of this result is presented in Appendix G.4, and its proof in Appendix G. After fine-grained pretraining, the network’s feature extractor gains a strong response to the fine-grained features, therefore its accuracy on the downstream hard test samples increases significantly.

Remark. One concern about the above theorems is that the neural networks are trained only on easy samples. As noted in Sections 1 and 3.2, *easy* samples should make up the *majority* of the training and testing samples. Pretraining at higher label granularities only improves network performance on *rare* samples. Our theoretical result presents the feature-learning bias of a neural network in an exaggerated fashion. Therefore, it is natural to start with the case of no hard training samples. In reality, even if a small portion of hard training samples is present, finite-sized training datasets can have many flaws that can cause the network to overfit severely before learning the fine-grained features, especially since rarer features are learnt more slowly and corrupted by greater amount of noise. We leave these deeper considerations for future theoretical work.

5.2 Proof strategy: representation-label correspondence

The key idea of the proof is to establish a correspondence between the *complexity of the labels* and *complexity of the network’s representations*. We show that when trained on coarse-grained labels (i.e. overly simple labels), the network only learns the common features well, so its representations of the data is overly simple. In contrast, training with fine-grained labels helps the network learn the fine-grained features well in addition to the common ones, so its representations are more complex.

We first sketch the proof of **baseline training** which uses *coarse-grained* labels.

Feature detector neurons. We show that, at initialization, with high probability, for every feature $\mathbf{v} \in \mathcal{V}$, there exists a small group of “lucky” neurons, denoted $S_y^{*(0)}(\mathbf{v})$ (with y indicating the superclass), that only activate on \mathbf{v} -dominated feature patches. We prove that if \mathbf{v} is a feature of class y , then with high probability,

the lucky neurons will remain activated on \mathbf{v} -dominated patches throughout training, and dominate the feature extractor's response to the feature \mathbf{v} . In particular, given any \mathbf{v} -dominated patch $\mathbf{x}_p = \alpha_p \mathbf{v} + \mathbf{c}_p$,

$$\underbrace{\sum_{r=1}^m \sigma \left(\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{x}_p \rangle + b_{y,r}^{(t)} \right)}_{\text{network representation of } \mathbf{v}\text{-dominated patch } \mathbf{x}_p} = \underbrace{\sum_{r \in S_y^{*(0)}(\mathbf{v})} \sigma \left(\langle \mathbf{w}_{y,r}^{(t)}, \mathbf{x}_p \rangle + b_{y,r}^{(t)} \right)}_{\text{detector neurons' response to } \mathbf{v}\text{-dominated patch } \mathbf{x}_p}, \quad t \in [0, \text{poly}(d)]. \quad (7)$$

Therefore, we call neurons in $S_y^{*(0)}(\mathbf{v})$ the *detector neurons* of feature \mathbf{v} .

The significance of equation 7 is that, we may now argue about the *network's representation of the input data* solely based on the *behavior of the feature detector neurons*.

Impartial representation at initialization. At initialization, the feature extractor's response to common and fine-grained features are *very close*. The reason is that, $|S_y^{*(0)}(\mathbf{v})| \approx |S_{y'}^{*(0)}(\mathbf{v}')|$ for all superclasses y, y' and features \mathbf{v}, \mathbf{v}' , and they all have a similar magnitude of activation strength. Written explicitly, given any **common-feature** patch $\mathbf{x}_{\text{com}} = \alpha \mathbf{v}_y + \mathbf{c}$ and **subclass-feature** patch $\mathbf{x}_{\text{sub}} = \alpha' \mathbf{v}_{y,c} + \mathbf{c}'$ (from the training or testing distribution), with high probability,

$$\underbrace{\sum_{r \in S_y^{(0)}(\mathbf{v}_y)} \sigma \left(\langle \mathbf{w}_{y,r}^{(0)}, \alpha \mathbf{v}_y + \mathbf{c} \rangle + b_{y,r}^{(0)} \right)}_{\text{network representation of common-feature patch } \mathbf{x}_{\text{com}} \text{ at } t=0} \approx \underbrace{\sum_{r \in S_y^{(0)}(\mathbf{v}_{y,c})} \sigma \left(\langle \mathbf{w}_{y,r}^{(0)}, \alpha' \mathbf{v}_{y,c} + \mathbf{c}' \rangle + b_{y,r}^{(0)} \right)}_{\text{network representation of subclass-feature patch } \mathbf{x}_{\text{sub}} \text{ at } t=0} \quad (8)$$

So what happened during training which caused a strong *imbalance* of representation of the common and fine-grained features in the end? The answer below is the core of the proof.

Overly simple labels = overly simple representations. The imbalance of growth is a result of the subclass-feature patches occurring with less frequency in the training set than the common-feature patches. Recall that the number of subclasses is k_y : for any subclass (y, c) , subclass-feature patches dominated by $\mathbf{v}_{y,c}$ are about k_y times *rarer* than the common feature patches. This has a direct impact on the growth speed of the common and fine-grained detector neurons: for any neuron $r_{\text{com}} \in S_y^{*(0)}(\mathbf{v}_y)$ and any $r_{\text{fine}} \in S_y^{*(0)}(\mathbf{v}_{y,c})$, $\Delta \mathbf{w}_{y,r_{\text{com}}}^{(t)}, \mathbf{v}_y = \Theta(k_y) \times \Delta \mathbf{w}_{y,r_{\text{fine}}}^{(t)}, \mathbf{v}_{y,c}$.

With careful arguments on the influence of noise and bias on the activation values, we can show that, for t *sufficiently large*, the fine-grained detector neurons are about $\Theta(k_y)$ times *weaker* in strength:

$$\underbrace{\sum_{r \in S_y^{(0)}(\mathbf{v}_y)} \sigma \left(\langle \mathbf{w}_{y,r}^{(t)}, \alpha \mathbf{v}_y + \mathbf{c} \rangle + b_{y,r}^{(t)} \right)}_{\text{network representation of common-feature patch } \mathbf{x}_{\text{com}} \text{ at large } t} = \Theta(k_y) \times \underbrace{\sum_{r \in S_y^{(0)}(\mathbf{v}_{y,c})} \sigma \left(\langle \mathbf{w}_{y,r}^{(t)}, \alpha' \mathbf{v}_{y,c} + \mathbf{c}' \rangle + b_{y,r}^{(t)} \right)}_{\text{network representation of subclass-feature patch } \mathbf{x}_{\text{sub}} \text{ at large } t} \quad (9)$$

Furthermore, we prove that, due to the exponential tail of cross-entropy, by the end of training,

$$\sum_{r \in S_y^{(0)}(\mathbf{v}_y)} \sigma \left(\langle \Delta \mathbf{w}_{y,r}^{(t)}, \alpha \mathbf{v}_y + \mathbf{c} \rangle + b_{y,r}^{(t)} \right) = \Theta(\log(d)) \quad (10)$$

which causes the representation of subclass-feature patches to be *vanishing* in strength:

$$\sum_{r \in S_y^{(0)}(\mathbf{v}_{y,c})} \sigma \left(\langle \mathbf{w}_{y,r}^{(t)}, \alpha' \mathbf{v}_{y,c} + \mathbf{c}' \rangle + b_{y,r}^{(t)} \right) = O\left(\frac{\log(d)}{k_y}\right) < o(1), \quad t \in \text{poly}(d). \quad (11)$$

In other words, the neural network almost cannot detect subclass features by the end of baseline training. Therefore, even though it can classify the easy test samples correctly since it learned the common features

well, it simply cannot classify the hard ones, which requires the model to solely rely on subclass-feature patches for inference.

Fine-grained training alleviates this issue.

Complex labels = complex representations. The proof of fine-grained training proceeds in a very similar fashion as the case of coarse-grained training. The main difference lies in the gradient updates. During training, for any neuron $r_{\text{com}} \in S_{(y,c)}^{*(0)}(\mathbf{v}_y)$ and any $r_{\text{fine}} \in S_{(y,c)}^{*(0)}(\mathbf{v}_{y,c})$,

$$\langle \Delta \mathbf{w}_{(y,c),r_{\text{com}}}^{(t)}, \mathbf{v}_y \rangle \quad \langle \Delta \mathbf{w}_{(y,c),r_{\text{fine}}}^{(t)}, \mathbf{v}_{y,c} \rangle. \quad (12)$$

In other words, the common and fine-grained detector neurons for each subclass grow at similar speeds now, because the common- and subclass-feature patches occur with similar frequency in each subclass. Again with careful analysis of how the noise and bias influence the activation values, we arrive at

$$\underbrace{\sum_{r \in S_{(y,c)}^{(0)}(\mathbf{v}_y)} \sigma \left(\langle \mathbf{w}_{(y,c),r}^{(t)}, \alpha \mathbf{v}_y + \cdot \rangle + b_{(y,c),r}^{(t)} \right)}_{\substack{\text{network representation of} \\ \text{common-feature patch } x_{\text{com}}, \text{ end of training}}} \quad \underbrace{\sum_{r \in S_{(y,c)}^{(0)}(\mathbf{v}_{y,c})} \sigma \left(\langle \mathbf{w}_{(y,c),r}^{(t)}, \alpha' \mathbf{v}_{y,c} + \cdot \rangle + b_{(y,c),r}^{(t)} \right)}_{\substack{\text{network representation of} \\ \text{subclass-feature patch } x_{\text{sub}}, \text{ end of training}}} \quad (13)$$

$\Omega(1)$ $\Omega(1)$

Therefore, both the common and fine-grained features are learnt well. It follows that the model can correctly utilize the common- and subclass-feature patches in the input, so it can classify easy and hard test samples with high accuracy.

6 Empirical Results

Building on our theoretical analysis in an idealized setting, this section discusses conditions on the source and target label functions that we observed to be important for fine-grained pretraining to work *in practice*, while remaining in the controlled setting described in Section 1 for the sake of tractability. We present the core experimental results obtained on ImageNet21k and iNaturalist 2021 in the main text, and leave the experimental details and ablation studies to Appendix A.

6.1 ImageNet21k ImageNet1k transfer experiment

This subsection provides more details about the experiment shown in Figure 2. Specifically, we show that the common practice of pretraining on ImageNet21k using leaf labels is indeed better than pretraining at lower granularities in the manual hierarchy.

Hierarchy definition. The label hierarchy in ImageNet21k is based on WordNet Miller (1995); Deng et al. (2009). To define fine-grained labels, we first define the leaf labels of the dataset as Hierarchy level 0. For each image, we trace the path from the leaf label to the root using the WordNet hierarchy. We then set the k -th synset (or the root synset, if it is higher in the hierarchy) as the level- k label of this image. This procedure also applies to the multi-label samples. This is how we generate the hierarchies shown in Figure 2.

Network choice and training. For this dataset, we use the more recent Vision Transformer ViT-B/16 Dosovitskiy et al. (2021). Our pretraining pipeline is almost identical to the one in Dosovitskiy et al. (2021). For fine-tuning, we experimented with several strategies and report only the best results in the main text; the finer details are discussed in Appendix A.1.2 and A.2. To ensure a fair comparison, we also used these strategies to find the best baseline result by using $D_{\text{train}}^{\text{tgt}}$ for pretraining.

6.2 Transfer experiment on iNaturalist 2021

We conduct a systematic study of the transfer method *within* the label hierarchies of iNaturalist 2021 (Horn & macaodha, 2021). This dataset is well-suited for our analysis because it has a manually defined label

hierarchy that is based on the biological traits of the creatures in the images. Additionally, the large sample size of this dataset reduces the likelihood of sample-starved pretraining on reasonably fine-grained hierarchy levels.

Our experiments on this dataset again demonstrate that, as long as the finer-grained labels contain little noise, are well-aligned with the target label space, and sample count per subclass is not too limited, then we observe improvement in the model’s generalization performance. However, we also show negative results outside of the aforementioned “nice regime”: when sample count per sub-class is limited, or the fine-grained labels are *noisy*, or potentially *misaligned* with the target label space’s, finer-grained labels do not necessarily improve generalization significantly.

Relevant datasets. We perform transfer experiments within iNaturalist2021. More specifically, we set $\mathcal{X}_{\text{train}}^{\text{src}}$ and $\mathcal{X}_{\text{train}}^{\text{tgt}}$ both equal to the training split of the input samples in iNaturalist2021, and set $\mathcal{X}_{\text{train}}^{\text{tgt}}$ to the testing split of the input samples in iNaturalist2021. To focus on the “fine-to-coarse” transfer setting, the *target problem* is to classify the root level of the manual hierarchy, which contains 11 superclasses. To generate a greater gap between the performance of different hierarchies and to shorten training time, we use the mini version of the training set in all our experiments.

Alternative hierarchies generation. To better understand the transfer method’s operating regime, we experiment with different ways of generating the fine-grained labels for pretraining: we perform kMeans clustering on the ViT-L/14-based CLIP embedding Radford et al. (2021); Dehghani et al. (2022) of every sample in the training set and use the cluster IDs as pretraining class labels. We carry out this experiment in two ways. The green curve in Figure 5 comes from performing kMeans clustering on the embedding of each superclass *separately*, while the purple one’s cluster IDs are from performing kMeans on the *whole dataset*. The former way preserves the implicit hierarchy of the superclasses in the cluster IDs: samples from superclass k cannot possibly share a cluster ID with samples belonging to superclass $k' = k$. Therefore, its label function is forced to align better with that of the 11 superclasses than the purple curve’s. We also assign random class IDs to samples.

Network choice and training. We experiment with ResNet 34 and 50 on this dataset. For pretraining on $D_{\text{train}}^{\text{src}}$ with fine-grained labels, we adopt a standard 90-epoch large-batch-size training procedure commonly used on ImageNet He et al. (2016); Goyal et al. (2017). Then we finetune the network for 90 epochs and test it on the 11-superclass $D_{\text{train}}^{\text{tgt}}$ and $D_{\text{test}}^{\text{tgt}}$, respectively, using the pretrained backbone $h(\Theta_{\text{src}}; \cdot)$. To ensure a fair comparison, we trained the baseline model using exactly the same training pipeline, except that the pretraining stage uses $D_{\text{train}}^{\text{tgt}}$. We observed that this “retraining” baseline consistently outperformed the naive one-pass 90-epoch training baseline on this dataset. Due to space limitations, we leave the results of ResNet50 to the appendix.

Interpretation of results. Figure 5 shows the validation errors of the resulting models on the 11-superclass problem. We make the following observations.

Reasonably fine-grained labels benefit generalization, but there is a catch. We can observe that in the blue curve of Figure 5 that, as long as the number of subclasses is less than 10^3 , we see obvious decline in the validation error on the target labels. In other words, reasonably fine-grained pretraining is indeed beneficial in this setting. We should note, however, the overall curve exhibits a U shape: overly fine-grained labels are not beneficial to downstream generalization. This is intuitive. If the pretraining granularity is too close to the target one, we should not expect improvement. On the other extreme, if we assign a unique label to *every* sample in the training data, it is highly likely that the only *differences* a model can find between each class would be frivolous details of the images, which would not be considered discriminative by the label function of the target coarse-label problem. In this case, the pretraining stage is almost meaningless and can be misleading, as evidenced by the very high label-per-sample error (red star in Figure 5).

High granularity can be helpful, but label-assignment consistency is critical. Random class ID pretraining (orange curve) performs the worst of all the alternatives. The label function of this type does not generate a *meaningful hierarchy* because it has no consistency in the features it considers discriminative when decomposing the superclasses. This is in stark contrast to the manual hierarchies, which decompose the superclasses based on the finer biological traits (mostly visual in nature) of the creatures in the image.

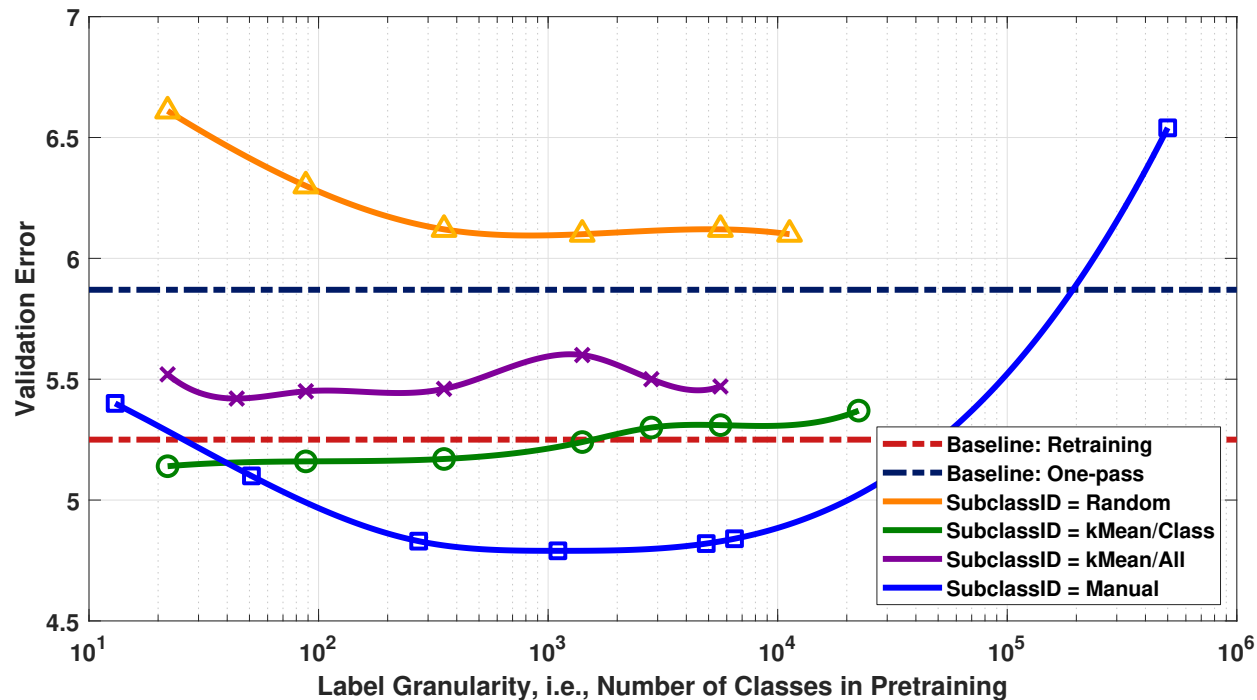


Figure 5: **In-dataset transfer**. ResNet34 validation error (with standard deviation) of finetuning on 11 superclasses of iNaturalist 2021, pretrained on various label hierarchies. The manual hierarchy outperforms the baseline and every other hierarchy, and exhibits a U-shaped curve.

Alignment between fine-grained and target label spaces are important. For fine-grained pretraining to be effective, the features that the pretraining label function considers discriminative must *align* well with those valued by the label function of the 11-superclass hierarchy. To see this point, observe that for models trained on cluster IDs obtained by performing kMeans on the CLIP embedding samples in each superclass *separately* (green curve in Figure 5), their validation errors are much lower than those trained on cluster IDs obtained by performing kMeans on the whole dataset (purple curve in Figure 5). As expected, the manually defined fine-grained label functions align best with that of the 11 superclasses, and the results corroborate this view.

7 Discussion

Q: Are there other reasons why fine-grained labels benefit neural network generalization?

A: Yes, it is possible, e.g., the optimization landscape induced by finer-grained labels could contain less saddle points, making it friendlier to SGD. We did not analyze this because our focus is primarily on the generalization instead of optimization aspect of the problem.

Q: Does a higher label granularity always imply better generalization?

A: No. There is an operating regime. Training a model with *pathologically* high label granularity is harmful. For example, if we assign a unique class to every sample in the dataset, the model will be forced to rely on the frivolous differences between each sample. We verify this intuition in Figure 5 in Appendix A.1.1 on the iNaturalist 2021 dataset. These extreme scenarios do not arise in common practice, so we do not focus on them in this paper.

Q: The theoretical setting appears restrictive.

A: Our theoretical setting is consistent with Cao et al. (2022); Allen-Zhu & Li (2023b; 2022); Shen et al. (2022b); Jelassi & Li (2022). With a limited number of available analytic tools in the literature, we believe these settings are necessary to keep things tractable.

8 Conclusion

In this paper, we formally studied the influence of pretraining label granularity on the generalization of DNNs, and performed large-scale experiments to complement our theoretical results. Under the new data model, hierarchical multi-view, we theoretically showed that higher label complexity leads to higher representation complexity, through which we explained why pretraining with fine-grained labels is beneficial to generalization. We complement our theory with experiments on ImageNet and iNaturalist, demonstrating that in the controlled setting of this paper, pretraining on reasonably fine-grained labels indeed benefits generalization.

Broader Impact Statement

This paper presents work whose goal is to advance the theory of deep learning. There are potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Zeyuan Allen-Zhu and Yuanzhi Li. Feature purification: How adversarial training performs robust deep learning. In *FOCS*, 2022.
- Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep (hierarchical) learning. In *COLT*, 2023a.
- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. In *ICLR*, 2023b.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *NeurIPS*, 2019.
- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional asymptotics of feature learning: How one gradient step improves the representation. In *NeurIPS*, 2022.
- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning in the presence of low-dimensional structure: A spiked random matrix perspective. In *NeurIPS*, 2023.
- Alsallakh Bilal, Amin Jourabloo, Mao Ye, Xiaoming Liu, and Liu Ren. Do convolutional neural networks learn class hierarchy? *IEEE transactions on visualization and computer graphics*, 2017.
- Yuan Cao, Zixiang Chen, Misha Belkin, and Quanquan Gu. Benign overfitting in two-layer convolutional neural networks. In *NeurIPS*, 2022.
- Zhuo Chen, Ruizhou Ding, Ting-Wu Chin, and Diana Marculescu. Understanding the impact of label granularity on cnn-based image classification. In *ICDMW*, 2018.
- Yin Cui, Yang Song, Chen Sun, Andrew Howard, and Serge Belongie. Large scale fine-grained categorization and domain-specific transfer learning. In *CVPR*, 2018.
- Yin Cui, Zeqi Gu, Dhruv Mahajan, Laurens Van Der Maaten, Serge Belongie, and Ser-Nam Lim. Measuring dataset granularity. *arXiv preprint arXiv:1912.10154*, 2019a.
- Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, 2019b.
- Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent. In *COLT*, 2022.
- Giacomo De Palma, Bobak Kiani, and Seth Lloyd. Random deep neural networks are biased towards simple functions. In *NeurIPS*, 2019.
- Mostafa Dehghani, Alexey Gritsenko, Anurag Arnab, Matthias Minderer, and Yi Tay. Scenic: A jax library for computer vision research and beyond. In *CVPR*, 2022.

- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A. Wichmann. Shortcut learning in deep neural networks. In *Nature Machine Intelligence*, 2020.
- Wonjoon Goo, Juyong Kim, Gunhee Kim, and Sung Ju Hwang. Taxonomy-regularized semantic deep convolutional neural networks. In *ECCV*, 2016.
- Palash Goyal and Shalini Ghosh. Hierarchical class-based curriculum loss. *arXiv preprint arXiv:2006.03629*, 2020.
- Priya Goyal, Piotr Dollar, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv:1706.02677*, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- Grant Van Horn and macaodha. inat challenge 2021 - fgvc8, 2021. URL <https://kaggle.com/competitions/inaturalist-2021>.
- Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. The low-rank simplicity bias in deep networks. *arXiv:2103.10427*, 2017.
- Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *NeurIPS*, 2018.
- Samy Jelassi and Yuanzhi Li. Towards understanding how momentum improves generalization in deep learning. In *ICML*, 2022.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In *ICLR*, 2019.
- Ralph P. Boas Jr. and Jr. John W. Wrench. Partial sums of the harmonic series. *The American Mathematical Monthly*, 1971.
- Peizhong Ju, Xiaojun Lin, and Ness Shroff. On the generalization power of overfitted two-layer neural tangent kernel models. In *ICML*, 2021.
- Peizhong Ju, Xiaojun Lin, and Ness Shroff. On the generalization power of the overfitted three-layer neural tangent kernel model. In *NeurIPS*, 2022.
- Da-Cheng Juan, Chun-Ta Lu, Zhen Li, Futang Peng, Aleksei Timofeev, Yi-Ting Chen, Yaxi Gao, Tom Duerig, Andrew Tomkins, and Sujith Ravi. Ultra fine-grained image semantic embedding. In *WSDM*, 2020.
- Dimitris Kalimeris, Gal Kaplun, Preetum Nakkiran, Benjamin Edelman, Tristan Yang, Boaz Barak, and Haofeng Zhang. Sgd on neural networks learns functions of increasing complexity. In *NeurIPS*, 2019.
- Stefani Karp, Ezra Winston, Yuanzhi Li, and Aarti Singh. Local signal adaptivity: Provable feature learning in neural networks beyond kernels. In *NeurIPS*, 2021.
- Kenji Kawaguchi. Deep learning without poor local minima. In *NeurIPS*, 2016.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012.
- Tanishq Kumar, Blake Bordelon, Samuel J. Gershman, and Cengiz Pehlevan. Grokking as the transition from lazy to rich training dynamics. *arXiv:2310.06110*, 2023.

- Béatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5), 2000.
- Kuang-Huei Lee, Anurag Arnab, Sergio Guadarrama, John Canny, and Ian Fischer. Compressive visual representations. In *NeurIPS*, 2021.
- Kaifeng Lyu, Zhiyuan Li, Runzhe Wang, and Sanjeev Arora. Gradient descent on two-layer nets: Margin maximization and simplicity bias. In *NeurIPS*, 2021.
- Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018.
- George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 1995.
- Jiquan Ngiam, Daiyi Peng, Vijay Vasudevan, Simon Kornblith, Quoc V Le, and Ruoming Pang. Domain adaptive transfer learning with specialist models. *arXiv preprint arXiv:1811.07056*, 2018.
- Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. In *NeurIPS*, 2021.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik. Imagenet-21k pretraining for the masses. In *NeurIPS Track on Datasets and Benchmarks*, 2021.
- Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. In *NeurIPS*, 2020.
- Ruoqi Shen, Sebastien Bubeck, and Suriya Gunasekar. Data augmentation as feature manipulation. In *ICML*, 2022a.
- Ruoqi Shen, Sebastien Bubeck, and Suriya Gunasekar. Data augmentation as feature manipulation. In *ICML*, 2022b.
- Sindi Shkodrani, Yu Wang, Marco Manfredi, and Nóra Baka. United we learn better: Harvesting learning improvements from class hierarchies across tasks. *arXiv preprint arXiv:2107.13627*, 2021.
- Eyal Shnarch, Ariel Gera, Alon Halfon, Lena Dankin, Leshem Choshen, Ranit Aharonov, and Noam Slonim. Cluster & tune: Boost cold start performance in text classification. *arXiv preprint arXiv:2203.10581*, 2022.
- Carlos N Silla and Alex A Freitas. A survey of hierarchical classification across different application domains. *Data Mining and Knowledge Discovery*, 2011.
- Mannat Singh, Laura Gustafson, Aaron Adcock, Vinicius de Freitas Reis, Bugra Gedik, Raj Prateek Kosaraju, Dhruv Mahajan, Ross Girshick, Piotr Dollár, and Laurens Van Der Maaten. Revisiting weakly supervised pre-training of visual perception models. In *CVPR*, 2022.
- Donghyun Son, Byounggyu Lew, Kwanghee Choi, Yongsu Baek, Seungwoo Choi, Beomjun Shin, Sungjoo Ha, and Buru Chang. Reliable decision from multiple subtasks through threshold optimization: Content moderation in the wild. In *WSDM*, 2023.
- Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *ICCV*, 2017.
- Xueting Yan, Ishan Misra, Abhinav Gupta, Deepti Ghadiyaram, and Dhruv Mahajan. Clusterfit: Improving generalization of visual representations. In *CVPR*, 2020.

Zhicheng Yan, Hao Zhang, Robinson Piramuthu, Vignesh Jagadeesh, Dennis DeCoste, Wei Di, and Yizhou Yu. Hd-cnn: hierarchical deep convolutional neural networks for large scale visual recognition. In *ICCV*, 2015.

Chuanguang Yang, Zhulin An, Linhang Cai, and Yongjun Xu. Hierarchical self-supervised augmented knowledge distillation. *arXiv preprint arXiv:2107.13715*, 2021.

Eric Zelikman, Jesse Mu, Noah D Goodman, and Yuhuai Tony Wu. Star: Self-taught reasoner bootstrapping reasoning with reasoning. In *NeurIPS*, 2022.

Xinqi Zhu and Michael Bain. B-cnn: branch convolutional neural network for hierarchical classification. *arXiv preprint arXiv:1709.09890*, 2017.

Appendix

A Additional Experimental Results

In this section, we present the full details of our experiments and relevant ablation studies. All of our experiments were performed using tools in the Scenic library Dehghani et al. (2022).

A.1 In-dataset transfer results

To clarify, in this transfer setting, we are essentially transferring *within* a dataset. More specifically, we set $X^{\text{src}} = X^{\text{tgt}}$ and only the label spaces Y^{src} and Y^{tgt} may differ (in distribution). The baseline in this setting is clear: train on $D_{\text{train}}^{\text{tgt}}$ and test on $D_{\text{test}}^{\text{tgt}}$. In contrast, after pretraining the backbone network $h(\Theta; \cdot)$ on Y^{src} , we finetune or linear probe it on $D_{\text{train}}^{\text{tgt}}$ using the backbone and then test on $D_{\text{test}}^{\text{tgt}}$.

A.1.1 iNaturalist 2021

iNaturalist 2021 is well-suited for our analysis because it has a high-quality, manually defined label hierarchy that is based on the biological traits of the creatures in the images. Additionally, the large sample size of this dataset reduces the likelihood of sample-starved pretraining on reasonably fine-grained hierarchy levels. We use the mini training dataset with size 500,000 instead of the full training dataset to show a greater gap between the results of different hierarchies and speed up training.

We use the architectures ResNet 34 and 50 He et al. (2016).

Training details. Our pretraining pipeline on iNaturalist is essentially the same as the standard large-batch-size ImageNet-type training for ResNets He et al. (2016); Goyal et al. (2017). The following pipeline applies to model pretraining on any hierarchy.

- Optimization: SGD with 0.9 momentum coefficient, 0.00005 weight decay, 4096 batch size, 90 epochs total training length. We perform 7 epochs of linear warmup in the beginning of training until the learning rate reaches $0.1 \times 4096/256 = 1.6$, and then apply the cosine annealing schedule. Each training instance is run on 16 TPU v4 chips, taking around 2 hours per run.
- Data augmentation: subtracting mean and dividing by standard deviation, image (original or its horizontal flip) resized such that its shorter side is 256 pixels, then a 224×224 random crop is taken.

For finetuning, we keep everything in the pipeline the same except setting the batch size to $4096/4 = 1024$ and base learning rate $1.6/4 = 0.4$. We found that finetuning at higher batch size and learning rate resulted in training instabilities and severely affected the final finetuned model’s validation accuracy, while finetuning at lower batch size and learning rate than the chosen one resulted in lower validation accuracy at the end even though their training dynamics was stabler.

For the baseline accuracy, as mentioned in the main text, to ensure fairness of comparison, in addition to only training the network on the target 11-superclass problem for 90 epochs (using the same pretraining pipeline), we also perform “retraining”: follow the exact training process of the models trained on the various hierarchies, but use $D_{\text{train}}^{\text{tgt}}$ as the training dataset in both the pretraining and finetuning stage. We observed consistent increase in the final validation accuracy of the model, so we report this as the baseline accuracy. Without retraining (so naive one-pass 90-epoch training on 11 superclasses), the average accuracy with standard deviation is 94.13, 0.025.

Clustering. To obtain the cluster-ID-based labels, we perform the following procedure.

1. For every sample X_n in the mini training dataset of iNaturalist 2021, obtain its ViT-L/14 CLIP embedding E_n .
2. Per-superclass kMeans clustering. Let C be the predefined number of clusters per class.

Table 1: **In-dataset transfer, iNaturalist 2021.** ResNet34 average finetuning validation error and standard deviation on 11 superclasses in iNaturalist 2021, pretrained on various label hierarchies with different label granularity. Baseline (11-superclass) and best performance are highlighted.

| Manual Hierarchy | $G(Y^{src})$ | 11 | 13 | 51 | 273 | 1103 | 4884 | 6485 |
|----------------------------|------------------|-------------------|------------|------------|------------|-------------------|------------|------------|
| Random class ID | Validation error | 5.25±0.051 | 5.40±0.075 | 5.10±0.038 | 4.83±0.041 | 4.79±0.045 | 4.82±0.056 | 4.84±0.033 |
| | $G(Y^{src})$ | 22 | 88 | 352 | 1,408 | 5,632 | 11,264 | 500,000 |
| | Validation error | 6.61±0.215 | 6.30±0.070 | 6.12±0.77 | 6.10±0.053 | 6.12±0.042 | 6.10±0.057 | 6.54±0.758 |
| CLIP+kMeans per superclass | $G(Y^{src})$ | 22 | 88 | 352 | 1408 | 2816 | 5632 | 22528 |
| | Validation error | 5.14±0.049 | 5.16±0.033 | 5.17±0.027 | 5.24±0.029 | 5.30±0.029 | 5.31±0.077 | 5.37±0.032 |
| C+k per supclass | $G(Y^{src})$ | 88 | 218 | 320 | 608 | 1040 | 1984 | |
| Class rebalanced | Validation error | 5.18±0.054 | 5.17±0.038 | 5.23±0.052 | 5.28±0.045 | 5.26±0.035 | 5.21±0.040 | |
| CLIP+kMeans whole dataset | $G(Y^{src})$ | 22 | 44 | 88 | 352 | 1408 | 2816 | 5632 |
| | Validation error | 5.52±0.015 | 5.42±0.047 | 5.45±0.049 | 5.46±0.019 | 5.60±0.029 | 5.50±0.029 | 5.47±0.029 |

- (a) For every superclass k , for the set of embedding $\{(\mathbf{E}_n, y_n = k)\}$ belonging to that superclass, perform kMeans clustering with cluster size set to C .
- (b) Given a sample with superclass ID $k \in \{1, 2, \dots, 11\}$ and cluster ID $c \in \{1, 2, \dots, C\}$, define its fine-grained ID as $C \times k + c$.

3. Whole-dataset kMeans clustering. Let C be the predefined number of clusters on the whole dataset.

- (a) Perform kMeans on the embedding of all the samples in the dataset, with the number of clusters set to C . Set the fine-grained class ID of a sample to its cluster ID.

Some might have the concern that having the same number of kMeans clusters per superclass could cause certain classes to have too few samples, which could be a reason for why the cluster ID hierarchies perform worse than the manual hierarchies. Indeed, the number of samples per superclass on iNaturalist is different, so in addition to the above “uniform-number-of-cluster-per-superclass” hierarchy, we add an extra label hierarchy by performing the following procedure to balance the sample size of each cluster:

1. Perform kMeans for each superclass with number of clusters set to 2, 8, 32, 64, 128, 256, 512, 1024 and save the corresponding image-ID-to-cluster-ID dictionaries (so we are basically reusing the clustering results of the CLIP+kMeans per superclass experiment)
2. For each superclass, find the image-ID-to-cluster-ID dictionary with the highest granularity while still keeping the minimum number of samples for each cluster $>$ predefined threshold (e.g. 1000 samples per subclass)
3. Now we have nonuniform granularity for each superclass while ensuring that the sample count per cluster is above some predefined threshold.

This simple procedure somewhat improves the balance of sample count per cluster, for example, Figure 6 shows the sample count per cluster for the cases of total number of clusters = 608 and 1984. Unfortunately, we do not observe any meaningful improvement on the model’s validation accuracy trained on this more refined hierarchy.

Experimental procedures. All the validation accuracies we report on ResNet34 are the averaged results of experiments performed on at least 6 random seeds: 2 random seeds for backbone pretraining and 3 random seeds for finetuning. We report the average accuracies with their standard deviation on various hierarchies in Table 1.

An additional experiment we performed with ResNet34 is a small grid search over what checkpoint of a pretrained backbone we should use for finetuning on the 11-superclass method; we tried the 50-, 70- and 90-epoch checkpoints of the backbone on the manual hierarchies. We report these results in Table 2. As we can see, 90-epoch checkpoints performs almost equally well as the 70-epoch checkpoints and better than the 50-epoch ones by a nontrivial margin. With this observation, we chose to use the end-of-pretraining 90-epoch checkpoints in all our other experiments without further ablation studies on those hierarchies.

Table 2: **In-dataset transfer, iNaturalist 2021.** ResNet34 average finetuned validation error and standard deviation on 11 superclasses in iNaturalist 2021, pretrained on the manual hierarchies, with different backbone checkpoints.

| | | | | | | | |
|---------------|---------------------|------------------|------------------|------------------|------------------|------------------|------------------|
| 90-Epoch ckpt | $G(Y^{\text{src}})$ | 13 | 51 | 273 | 1103 | 4884 | 6485 |
| | Validation error | 5.40 ± 0.075 | 5.10 ± 0.038 | 4.83 ± 0.041 | 4.79 ± 0.045 | 4.82 ± 0.056 | 4.84 ± 0.033 |
| 70-Epoch ckpt | $G(Y^{\text{src}})$ | 13 | 51 | 273 | 1103 | 4884 | 6485 |
| | Validation error | 5.43 ± 0.055 | 5.08 ± 0.029 | 4.86 ± 0.037 | 4.82 ± 0.034 | 4.83 ± 0.064 | 4.85 ± 0.018 |
| 50-Epoch ckpt | $G(Y^{\text{src}})$ | 13 | 51 | 273 | 1103 | 4884 | 6485 |
| | Validation error | 5.53 ± 0.036 | 5.2 ± 0.031 | 4.90 ± 0.038 | 4.9 ± 0.042 | 4.91 ± 0.020 | 4.95 ± 0.026 |

Table 3: **In-dataset transfer, iNaturalist 2021.** ResNet50 finetuned average validation error and standard deviation on 11 superclasses in iNaturalist 2021, pretrained on label hierarchies with different label granularity.

| | | | | | | | | |
|------------------|---------------------|------------------------------------|------------------|------------------|------------------|------------------------------------|------------------|------------------|
| Manual Hierarchy | $G(Y^{\text{src}})$ | 11 | 13 | 51 | 273 | 1103 | 4884 | 6485 |
| | Validation error | 4.43 ± 0.029 | 4.44 ± 0.063 | 4.36 ± 0.062 | 4.22 ± 0.021 | 4.20 ± 0.035 | 4.23 ± 0.054 | 4.33 ± 0.037 |
| Random class ID | $G(Y^{\text{src}})$ | 22 | 88 | 352 | 1,408 | 5,632 | 11,264 | 500,000 |
| | Validation error | 5.36 ± 0.111 | 5.31 ± 0.079 | 5.24 ± 0.093 | 5.38 ± 0.052 | 5.37 ± 0.033 | 5.40 ± 0.033 | 5.13 ± 0.072 |

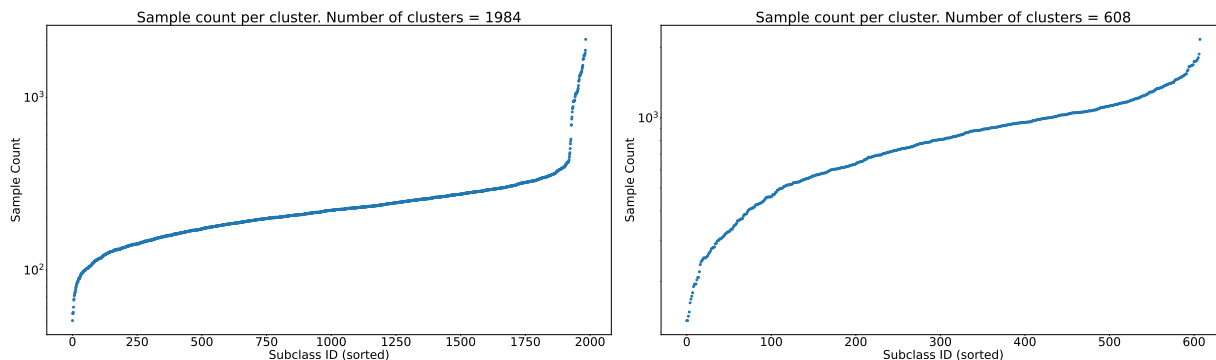


Figure 6: **In-dataset transfer, iNaturalist 2021.** Number of samples per cluster in the case of 608 and 1984 total clusters, after applying the sample size rebalancing procedure described in subsection A.1.1. Observe that the sample sizes are reasonably balanced across almost all the subclasses.

Table 4: **In-dataset transfer.** ViT-B/16 validation error on the binary problem “is this object a *living thing*?” of ImageNet21k. Pretrained on various hierarchy levels of ImageNet21k, finetuned on the binary problem. Observe that the maximal improvement appears at the leaf labels, and as $G(Y^{\text{src}})$ approaches 2, the percentage improvement approaches 0.

| Hierarchy level | $G(Y^{\text{src}})$ | Validation error |
|-----------------|---------------------|------------------|
| Baseline | 2 | 7.90 |
| 0 (leaf) | 21843 | 6.56 |
| 1 | 5995 | 6.76 |
| 2 | 2281 | 6.70 |
| 4 | 519 | 6.97 |
| 6 | 160 | 7.31 |
| 9 | 38 | 7.55 |

Our ResNet50 results are not as extensive as those on ResNet34. We present the average accuracies and standard deviations in Table 3.

A.1.2 ImageNet21k

The ImageNet21k dataset we experiment on contains a total of 12,743,321 training samples and 102,400 validation samples, with 21843 leaf labels. A small portion of samples have multiple labels.

Caution: due to the high demand on computational resources of training ViT models on ImageNet21k, all of our experiments that require (pre-)training or finetuning/linear probing on this dataset were performed with one random seed.

Hierarchy generation. To define fine-grained labels, we start by defining the leaf labels of the dataset to be Hierarchy level 0. For every image, we trace from the leaf synset to the root synset relying on the WordNet hierarchy, and set the k -th synset (or the root synset, whichever is higher in level) as the level- k label of this image; this procedure also applies to the multi-label samples. This is the way we generate the manual hierarchies shown in the main text.

Due to the lack of a predefined coarse-label problem, we manually define our target problem to be a binary one: given an image, if the synset “Living Thing” is present on the path tracing from the leaf label of the image to the root, assign label 1 to this image; otherwise, assign 0. This problem almost evenly splits the training and validation sets of ImageNet21k: 5,448,549:7,294,772 for training, 43,745:58,655 for validation.

Network choice and pretraining pipeline. We experiment with the ViT-B/16 model Dosovitskiy et al. (2021). The pretraining pipeline of this model follows the one in Dosovitskiy et al. (2021) exactly: we train the model for 90 epochs using the Adam optimizer, with $\beta_1 = 0.9, \beta_2 = 0.999$, weight decay coefficient equal to 0.03 and a batch size of 4096; we let the dropout rate be 0.1; the output dense layer’s bias is initialized to -10.0 to prevent huge loss value coming from the off-diagonal classes near the beginning of training Cui et al. (2019b); for learning rate, we perform linear warmup for 10,000 steps until the learning rate reaches 10^{-3} , then it is linearly decayed to 10^{-5} . The data augmentations are the common ones in ImageNet-type training Dosovitskiy et al. (2021); He et al. (2016): random cropping and horizontal flipping. Note that we use the sigmoid cross-entropy for training since the dataset has multi-label samples.

Each training instance (90 epochs) is run on 64 TPU v4 chips, taking approximately 1.5 to 2 days.

Evaluation on the binary problem. After the 90-epoch pretraining on the manual hierarchies, we evaluate the model on the binary problem. We report the best accuracies on each hierarchy level in Table 4. To get a sense of how the relevant hyperparameters influence final accuracy of the model, we try out the following finetuning/linear probing strategies on the backbone trained on the *leaf labels* and the target *binary problem* of the dataset, and report the results in Table 5 (similar to our experiments on iNaturalist, we include the backbone trained on the binary problem in these ablation studies to ensure that our comparisons against the baseline are fair) :

1. 90-epochs finetuning in the same fashion as the pretraining stage, but with a small grid search over

$$(\text{batch size, base learning rate}) = \{(4096, 0.001), (4096/4 = 1024, 0.001/4 = 0.00025), \\ (4096/8 = 512, 0.001/8 = 0.000125)\}.$$

2. Linear probing with 20 epochs training length, using exactly the same training pipeline as in pretraining. We ran a small grid search over $(\text{batch size, base learning rate}) = \{(4096, 0.001), (4096/8 = 512, 0.001/8 = 0.000125)\}$.
3. 10-epochs finetuning, no linear warmup, 3 epochs of constant learning rate in the beginning followed by 7 epochs of linear decay, with a small grid search over $(\text{batch size, base learning rate}) = \{(4096, 0.001), (4096/8 = 512, 0.001/8 = 0.000125)\}$.

Table 5 helps us decide the best accuracies to report. First, as expected the linear probing results are much worse than the finetuning ones. Second, the “retraining” accuracy of 92.102 is the best baseline we can report (the same thing happened in the iNaturalist case) — if we only train the model for 90 epochs (the naive one-pass training) on the binary problem, then the model’s final validation accuracy is 91.746%, which is lower than 92.102% by a nontrivial margin. In contrast, the short 10-epoch finetuning strategy

Table 5: **In-dataset transfer, ImageNet21k.** ViT-B/16 validation accuracy on the binary problem “Is the object a Living Thing” on ImageNet21k. Ablation study on the exact finetuning/linear probing strategy.

| | Eval strategy | 90-epoch finetune | | | Linear probe | | 10-epoch finetune | |
|-----------------|-----------------------|-------------------|---------------|---------------|--------------|----------------|-------------------|---------------|
| Leaf-pretrained | (Batch size, base lr) | (4096,1e-3) | (1024,2.5e-4) | (512,1.25e-4) | (4096, 1e-3) | (512, 1.25e-4) | (4096,1e-3) | (512,1.25e-4) |
| | Validation error | 92.782 | 93.177 | 93.295 | 87.497 | 87.493 | 92.294 | 93.439 |
| Baseline | (Batch size, base lr) | (4096,1e-3) | (1024,2.5e-4) | (512,1.25e-4) | (4096, 1e-3) | (512, 1.25e-4) | (4096,1e-3) | (512,1.25e-4) |
| | Validation error | 92.102 | 91.971 | 91.939 | 91.703 | 91.719 | 92.002 | 91.856 |

Table 6: **In-dataset transfer, ImageNet1k.** ResNet50 finetuned average validation error and standard deviation on the vanilla 1000 classes, pretrained on label hierarchies with different label granularity.

| | | | | |
|----------------------|---------------------|--------------|---------------|---------------|
| ResNet50 CLIP+kMeans | $G(Y^{\text{src}})$ | 2000 | 4000 | 8000 |
| per-class | Validation error | 23.4 ± 0.13 | 23.48 ± 0.098 | 23.49 ± 0.204 |
| ViT-L/14 CLIP+kMeans | $G(Y^{\text{src}})$ | 2000 | 4000 | 8000 |
| per-class | Validation error | 23.4 ± 0.127 | 23.47 ± 0.074 | 23.78 ± 0.048 |
| Random ID | $G(Y^{\text{src}})$ | 2000 | 4000 | 8000 |
| per-class | Validation error | 23.4 ± 0.068 | 23.4 ± 0.070 | 23.65 ± 0.071 |

works best for the backbone trained on the leaf labels, therefore, we also use this strategy to evaluate the backbones trained on all the other manual hierarchies. A peculiar observation we made was that, finetuning the leaf-labels-pretrained backbone for extended period of time on the binary problem caused it to overfit severely: for batch size and base learning rate in the set $\{(4096, 0.001), (1024, 0.00025), (512, 0.000125)\}$, throughout the 90 epochs of finetuning, although its training loss exhibits the normal behavior of staying mostly monotonically decreasing, its validation accuracy actually reached its peak during the linear warmup period!

A.1.3 ImageNet1k

Our ImageNet1k in-dataset transfer experiments are done in a very similar fashion to the iNaturalist ones. In particular, the pretraining and finetuning pipeline for ResNet50 is exactly the same as the one in the iNaturalist case, so we do not repeat it here.

Due to a lack of more fine-grained manual label on this dataset, we generate fine-grained labels by performing kMeans on the ViT-L/14 CLIP embedding of the dataset separately for each class; the exact procedure is also identical to the iNaturalist case. The CLIP backbones we use here are the ResNet50 version and the ViT-L/14 version. We report the average accuracies and their standard deviation in Table 6. All results are obtained from at least one random seed during pretraining and 3 random seeds during finetuning.

The best baseline we report is the one using retraining: if we adopt the pretrain-then-finetune procedure but with $D_{\text{train}}^{\text{tgt}}$ (i.e. the vanilla 1000-class labels) set as the pretraining dataset, then we obtain an average validation error of 23.28% with standard deviation of 0.103, averaged over results of 3 random seeds. In comparison, if we only perform the naive one-pass 90-epoch training, we obtain average validation error 24.04%, with standard deviation 0.057.

From Table 6, we see that there is virtually no difference between the baseline and the best errors obtained by the models trained on the custom hierarchies: they are almost equally bad. Noting that the sample size of each class in ImageNet1k is only around 10^3 , and the fact that ImageNet1k classification is a “hard problem” — it is a problem of high sample complexity — further decomposing the classes causes each fine-grained class to have too few samples, leading to the above negative results. This reflects the intuition that higher label granularity does not necessarily mean better model generalization, since the sample size per class might become too small.

Table 7: **Cross-dataset transfer.** ViT-B/16 average *finetuning* validation accuracy on ImageNet1k along with standard deviation, pretrained on various hierarchy levels of ImageNet21k, and a small grid search over the base learning rate.

| Pretrained on / Base lr | 3×10^{-3} | 3×10^{-2} | 6×10^{-2} | 3×10^{-1} |
|--------------------------|--------------------|--------------------|--------------------|--------------------|
| ImageNet21k, Hier. lv. 0 | 80.87±0.012 | 82.48±0.005 | 82.51±0.042 | 81.40±0.041 |
| ImageNet21k, Hier. lv. 1 | 77.38±0.037 | 81.03±0.054 | 81.28±0.045 | 80.40±0.087 |
| ImageNet21k, Hier. lv. 2 | 74.91±0.012 | 79.76±0.021 | 80.26±0.05 | 79.7±0.019 |
| ImageNet21k, Hier. lv. 4 | 63.65±0.052 | 76.43±0.033 | 77.32±0.088 | 77.53±0.078 |
| ImageNet21k, Hier. lv. 6 | 62.17±0.012 | 73.65±0.033 | 73.92±0.073 | 75.53±0.024 |
| ImageNet21k, Hier. lv. 9 | 53.68±0.034 | 69.33±0.045 | 71.08±0.068 | 72.75±0.071 |

Table 8: **Cross-dataset transfer.** ViT-B/16 average *linear-probing* validation accuracy on ImageNet1k along with standard deviation, pretrained on various hierarchy levels of ImageNet21k.

| Pretrained on | Hier. lv | $G(Y^{\text{src}})$ | Validation acc. |
|---------------|----------|---------------------|--------------------|
| IM21k | 0 (leaf) | 21843 | 81.45±0.021 |
| | 1 | 5995 | 78.33±0.018 |
| | 2 | 2281 | 75.66±0.005 |
| | 4 | 519 | 68.95±0.051 |
| | 6 | 160 | 63.65±0.035 |
| | 9 | 38 | 57.35±0.016 |

A.2 Cross-dataset transfer, ImageNet21k ImageNet1k

In this subsection, we report the average validation accuracy and standard deviation of the cross-dataset transfer experiment from ImageNet21k to ImageNet1k, as discussed in Figure 2 and Section 1 in the main text.

Network choice. We use the same architecture ViT-B/16 as the one in the in-dataset ImageNet21k transfer experiment and follow the same training procedure, which we repeat here for the reader’s convenience. The pretraining pipeline of this model follows the one in Dosovitskiy et al. (2021): we train the model for 90 epochs using the Adam optimizer, with $\beta_1 = 0.9, \beta_2 = 0.999$, weight decay coefficient equal to 0.03 and a batch size of 4096; we let the dropout rate be 0.1; the output dense layer’s bias is initialized to -10.0 to prevent huge loss value coming from the off-diagonal classes near the beginning of training Cui et al. (2019b); for learning rate, we perform linear warmup for 10,000 steps until the learning rate reaches 10^{-3} , then it is linearly decayed to 10^{-5} . The data augmentations are the common ones in ImageNet-type training Dosovitskiy et al. (2021); He et al. (2016): random cropping and horizontal flipping. Note that we use the sigmoid cross-entropy for training since the dataset has multi-label samples.

Additionally, each training instance (90 epochs) is run on 64 TPU v4 chips, taking approximately 1.5 to 2 days.

Finetuning. For finetuning on ImageNet1k, our procedure is very similar to the one in the original ViT paper Dosovitskiy et al. (2021), described in its Appendix B.1.1. We optimize the network for 8 epochs using SGD with momentum factor set to 0.9, zero weight decay, and batch size of 512. The dropout rate, unlike in pretraining, is set to 0. Gradient clipping at 1.0 is applied. Unlike Dosovitskiy et al. (2021), we still finetune at the resolution of 224×224 . For learning rate, we apply linear warmup for 500 epochs until it reaches the base learning rate, then cosine annealing is applied; we perform a small grid search of base learning rate = $\{3 \times 10^{-3}, 3 \times 10^{-2}, 6 \times 10^{-2}, 3 \times 10^{-1}\}$. Every one of these grid search is repeated over

3 random seeds. We report the ImageNet1k validation accuracies and their standard deviations in Table 7. In the main text, we report the best accuracy for each hierarchy level.

Linear probing. For linear probing, we use the following procedure. We optimize the linear classifier for 40 epochs (similar to Lee et al. (2021)) using SGD with Nesterov momentum factor set to 0.9, a small weight decay coefficient 10^{-6} , and batch size 512. We start with a base learning rate of 0.9, and multiply it by 0.97 per 0.5 epoch. In terms of data augmentation, we adopt the standard ones like before: horizontal flipping and random cropping of size 224×224 . We repeat this linear probing procedure over 3 random seeds given the pretrained backbone, and report the average validation accuracy and standard deviation in Table 8.

Baseline. The baseline accuracy on ImageNet1k is directly taken from the ViT paper Dosovitskiy et al. (2021) (see Table 5 in it), in which the ViT-B/16 model is trained for 300 epochs on ImageNet1k.

B Theory, Problem Setup

B.1 Data Properties

1. Coarse classification: a binary task, +1 vs. -1.
2. An input sample $\mathbf{X} \in \mathbb{R}^{d \times P}$ consists of P patches, each with dimension d . In this work, always assume d is sufficiently large¹;
3. Assume there exists k_+ subclasses of the superclass “+”, and k_- subclasses of the superclass “-”. Let $k_+ = k_-$.
4. Assume orthonormal dictionary $V = \{\mathbf{v}_1, \dots, \mathbf{v}_d\} \in \mathbb{R}^d$, which forms an orthonormal basis of \mathbb{R}^d . Define $\mathbf{v}_+ \in V$ to be the common feature of class “+”. For each subclass $(+, c)$ (where $c \in [k_+]$), denote the subclass feature of it as $\mathbf{v}_{+,c} \in V$. Similar for the “-” class.
5. For an easy sample \mathbf{X} belonging to the $(+, c)$ class (for $c \in [k_+]$), we sample its patches as follows:

Definition: we define the function $P : \mathbb{R}^{d \times P} \times V \rightarrow [P]$ (so $(\mathbf{X}; \mathbf{v}) \rightarrow I \in [P]$) to extract, from sample \mathbf{X} , the indices of the patches on which the dictionary word $\mathbf{v} \in V$ dominates.

 - (a) (Common-feature patches) With probability $\frac{s}{P}$, a patch \mathbf{x}_p in \mathbf{X} is a common-feature patch, on which $\mathbf{x}_p = \alpha_p \mathbf{v}_+ + \mathbf{p}$ for some (random) $\alpha_p \in [\frac{1-l}{1+l}, \frac{1+l}{1-l}]$;
 - (b) (Subclass-feature patches) With probability $\frac{s}{P - |P(\mathbf{X}; \mathbf{v}_+)|}$, a patch with index $p \in ([P] - P(\mathbf{X}; \mathbf{v}_+))$ is a subclass-feature patch, on which $\mathbf{x}_p = \alpha_p \mathbf{v}_{+,c} + \mathbf{p}$, for random $\alpha_p \in [\frac{1-l}{1+l}, \frac{1+l}{1-l}]$;
 - (c) (Noise patches) For the remaining $P - |P(\mathbf{X}; \mathbf{v}_+)| - |P(\mathbf{X}; \mathbf{v}_{+,c})|$ patches, $\mathbf{x}_p = \mathbf{p}$.
6. A hard sample \mathbf{X}_{hard} for class $(+, c)$ is exactly the same as an easy one except:
 - (a) Its common-feature patches are replaced by noise patches;
 - (b) (Feature noise patches) With probability $\frac{s}{P - |P(\mathbf{X}; \mathbf{v}_{+,c})|}$, a patch with index $p \in ([P] - P(\mathbf{X}; \mathbf{v}_{+,c}))$ is a feature-noise patch, on which $\mathbf{x}_p = \alpha_p^\dagger \mathbf{v}_- + \mathbf{p}$ for some (random) $\alpha_p^\dagger \in [t_{\text{lower}}^\dagger, t_{\text{upper}}^\dagger]$;
 - (c) Set one of the noise patches to $\mathbf{x} \sim N(\mathbf{0}, \sigma_\zeta^2 \mathbf{I}_d)$.
7. A sample \mathbf{X} belongs to the “+” superclass if $|P(\mathbf{X}; \mathbf{v}_+)| > 0$ or $|P(\mathbf{X}; \mathbf{v}_{+,c})| > 0$ for any c (excluding feature-noise patches).
8. The above sample definitions also apply to the “-” classes by switching the class signs.
9. A training batch of samples contains exactly $N/2k_+$ samples for each $(+, c)$ and $(-, c)$ subclass. This also means that each training batch contains exactly $N/2$ samples belonging to the +1 superclass, and $N/2$ samples for the -1 superclass.
10. As discussed in the main text, for both coarse-grained (baseline) and fine-grained training, we only train on *easy* samples.

B.2 Learner Assumptions and Training Algorithm

Assume the learner is a two-layer convolutional ReLU network:

$$F_c(\mathbf{X}) = \sum_{r=1}^m a_{c,r} \sum_{p=1}^P \sigma(\mathbf{w}_{c,r} \cdot \mathbf{x}_p + b_{c,r}) \quad (14)$$

To simplify analysis and only focus on the learning of the feature extractor, we freeze $a_{c,r} = 1$ throughout training. The nonlinear activation $\sigma(\cdot) = \max(0, \cdot)$ is ReLU. Note that the convolution kernels have dimension d and stride d .

¹Consider each d -dimensional patch of the input as an embedding of the input image generated by, for instance, an intermediate layer of a DNN.

Remark. One difference between this architecture and a CNN used in practice is that we do not allow feature sharing across classes: for each class c , we are assigning a disjoint group of neurons $\mathbf{w}_{c,r}$ to it. Separating neurons for each class is a somewhat common trick to lower the complexity of analysis in DNN theory literature Allen-Zhu & Li (2023b); Karp et al. (2021); Cao et al. (2022), as it reduces complex coupling between neurons *across* classes which is not the central focus of our study in this paper.

Now we discuss the **training algorithm**.

Initialization.

Sample $\mathbf{w}_{c,r}^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 I_d)$, and set $b_{c,r}^{(0)} = -\sigma_0 c_b \sqrt{\log(d)}$.

Training.

We adopt the standard cross-entropy training:

$$L(F) = \sum_{n=1}^N L(F; \mathbf{X}_n, y_n) = - \sum_{n=1}^N \log \left(\frac{\exp(F_{y_n}(\mathbf{X}_n))}{\sum_{c=1}^C \exp(F_c(\mathbf{X}_n))} \right) \quad (15)$$

This induces the stochastic gradient descent update for each hidden neuron ($c \in [k], r \in [m]$) per minibatch of N iid samples:

$$\begin{aligned} \mathbf{w}_{c,r}^{(t+1)} = \mathbf{w}_{c,r}^{(t)} + \eta \frac{1}{NP} \sum_{n=1}^N \left(\mathbb{1}\{y_n = c\} [1 - \text{logit}_c^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{c,r}^{(t)} \cdot \mathbf{x}_{n,p}^{(t)} + b_{c,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \right. \\ \left. \mathbb{1}\{y_n = c\} [-\text{logit}_c^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{c,r}^{(t)} \cdot \mathbf{x}_{n,p}^{(t)} + b_{c,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \right) \end{aligned} \quad (16)$$

where

$$\text{logit}_c^{(t)}(\mathbf{X}) = \frac{\exp(F_c(\mathbf{X}))}{\sum_{y=1}^C \exp(F_y(\mathbf{X}))} \quad (17)$$

As for the bias,

$$b_{c,r}^{(t+1)} = b_{c,r}^{(t)} - \frac{\mathbf{w}_{c,r}^{(t+1)} - \mathbf{w}_{c,r}^{(t)}}{\log^5(d)} \quad (18)$$

Remark. 1. The initialization strategy is similar to the one in Allen-Zhu & Li (2022).

2. Since the only difference between the training samples of coarse and fine-grained pretraining is the label space, the form of SGD update is identical. The only difference is the number of output nodes of the network: for coarse training, the output nodes are just F_+ and F_- (binary classification), while for fine-grained training, the output nodes are $F_{+,1}, F_{+,2}, \dots, F_{+,k_+}, F_{-,1}, F_{-,2}, \dots, F_{-,k_-}$, a total of $k_+ + k_-$ nodes.
3. The bias is for thresholding out the neuron's noisy activations that grow slower than $1/\log^5(d)$ times the activations on the features which the neuron detects. This way, the bias does not really influence updates to the neuron's response to the (common and/or fine-grained) features which it activates strongly on, since $1 - \frac{1}{\log^5(d)} \approx 1$, while it removes useless low-magnitude noisy activations. This in fact creates a (generalization) gap between the nonlinear model that we are studying and linear models. Due to our parameter choices (as discussed below), if the model has no nonlinearity (remove the ReLU activations), then even if the model can be written as $F_+(\mathbf{X}) = \sum_{p \in [P]} c_+ \mathbf{v}_+ \cdot \mathbf{x}_p + c_{+,1} \mathbf{v}_{+,1} \cdot \mathbf{x}_p + \dots + c_{+,k_+} \mathbf{v}_{+,k_+} \cdot \mathbf{x}_p$ and $F_-(\mathbf{X}) = \sum_{p \in [P]} c_- \mathbf{v}_- \cdot \mathbf{x}_p + c_{-,1} \mathbf{v}_{-,1} \cdot \mathbf{x}_p + \dots + c_{-,k_-} \mathbf{v}_{-,k_-} \cdot \mathbf{x}_p$ for any sequence of nonnegative real numbers $c_+, c_-, \{c_{+,j}\}_{j=1}^{k_+}, \{c_{-,j}\}_{j=1}^{k_-}$ (which is the ideal situation since the true features are not corrupted by anything), it is impossible for the model to reach $o(1)$ error on the input samples, because the number of noise patches will accumulate to a variance of $(P - O(s^*)) \sigma_\zeta^2 = O(s^*)$, which significantly overwhelms the signal from the true features. On the

other hand, each noise patch is sufficiently small in magnitude with high probability (their strength is $o(1/\log^5(d))$), so a slightly negative bias, as described above, can threshold out these noise-based signals and prevent them from accumulating across the patches.

An important difference between our bias update rule and the one in Allen-Zhu & Li (2022) is that, our rule depends on the ℓ_2 norm of the neuron’s update, while the one in Allen-Zhu & Li (2022) is hard-coded and not dependent on the neuron weights. The reason that we should not hard code the bias update rate is that, the neurons that are responsible for detecting the common features will grow more quickly in norm than those responsible for detecting the fine-grained features, therefore, to ensure fairness between the different groups of neurons (i.e. only using the bias to remove useless activations on the noise patches while creating minimal disturbance to the neurons’ activation on feature-dominated patches), we rely on our neuron-dependent bias update rule.

B.3 Parameter Choices

The following are fixed choices of parameters for the sake of simplicity in our proofs.

1. Always assume d is sufficiently large. All of our asymptotic results are presented with respect to d ;
2. $\text{poly}(d)$ denotes the asymptotic order “polynomial in d ”;
3. $\text{polylog}(d)$ asymptotic order “polylogarithmic in d ”;
4. $\text{polylog}(d)$ $k_+ = k_- d^{0.4}$ and $s^* \log^5(d)$ k_+ (i.e. k_+ lower bounded by polynomial of $\log(d)$ of sufficiently high degree);
5. Small positive constant $c_0 \in (0, 0.1)$;
6. For coarse-grained (baseline) training, set $c_b = \frac{4 + 2c_0}{2 + 2c_0}$, and for fine-grained training, set $c_b = \frac{1}{2 + 2c_0}$;
7. $0 < \iota \leq \frac{1}{\text{polylog}(d)}$;
8. $\iota_{lower}^\dagger \leq \frac{1}{\log^4(d)}$, and $s^\dagger \iota_{upper}^\dagger \leq O\left(\frac{1}{\log(d)}\right)$;
9. $s^\dagger \leq 1$;
10. $s^* \leq \text{polylog}(d)$ with a degree > 15 ;
11. $\sigma_\zeta = \frac{1}{\log^{10}(d)\sqrt{d}}$;
12. $\sigma_\zeta \in \left[\omega\left(\frac{\text{polylog}(d)}{\sqrt{d}}\right), O\left(\frac{1}{\text{polylog}(d)}\right)\right]$;
13. $P\sigma_\zeta \leq \omega(\text{polylog}(d))$, and $P \leq \text{poly}(d)$;
14. $\sigma_0 \leq O\left(\frac{1}{d^{3s} \log(d)}\right)$, and set $\eta = \Theta(\sigma_0)$ for simplicity;
15. Batch of samples $B^{(t)}$ at every iteration has a deterministic size of $N \in (\Omega(\text{polylog}(d)k_+d), \text{poly}(d))$.
16. Note: we sometimes abuse the notation $x = a \pm b$ as an abbreviation for $x \in [a - b, a + b]$.

Remark. We believe the range of parameter choice can be (asymptotically) wider than what is considered here, but for the purpose of illustrating the main messages of the paper, we do not consider a more general set of parameter choice necessary because having a wider range of it can significantly complicate and obscure the already lengthy proofs without adding to the core messages.

B.4 Plan of presentation and central ideas

We shall devote the majority of our effort to proving results for the coarse-label learning dynamics, starting with appendix section C and ending on E, and only devote section G to the fine-grained-label learning dynamics, since the analysis of fine-grained training overlaps significantly with the coarse-grained one.

One technical difficulty in making the above ideas rigorous lies in the ReLU activation (with time-dependent bias): due to randomness in the gradient updates and the initialization, it is possible for individual hidden neurons that activate on \mathcal{V} -dominated patches at one time iterate to no longer do so at the next iterate, and the opposite can happen. This can be problematic: for instance, it is possible that certain “lucky” neurons for \mathcal{V}_+ at one iterate become dead on \mathcal{V}_+ -dominated patches at the next iterate, while some “unlucky” neurons that were dead on $\mathcal{V}_{+,c}$ -dominated patches before start activating on these patches at the current iterate. In our proof, we show that this kind of situation does not happen too frequently nor do they contribute too much to the overall behavior of the neural network, by carefully keeping track of each hidden neuron’s response to feature vectors and noise vectors throughout training.

C Coarse-grained training, Initialization Geometry

For coarse-grained training, assume $m = \Theta(d^{2+2c_0})$.

Definition C.1. Define the following sets of interest of the hidden neurons:

1. $U_{+,r}^{(0)} = \{\mathbf{v} \in \mathcal{V} : \mathbf{w}_{+,r}^{(0)}(\mathbf{v}) \geq \sigma_0 \sqrt{\frac{1}{4+2c_0} \log(d) - \frac{1}{\log^5(d)}}\}$
2. Given $\mathbf{v} \in \mathcal{V}$, $S_+^{*(0)}(\mathbf{v}) \subset \mathcal{V} \times [m]$ satisfies:
 - (a) $\mathbf{w}_{+,r}^{(0)}(\mathbf{v}) \geq \sigma_0 \sqrt{\frac{1}{4+2c_0} \log(d) + \frac{1}{\log^5(d)}}$
 - (b) $\mathbf{v}' \in \mathcal{V}$ s.t. $\mathbf{w}_{+,r}^{(0)}(\mathbf{v}') < \sigma_0 \sqrt{\frac{1}{4+2c_0} \log(d) - \frac{1}{\log^5(d)}}$
3. Given $\mathbf{v} \in \mathcal{D}$, $S_+^{(0)}(\mathbf{v}) \subset \mathcal{V} \times [m]$ satisfies:
 - (a) $\mathbf{w}_{+,r}^{(0)}(\mathbf{v}) \geq \sigma_0 \sqrt{\frac{1}{4+2c_0} \log(d) - \frac{1}{\log^5(d)}}$
4. For any $(+, r) \in S_{+,reg}^{*(0)} \subset \mathcal{V} \times [m]$:
 - (a) $\mathbf{w}_{+,r}^{(0)}(\mathbf{v}) \geq \sigma_0 \sqrt{\frac{1}{10} \log(d)}$ $\forall \mathbf{v} \in \mathcal{V}$
 - (b) $|U_{+,r}^{(0)}| = O(1)$

Proposition 1. Assume $m = \Theta(d^{2+2c_0})$, i.e. the number of neurons assigned to the + and - class are equal and set to $\Theta(d^{2+2c_0})$.

At $t = 0$, for all $\mathbf{v} \in \mathcal{V}$, the following properties are true with probability at least $1 - d^{-2}$ over the randomness of the initialized kernels:

1. $|S_+^{*(0)}(\mathbf{v})|, |S_+^{(0)}(\mathbf{v})| = \Theta\left(\frac{1}{\log(d)}\right) d^{c_0}$
2. In particular, for any $\mathbf{v}, \mathbf{v}' \in \mathcal{V}$, $\left| \frac{|S_+^{*(0)}(\mathbf{v})|}{|S_+^{*(0)}(\mathbf{v}')|} - 1 \right|, \left| \frac{|S_+^{(0)}(\mathbf{v})|}{|S_+^{(0)}(\mathbf{v}')|} - 1 \right| = O\left(\frac{1}{\log^5(d)}\right)$
3. $S_{+,reg}^{(0)} = [m]$

Proof. Recall the tail bound of $g \sim \mathcal{N}(0, 1)$ for every $\epsilon > 0$:

$$\frac{1}{2} \frac{1}{2\pi} \frac{\epsilon}{\epsilon^2 + 1} e^{-\epsilon^2/2} \leq \mathbb{P}[g \leq -\epsilon] \leq \frac{1}{2} \frac{1}{2\pi} \frac{1}{\epsilon} e^{-\epsilon^2/2} \quad (19)$$

First note that for any $r \in [m]$, $\{\mathbf{w}_{+,r}^{(0)}(\mathbf{v})\}_{\mathbf{v} \in \mathcal{V}}$ is a sequence of iid random variables with distribution $\mathcal{N}(0, \sigma_0^2)$.

The proof of the first point proceeds in two steps.

1. The following properties hold at $t = 0$:

$$\begin{aligned}
p_1 &:= \mathbb{P} \left[\mathcal{W}_{+,r}^{(0)}, \mathbf{v} \quad \sigma_0 \quad \frac{1}{4+2c_0} \sqrt{\log(d) + \frac{1}{\log^5(d)}} \right] \\
&\quad \frac{1}{8\pi} d^{-2-c_0} e^{(-2-c_0)/\log^5(d)} \\
&\quad \times \left[\frac{\sqrt{(4+2c_0) \left(\log(d) + \frac{1}{\log^5(d)} \right)}}{(4+2c_0) \left(\log(d) + \frac{1}{\log^5(d)} \right) + 1}, \frac{1}{\sqrt{(4+2c_0) \left(\log(d) + \frac{1}{\log^5(d)} \right)}} \right] \\
&= \Theta \left(\frac{1}{\sqrt{\log(d)}} \right) d^{-2-c_0}
\end{aligned} \tag{20}$$

and

$$\begin{aligned}
p_2 &:= \mathbb{P} \left[\mathcal{W}_{+,r}^{(0)}, \mathbf{v} \quad \sigma_0 \quad \frac{1}{4+2c_0} \sqrt{\log(d) - \frac{1}{\log^5(d)}} \right] \\
&\quad \frac{1}{8\pi} d^{-2-c_0} e^{-(-2-c_0)/\log^5(d)} \\
&\quad \times \left[\frac{\sqrt{(4+2c_0) \left(\log(d) - \frac{1}{\log^5(d)} \right)}}{(4+2c_0) \left(\log(d) - \frac{1}{\log^5(d)} \right) + 1}, \frac{1}{\sqrt{(4+2c_0) \left(\log(d) - \frac{1}{\log^5(d)} \right)}} \right] \\
&= \Theta \left(\frac{1}{\sqrt{\log(d)}} \right) d^{-2-c_0}
\end{aligned} \tag{21}$$

Therefore, for any $r \in [m]$, the random event described in $S_+^{*(0)}$ holds with probability

$$\begin{aligned}
p_1 \times (1 - p_2)^{d-1} &= \Theta \left(\frac{1}{\sqrt{\log(d)}} \right) d^{-2-c_0} \times \left(1 - \Theta \left(\frac{1}{\sqrt{\log(d)}} \right) d^{-2-c_0} \right)^{d-1} \\
&= \Theta \left(\frac{1}{\sqrt{\log(d)}} \right) d^{-2-c_0}.
\end{aligned} \tag{22}$$

The last equality holds because defining $f(d) = d^{-2-c_0}$ and d being sufficiently large,

$$g(d) := \frac{(d-1) \log(1-f(d))}{(d-1) \times (f(d) + O(f(d)^2))} = O(d^{-1}) \tag{23}$$

which means

$$(1 - f(d))^{d-1} = e^{-g(d)} = (1 - O(d^{-1}))^{d-1} \tag{24}$$

2. Given $\mathbf{v} \in \mathcal{V}$, $|S_+^{*(0)}(\mathbf{v})|$ is a binomial random variable, with each Bernoulli trial (ranging over $r \in [m]$) having success probability $p_1(1 - p_2)^{d-1}$. Therefore, $\mathbb{E} \left[|S_+^{*(0)}(\mathbf{v})| \right] = mp_1(1 - p_2)^{d-1} = \Theta \left(\frac{1}{\log(d)} \right) d^{c_0}$.

Now recall the Chernoff bound of binomial random variables. Let $\{X_n\}_{n=1}^m$ be an iid sequence of Bernoulli random variable with success rate p , and $S_n = \sum_{n=1}^m X_n$. Then for any $\delta \in (0, 1)$,

$$\begin{aligned}
\mathbb{P}[S_n \geq (1 + \delta)mp] &\leq \exp \left(-\frac{\delta^2 mp}{3} \right) \\
\mathbb{P}[S_n \leq (1 - \delta)mp] &\leq \exp \left(-\frac{\delta^2 mp}{2} \right)
\end{aligned} \tag{25}$$

It follows that, for each $\mathbf{v} \in V$, $|S_+^{*(0)}(\mathbf{v})| = \Theta\left(\frac{1}{\log(d)}\right) d^{c_0}$ with probability at least $1 - \exp(-\Omega(\log^{-1/2}(d))d^{c_0})$. Taking union bound over all possible $\mathbf{v} \in D$, the random event still holds with probability at least $1 - \exp(-\Omega(\log^{-1/2}(d))d^{c_0} + O(\log(d))) = 1 - \exp(-\Omega(d^{0.5c_0}))$ (in sufficiently high dimension).

The proof for $S_+^{(0)}(\mathbf{v})$ proceeds in virtually the same way, so we omit the calculations here.

To show the second point, in particular $\left| \frac{|S_+^{(0)}(\mathbf{v})|}{|S_+^{*(0)}(\mathbf{v})|} - 1 \right| = O\left(\frac{1}{\log^5(d)}\right)$, we need to be a bit more careful in our bounds of the relevant sets. In particular, we need to directly use the CDF of gaussian random variables:

$$\begin{aligned}
& \left| \mathbb{P} \left[\mathbf{w}_{+,r}^{(0)}, \mathbf{v} \in \sigma_0 \sqrt{\frac{1}{4+2c_0} \log(d) + \frac{1}{\log^5(d)}} \right] (1 \pm O(d^{-1})) \right. \\
& \quad \left. - \mathbb{P} \left[\mathbf{w}_{+,r}^{(0)}, \mathbf{v} \in \sigma_0 \sqrt{\frac{1}{4+c_0} \log(d) - \frac{1}{\log^5(d)}} \right] \right| \\
& \leq \frac{1}{2} \frac{1}{2\pi} \int_{\sqrt{\frac{1}{4+2c_0} \log(d) - \frac{1}{\log^5(d)}}}^{\sqrt{\frac{1}{4+2c_0} \log(d) + \frac{1}{\log^5(d)}}} e^{-\epsilon^2/2} d\epsilon + O\left(\frac{1}{d^{3+c_0} \sqrt{\log(d)}}\right) \\
& \leq \frac{1}{2} \frac{1}{2\pi} d^{-2-c_0} e^{(2+c_0)/\log^5(d)} \frac{1}{4+2c_0} \left(\sqrt{\log(d) + \frac{1}{\log^5(d)}} - \sqrt{\log(d) - \frac{1}{\log^5(d)}} \right) \\
& \quad + O\left(\frac{1}{d^{3+c_0} \sqrt{\log(d)}}\right) \\
& = \frac{1}{2} \frac{1}{2\pi} d^{-2-c_0} e^{(2+c_0)/\log^5(d)} \frac{1}{4+2c_0} \frac{\frac{2}{\log^5(d)}}{\sqrt{\log(d) + \frac{1}{\log^5(d)}} + \sqrt{\log(d) - \frac{1}{\log^5(d)}}} + O\left(\frac{1}{d^{3+c_0} \sqrt{\log(d)}}\right)
\end{aligned} \tag{26}$$

The expected difference in number between the two sets is just the above expression multiplied by $m = \Theta(d^{2+2c_0})$, and with probability at least $1 - \exp(-\Omega(d^{-c_0/4}))$, the difference term satisfies

$$\begin{aligned}
& \frac{1}{2} \frac{1}{2\pi} (1 \pm d^{-c_0/2}) \Theta(d^{c_0}) e^{(2+c_0)/\log^5(d)} \frac{1}{4+2c_0} \frac{\frac{2}{\log^5(d)}}{\sqrt{\log(d) + \frac{1}{\log^5(d)}} + \sqrt{\log(d) - \frac{1}{\log^5(d)}}} \\
& \pm O\left(\frac{d^{2+2c_0}}{d^{3+c_0} \sqrt{\log(d)}}\right) \\
& \Theta\left(\frac{1}{\sqrt{\log(d)}}\right) d^{c_0} \times \frac{1}{\log^5(d)}
\end{aligned} \tag{27}$$

By further noting from before that $|S_+^{(0)}(\mathbf{v})| = \Theta\left(\frac{1}{\log(d)}\right) d^{c_0}$, $\left| \frac{|S_+^{(0)}(\mathbf{v})|}{|S_+^{*(0)}(\mathbf{v})|} - 1 \right| = O\left(\frac{1}{\log^5(d)}\right)$ follows. The proof of $\left| \frac{|S_+^{(0)}(\mathbf{v})|}{|S_+^{*(0)}(\mathbf{v})|} - 1 \right| = O\left(\frac{1}{\log^5(d)}\right)$ follows a very similar argument, so we omit the calculations here.

Now, as for the set $S_{reg}^{(0)}$, we know for any $r \in [m]$ and $\mathbf{v}_i \in D$,

$$\mathbb{P} \left[\mathbf{w}_{+,r}^{(0)}, \mathbf{v}_i \in \sigma_0 \sqrt{\frac{1}{10} \log(d)} \right] = O\left(\frac{1}{\sqrt{\log(d)}}\right) d^{-5}. \tag{28}$$

Taking the union bound over r and i yields

$$\mathbb{P} \left[r \text{ and } i \text{ s.t. } \mathbf{w}_{+,r}^{(0)}, \mathbf{v}_i \in \sigma_0 \sqrt{\frac{1}{10} \log(d)} \right] \leq mdO\left(\frac{1}{\sqrt{\log(d)}}\right) d^{-5} < d^{-2}. \tag{29}$$

Finally, to show $\left|U_{+,r}^{(0)}\right| = O(1)$ holds for every $(+, r)$, we just need to note that for any arbitrary $(+, r)$ neuron, the probability of $\left|U_{+,r}^{(0)}\right| > 4$ is no greater than

$$p_2^4 \binom{d}{4} = O\left(\frac{1}{\log^2 d}\right) d^{-8-4c_0} \times d^4 = O\left(\frac{1}{\log^2 d}\right) d^{-4-4c_0} \quad (30)$$

Taking union bound over all $m = O(d^{2+2c_0})$ neurons yields the desired result.

□

D Coarse-grained SGD Phase I: (Almost) Constant Loss, Neurons Diversify

Definition D.1. We define T_0 to be the first time which there exists some sample n such that

$$F_c^{(T_0)}(\mathbf{X}_n^{(T_0)}) \leq d^{-1} \quad (31)$$

Without loss of generality assume $c = +$. Define phase I to be the time $t \in [0, T_0)$.

D.1 Main results

Theorem D.1 (Phase 1 SGD update properties). *The following properties hold with probability at least $1 - O\left(\frac{mNPk_+ t}{\text{poly}(d)}\right) - O(e^{-\Omega(\log^2(d))})$ for every $t \in [0, T_0)$.*

1. (On-diagonal common-feature neuron growth) For every $(+, r), (+, r') \in S_+^{*(0)}(\mathbf{v}_+)$,

$$\mathbf{w}_{+,r}^{(t)} - \mathbf{w}_{+,r}^{(0)} = \mathbf{w}_{+,r'}^{(t)} - \mathbf{w}_{+,r'}^{(0)} \quad (32)$$

Moreover,

$$\Delta \mathbf{w}_{+,r}^{(t)} = \eta \left(\left(\frac{1}{2} \pm \psi_1 \right) \frac{1}{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \pm O\left(\frac{1}{\log^{10}(d)}\right) \right) \frac{s^*}{2P} \mathbf{v}_+ + \Delta_{+,r}^{(t)} \quad (33)$$

where $\Delta_{+,r}^{(t)} \sim N(\mathbf{0}, \sigma_{\zeta_{+,r}}^{(t)2} \mathbf{I})$, $\sigma_{\zeta_{+,r}}^{(t)} = \eta \sigma_\zeta \left(\left(\frac{1}{2} \pm \psi_1 \right) \frac{1}{1 \pm s^{*-1/3}} \right) \frac{\sqrt{s}}{P\sqrt{2N}}$, and $|\psi_1| \leq d^{-1}$.

Furthermore, every $(+, r) \in S_+^{*(0)}(\mathbf{v}_+)$ activates on \mathbf{v}_+ -dominated patches at time t .

2. (On-diagonal finegrained-feature neuron growth) For every possible choice of c and every $(+, r), (+, r') \in S_+^{*(0)}(\mathbf{v}_{+,c})$,

$$\mathbf{w}_{+,r}^{(t)} - \mathbf{w}_{+,r}^{(0)} = \mathbf{w}_{+,r'}^{(t)} - \mathbf{w}_{+,r'}^{(0)} \quad (34)$$

Moreover,

$$\Delta \mathbf{w}_{+,r}^{(t)} = \eta \left(\left(\frac{1}{2} \pm \psi_1 \right) \frac{1}{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \pm O\left(\frac{1}{\log^{10}(d)}\right) \right) \frac{s^*}{2k_+ P} \mathbf{v}_{+,c} + \Delta_{+,r}^{(t)} \quad (35)$$

where $\Delta_{+,r}^{(t)} \sim N(\mathbf{0}, \sigma_{\zeta_{+,r}}^{(t)2} \mathbf{I})$, and $\sigma_{\zeta_{+,r}}^{(t)} = \eta \sigma_\zeta \left(\left(\frac{1}{2} \pm \psi_1 \right) \frac{1}{1 \pm s^{*-1/3}} \right) \frac{\sqrt{s}}{P \cdot 2Nk_+}$.

Furthermore, every $(+, r) \in S_+^{*(0)}(\mathbf{v}_{+,c})$ activates on $\mathbf{v}_{+,c}$ -dominated patches at time t .

3. The above results also hold with the “+” and “−” signs flipped.

Proof. The SGD update rule produces the following update:

$$\mathbf{w}_{+,r}^{(t+1)} = \mathbf{w}_{+,r}^{(t)} + \eta \frac{1}{NP} \times \quad (36)$$

$$\sum_{n=1}^N \left(\mathbb{1}\{y_n = +\} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{+,r}^{(t)} \cdot \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \right) \quad (37)$$

$$+ \mathbb{1}\{y_n = -\} [-\text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{+,r}^{(t)} \cdot \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \quad (38)$$

In particular,

$$\begin{aligned}
\text{equation 37} &= \sum_{n=1}^N \mathbb{1}\{y_n = +\} \left(\frac{1}{2} \pm \psi_1 \right) \times \\
&\quad \left\{ \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v}_+)| > 0\} \left[\sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \sigma'(\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} + b_{+,r}^{(t)}) \left(\alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} \right) \right. \right. \\
&\quad \left. \left. + \sum_{p \notin \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \sigma'(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \right] \right. \\
&\quad \left. + \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v}_+)| = 0\} \sum_{p \in [P]} \sigma'(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \right\} \\
&= \sum_{n=1}^N \mathbb{1}\{y_n = +\} \left(\frac{1}{2} \pm \psi_1 \right) \times \\
&\quad \left\{ \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v}_+)| > 0\} \left[\sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \mathbb{1}\left\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} \quad b_{+,r}^{(t)} \right\} \left(\alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} \right) \right. \right. \\
&\quad \left. \left. + \sum_{p \notin \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \mathbb{1}\left\{ \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} \quad b_{+,r}^{(t)} \right\} \mathbf{x}_{n,p}^{(t)} \right] \right. \\
&\quad \left. + \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v}_+)| = 0\} \sum_{p \in [P]} \mathbb{1}\left\{ \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} \quad b_{+,r}^{(t)} \right\} \mathbf{x}_{n,p}^{(t)} \right\}
\end{aligned} \tag{39}$$

The rest of the proof proceeds by induction (in Phase 1).

First, recall that we set $b_{c,r}^{(0)} = -\frac{1}{4+2c_0} \sqrt{\log(d)}$, and $\Delta b_{c,r}^{(t)} = -\frac{\|w_{c,r}^{(t)}\|_2}{\log^5(d)}$ for all t in phase 1, and for any $+$ -class sample \mathbf{X}_n with $p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)$, $\alpha_{n,p}^{(t)} = \frac{1}{1 \pm \iota}$ by our data assumption.

Base case $t = 0$.

1. (On-diagonal common-feature neuron growth)

The base case for the neuron expression of point 1. is trivially true.

We show that the neurons $(+, r) \in S_+^{*(0)}(\mathbf{v}_+)$ only activate on \mathbf{v}_+ -dominated patches at time $t = 0$.

With probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, by Lemma H.3, we have for all possible choices of r, n, p :

$$\left| \mathbf{w}_{+,r}^{(0)}, \alpha_{n,p}^{(0)} \mathbf{v}_+ + \frac{(0)}{n,p} \right| = O(\sigma_0 \sigma_\zeta \sqrt{d \log(d)}) = O\left(\frac{\sigma_0}{\log^9(d)}\right) \tag{40}$$

It follows that

$$\begin{aligned}
&\mathbf{w}_{+,r}^{(0)}, \alpha_{n,p}^{(0)} \mathbf{v}_+ + \frac{(0)}{n,p} \\
&= \sigma_0 \left\{ \frac{1}{1 \pm \iota} \times \left(\frac{1}{4+2c_0} \sqrt{\log(d) + 1/\log^5(d)}, \frac{1}{10} \sqrt{\log(d)} \right) \pm \frac{1}{\log^9(d)} \right\} \\
&= \sigma_0 \left\{ \left(\frac{1}{1 - \iota} \frac{1}{4+2c_0} \sqrt{\log(d) + 1/\log^5(d)}, \frac{1}{1 + \iota} \frac{1}{10} \sqrt{\log(d)} \right) \pm \frac{1}{\log^9(d)} \right\}
\end{aligned} \tag{41}$$

Employing the basic identity $a - b = \frac{a^2 - b^2}{a + b}$, we have the lower bound

$$\begin{aligned}
& \sigma_0^{-1} \left(\mathbf{w}_{+,r}^{(0)}, \alpha_{n,p}^{(0)} \mathbf{v}_+ + \frac{(0)}{n,p} + b_{+,r}^{(0)} \right) \\
& \quad \sqrt{(1 - \iota)(4 + 2c_0)(\log(d) + 1/\log^5(d))} - \sqrt{(4 + 2c_0)\log(d)} - O\left(\frac{1}{\log^9(d)}\right) \\
& = \frac{(1 - \iota)(4 + 2c_0)(\log(d) + 1/\log^5(d)) - (4 + 2c_0)\log(d)}{\sqrt{(1 - \iota)(4 + 2c_0)(\log(d) + 1/\log^5(d))} + \sqrt{(4 + 2c_0)\log(d)}} - O\left(\frac{1}{\log^9(d)}\right) \\
& = \frac{(4 + 2c_0)(-\iota\log(d) + (1 - \iota)/\log^5(d))}{\sqrt{(1 - \iota)(4 + 2c_0)(\log(d) + 1/\log^5(d))} + \sqrt{(4 + 2c_0)\log(d)}} - O\left(\frac{1}{\log^9(d)}\right) \\
& > 0
\end{aligned} \tag{42}$$

The last inequality holds since $\iota = \frac{1}{\text{polylog}(d)}$ and d is sufficiently large such that $\frac{1}{\log^9(d)}$ does not drive the positive term down past 0.

Therefore, the neurons in $S_+^{*(0)}(\mathbf{v}_+)$ indeed activate on the \mathbf{v}_+ -dominated patches at $t = 0$.

The rest of the patches $\mathbf{x}_{n,p}^{(0)}$ is either a feature patch (not dominated by \mathbf{v}_+) or a noise patch. By definition, $(+, r) \in S_+^{*(0)}(\mathbf{v}_+) \iff (+, r) \in S_+^{(0)}(\mathbf{v}_+)$. Therefore, by Theorem F.1, with probability at least $1 - O\left(\frac{mk_+ NP}{\text{poly}(d)}\right)$, at time $t = 0$, the $(+, r) \in S_+^{*(0)}(\mathbf{v}_+)$ neurons we are considering cannot activate on any feature patch dominated by $\mathbf{v} \neq \mathbf{v}_+$, nor on any noise patches.

It follows that the expression equation 37 at time $t = 0$ is as follows:

$$\begin{aligned}
\text{equation 37} & = \sum_{n=1}^N \mathbb{1}\{y_n = +\} \left(\frac{1}{2} \pm \psi_1 \right) \times \\
& \quad \left\{ \mathbb{1}\{|P(\mathbf{X}_n^{(0)}; \mathbf{v}_+)| > 0\} \left[\sum_{p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+)} \left(\overline{1 \pm \iota} \mathbf{v}_+ + \frac{(0)}{n,p} \right) + \sum_{p \notin \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+)} 0 \right] \right. \\
& \quad \left. + \mathbb{1}\{|P(\mathbf{X}_n^{(0)}; \mathbf{v}_+)| = 0\} \sum_{p \in [P]} 0 \right\} \\
& = \left(\frac{1}{2} \pm \psi_1 \right) \sum_{n=1}^N \mathbb{1}\{y_n = +, |P(\mathbf{X}_n^{(0)}; \mathbf{v}_+)| > 0\} \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+)} \left(\overline{1 \pm \iota} \mathbf{v}_+ + \frac{(0)}{n,p} \right) \\
& = \left(\frac{1}{2} \pm \psi_1 \right) \times \\
& \quad \left| \left\{ (n, p) \in [N] \times [P] : y_n = +, |P(\mathbf{X}_n^{(0)}; \mathbf{v}_+)| > 0, p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+) \right\} \right| \left(\overline{1 \pm \iota} \mathbf{v}_+ \right) \\
& \quad + \sum_{n=1}^N \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+)} \mathbb{1}\{y_n = +\} \left(\frac{1}{2} \pm \psi_1 \right) \frac{(0)}{n,p}
\end{aligned} \tag{43}$$

On average,

$$\begin{aligned}
& \mathbb{E} \left[\left| \left\{ (n, p) \in [N] \times [P] : y_n = +, |P(\mathbf{X}_n^{(0)}; \mathbf{v}_+)| > 0, p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+) \right\} \right| \right] \\
& = \frac{s^*}{P} \times P \times \frac{N}{2} = \frac{s^* N}{2}
\end{aligned} \tag{44}$$

Furthermore, with our parameter choices, and by concentration of binomial random variables, with probability at least $1 - e^{-\text{polylog}(d)}$,

$$\left| \left\{ (n, p) \in [N] \times [P] : y_n = +, |P(\mathbf{X}_n^{(0)}; \mathbf{v}_+)| > 0, p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+) \right\} \right| = \frac{s^* N}{2} \left(1 \pm s^{*-1/3} \right) \tag{45}$$

must be true.

It follows that

$$\begin{aligned} \text{equation 37} &= \left(\frac{1}{2} \pm \psi_1\right) \times \frac{s^*N}{2} \left(1 \pm s^{*-1/2}\right) \times \left(\overline{1 \pm \iota} \mathbf{v}_+\right) \\ &+ \sum_{n=1}^N \sum_{p \in \mathcal{P}(X_n^{(0)}; \mathbf{v}_+)} \{y_n = +\} \left(\frac{1}{2} \pm \psi_1\right) \begin{matrix} (0) \\ n,p \end{matrix} \end{aligned} \quad (46)$$

The other component expression equation 38 is zero with probability at least $1 - O\left(\frac{mk_+ NP}{\text{poly}(d)}\right)$ by Theorem F.1.

By noting that

$$\begin{aligned} \text{Var}\left(\Delta_{+,r}^{(0)}\right) &= \text{Var}\left(\frac{\eta}{NP} \sum_{n=1}^N \sum_{p \in \mathcal{P}(X_n^{(0)}; \mathbf{v}_+)} \{y_n = +\} \left(\frac{1}{2} \pm \psi_1\right) \begin{matrix} (0) \\ n,p \end{matrix}\right) \\ &= \eta^2 \left(\frac{1}{2} \pm \psi_1\right)^2 \frac{s^*}{2NP^2} \left(1 \pm s^{*-1/3}\right) \sigma_\zeta^2, \end{aligned} \quad (47)$$

and

$$\mathbb{E}\left[\Delta_{+,r}^{(0)}\right] = \mathbb{E}\left[\frac{\eta}{NP} \sum_{n=1}^N \sum_{p \in \mathcal{P}(X_n^{(0)}; \mathbf{v}_+)} \{y_n = +\} \left(\frac{1}{2} \pm \psi_1\right) \begin{matrix} (0) \\ n,p \end{matrix}\right] = \mathbf{0}, \quad (48)$$

we finish the proof of the base case for point 1.

2. (On-diagonal finegrained-feature neuron growth)

The proof of the base case of point 2. is virtually identical to point 1, so we omit the computations here.

Inductive step: We condition on the high probability events of the induction hypothesis for $t \in [0, T]$ (with $T < T_0$ of course), and prove the statements for $t = T + 1$.

1. (On-diagonal common-feature neuron growth)

By the induction hypothesis, up to time $t = T$, with probability at least $1 - O\left(\frac{mk_+ NPT}{\text{poly}(d)}\right)$, for all $(+, r)$ $S_+^{*(T)}(\mathbf{v}_+)$,

$$\Delta \mathbf{w}_{+,r}^{(t)} = \eta \left(\left(\frac{1}{2} \pm \psi_1\right) \overline{1 \pm \iota} \left(1 \pm s^{*-1/3}\right) \right) \frac{s^*}{2P} \mathbf{v}_+ + \Delta_{+,r}^{(t)} \quad (49)$$

where $\Delta_{+,r}^{(t)} \sim \mathcal{N}(\mathbf{0}, \sigma_\zeta^{(t)2} \mathbf{I})$, $\sigma_\zeta^{(t)} = \eta \sigma_\zeta \left(\left(\frac{1}{2} \pm \psi_1\right) \overline{1 \pm s^{*-1/3}} \right) \frac{\sqrt{s}}{P\sqrt{2N}}$.

Expression of $\mathbf{w}_{+,r}^{(T+1)}$.

Conditioning on the high-probability event of the induction hypothesis, at time $t = T + 1$,

$$\begin{aligned} \mathbf{w}_{+,r}^{(T+1)} &= \mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^T \Delta \mathbf{w}_{+,r}^{(\tau)} \\ &= \eta T \left(\left(\frac{1}{2} \pm \psi_1\right) \overline{1 \pm \iota} \left(1 \pm s^{*-1/3}\right) \right) \frac{s^*}{2P} \mathbf{v}_+ + \begin{matrix} (t) \\ +,r \end{matrix} \end{aligned} \quad (50)$$

where $\begin{matrix} (t) \\ +,r \end{matrix} \sim \mathcal{N}(\mathbf{0}, \sigma_\zeta^{(t)2} \mathbf{I})$, $\sigma_\zeta^{(t)} = \eta \sigma_\zeta \overline{T} \left(\left(\frac{1}{2} \pm \psi_1\right) \overline{1 \pm s^{*-1/3}} \right) \frac{\sqrt{s}}{P\sqrt{2N}}$.

Let us compute $\Delta \mathbf{w}_{+,r}^{(T+1)}$.

We first want to show that $\mathbf{w}_{+,r}^{(T+1)}$ activates on \mathbf{v}_+ -dominated patches $\mathbf{x}_{n,p}^{(T+1)} = \overline{1 \pm \iota} \mathbf{v}_+ + \frac{(T+1)}{n,p}$. We need to show that the following expression is above 0:

$$\begin{aligned} & \mathbf{w}_{+,r}^{(T+1)}, \mathbf{x}_{n,p}^{(T+1)} + b_{+,r}^{(T+1)} \\ = & \mathbf{w}_{+,r}^{(0)}, \overline{1 \pm \iota} \mathbf{v}_+ + \frac{(T+1)}{n,p} + b_{+,r}^{(0)} \\ & + \left\langle \eta T \left(\left(\frac{1}{2} \pm \psi_1 \right) \overline{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \pm O \left(\frac{1}{\log^{10}(d)} \right) \right) \frac{s^*}{2P} \mathbf{v}_+ + \frac{(T+1)}{+,r}, \overline{1 \pm \iota} \mathbf{v}_+ + \frac{(T+1)}{n,p} \right\rangle \quad (51) \\ & + \sum_{\tau=0}^T \Delta b_{+,r}^{(\tau)} \end{aligned}$$

Let us treat the three terms (on three lines) separately.

First, following virtually the same argument as in the base case, the following lower bound holds with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$ for all n, p and $(+, r)$ $S_+^{*(T)}(\mathbf{v}_+)$:

$$\begin{aligned} & \mathbf{w}_{+,r}^{(0)}, \overline{1 \pm \iota} \mathbf{v}_+ + \frac{(T+1)}{n,p} + b_{+,r}^{(0)} \\ \sigma_0 \left\{ \sqrt{(1 - \iota)(4 + 2c_0)(\log(d) + 1/\log^5(d))} - \sqrt{(4 + 2c_0)\log(d)} - O\left(\frac{1}{\log^9(d)}\right) \right\} \quad (52) \\ & > 0 \end{aligned}$$

Now consider the second term.

We know, with probability at least $1 - e^{-\Omega(d)}$, for all n and p ,

$$\left| \frac{(T+1)}{n,p}, \mathbf{v}_+ \right| = O\left(\frac{1}{\log^{10}(d)}\right), \quad (53)$$

therefore,

$$\begin{aligned} & \left| \eta T \left(\left(\frac{1}{2} \pm \psi_1 \right) \overline{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \pm O\left(\frac{1}{\log^{10}(d)}\right) \right) \frac{s^*}{2P} \mathbf{v}_+, \frac{(T+1)}{n,p} \right| \quad (54) \\ & \eta T \frac{s^*}{2P} O\left(\frac{1}{\log^{10}(d)}\right). \end{aligned}$$

Moreover, with probability at least $1 - e^{-\Omega(d)}$,

$$\left| \frac{(T+1)}{+,r}, \mathbf{v}_+ \right| = \eta \frac{\bar{T}}{P} \frac{\bar{s}^*}{2N} \times O\left(\frac{1}{\log^{10}(d)}\right) \quad (55)$$

and with probability at least $1 - e^{-\Omega(d)}$,

$$\left| \frac{(T)}{+,r}, \frac{(T+1)}{n,p} \right| = O\left(\sigma_\zeta \sigma_\zeta^{(T)} d\right) = O\left(\eta \frac{\bar{T}}{P} \frac{\bar{s}^*}{2N} \frac{1}{\log^{20}(d)d}\right) = \eta \frac{\bar{T}}{P} \frac{\bar{s}^*}{2N} \frac{1}{\log^{19}(d)} \quad (56)$$

therefore

$$\eta T \frac{(T+1)}{+,r}, \overline{1 \pm \iota} \mathbf{v}_+ + \frac{(T+1)}{n,p} = \eta \frac{\bar{T}}{P} \frac{\bar{s}^*}{2N} O\left(\frac{1}{\log^{10}(d)}\right). \quad (57)$$

It follows that with probability at least $1 - O(e^{-\Omega(d)})$,

$$\begin{aligned}
& \left\langle \eta T \left(\left(\frac{1}{2} \pm \psi_1 \right) \overline{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \right) \frac{s^*}{2P} \mathbf{v}_{+,r} + \frac{(T+1)}{n,p}, \overline{1 \pm \iota} \mathbf{v}_{+,r} + \frac{(T+1)}{n,p} \right\rangle \\
&= \eta T \left(\left(\frac{1}{2} \pm \psi_1 \right) \overline{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \right) \frac{s^*}{2P} \mathbf{v}_{+,r}, \overline{1 \pm \iota} \mathbf{v}_{+,r} \\
&+ \eta T \left(\left(\frac{1}{2} \pm \psi_1 \right) \overline{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \right) \frac{s^*}{2P} \mathbf{v}_{+,r}, \frac{(T+1)}{n,p} \\
&+ \eta \frac{(T+1)}{n,p}, \overline{1 \pm \iota} \mathbf{v}_{+,r} + \frac{(T+1)}{n,p} \\
&\eta T \left(\frac{1}{2} - \psi_1^{(T+1)} \right) (1 - \iota) \left(1 - s^{*-1/3} \right) \frac{s^*}{2P} - \eta \frac{\overline{s^*}}{P} \frac{1}{2N} O \left(\frac{1}{\log^{10}(d)} \right).
\end{aligned} \tag{58}$$

Now we compute the third term. By the induction hypothesis,

$$\begin{aligned}
& \sum_{t=0}^T \Delta b_{+,r}^{(t)} \\
&= \sum_{t=0}^T \frac{\Delta \mathbf{w}_{+,r}^{(t)}}{\log^5(d)} \\
&= \sum_{t=0}^T \frac{1}{\log^5(d)} \left\| \eta \left(\frac{1}{2} \pm \psi_1 \right) \overline{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \frac{s^*}{2P} \mathbf{v}_{+,r} + \Delta_{+,r}^{(t)} \right\|_2 \\
& \sum_{t=0}^T \frac{1}{\log^5(d)} \eta \left(\frac{1}{2} + \psi_1 \right) \overline{1 + \iota} \left(1 + s^{*-1/3} \right) \frac{s^*}{2P} \mathbf{v}_{+,r} + \sum_{t=0}^T \frac{1}{\log^5(d)} \left\| \Delta_{+,r}^{(t)} \right\|_2 \\
&= \frac{1}{\log^5(d)} \eta T \left(\frac{1}{2} + \psi_1 \right) \overline{1 + \iota} \left(1 + s^{*-1/3} \right) \frac{s^*}{2P} + \sum_{t=0}^T \frac{1}{\log^5(d)} \left\| \Delta_{+,r}^{(t)} \right\|_2
\end{aligned} \tag{59}$$

With probability at least $1 - O\left(\frac{mT}{\text{poly}(d)}\right)$, for all $t \in [0, T]$ and r in consideration,

$$\left\| \Delta_{+,r}^{(t)} \right\|_2 \leq \eta \frac{\overline{s^*}}{P} \frac{1}{2N} O \left(\frac{1}{\log^{10}(d)} \right) \tag{60}$$

Therefore,

$$\begin{aligned}
& \sum_{t=0}^T \Delta b_{+,r}^{(t)} \\
& \frac{1}{\log^5(d)} \left(\eta T \left(\frac{1}{2} + \psi_1 \right) \overline{1 + \iota} \left(1 + s^{*-1/3} \right) \frac{s^*}{2P} + \eta T \frac{\overline{s^*}}{P} \frac{1}{2N} O \left(\frac{1}{\log^{10}(d)} \right) \right)
\end{aligned} \tag{61}$$

Combining our calculations of the three terms from above, we find the following estimate:

$$\begin{aligned}
& \mathbf{w}_{+,r}^{(T+1)}, \mathbf{x}_{n,p}^{(T+1)} + b_{+,r}^{(T+1)} \\
& > 0 \\
& + \eta T \left(\frac{1}{2} - \psi_1 \right) (1 - \iota) \left(1 - s^{*-1/3} \right) \frac{s^*}{2P} - \eta \frac{\bar{T}}{P} \frac{\bar{s}^*}{2N} O \left(\frac{1}{\log^{10}(d)} \right) \\
& - \frac{1}{\log^5(d)} \left(\eta T \left(\frac{1}{2} + \psi_1 \right) \frac{1}{1 + \iota} \left(1 + s^{*-1/3} \right) \frac{s^*}{2P} + \eta T \frac{\bar{s}^*}{P} \frac{1}{2N} O \left(\frac{1}{\log^{10}(d)} \right) \right) \\
& > \eta T \left(\left(\frac{1}{2} - \psi_1 \right) (1 - \iota) \left(1 - s^{*-1/3} \right) - O \left(\frac{1}{\log^4(d)} \right) \right) \frac{s^*}{2P} \\
& > 0
\end{aligned} \tag{62}$$

On the other hand, by Theorem F.1, with probability at least $1 - O \left(\frac{mk_+ NPT}{\text{poly}(d)} \right)$, none of the $(+, r)$ $S_+^{*(T)}(\mathbf{v}_+)$ can activate on $\mathbf{x}_{n,p}^{(T+1)}$ that are feature-patches dominated by \mathbf{v}_+ or noise patches.

Combining the above observations, with probability at least $1 - O \left(\frac{mk_+ NP(T+1)}{\text{poly}(d)} \right)$, the update expressions up to time $t = T + 1$ can be written as follows:

$$\begin{aligned}
\Delta \mathbf{w}_{+,r}^{(t)} &= \left(\frac{1}{2} \pm \psi_1 \right) \\
& \times \left\{ \left| \left\{ (n, p) \quad [N] \times [P] : y_n = +, |P(\mathbf{X}_n^{(t)}; \mathbf{v}_+)| > 0, p \quad P(\mathbf{X}_n^{(t)}; \mathbf{v}_+) \right\} \right| \left(\frac{1}{2} \pm \iota \mathbf{v}_+ \right) \right. \\
& \left. + \sum_{n=1}^N \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+)} \{y_n = +\} \left(\frac{1}{2} \pm \psi_1 \right) \frac{1}{n,p} \right\}
\end{aligned} \tag{63}$$

The rest of the derivations proceeds virtually the same as in the base case; we just need to rely on the concentration of binomial random variables to calculate

$$\left| \left\{ (n, p) \quad [N] \times [P] : y_n = +, |P(\mathbf{X}_n^{(0)}; \mathbf{v}_+)| > 0, p \quad P(\mathbf{X}_n^{(0)}; \mathbf{v}_+) \right\} \right| = \frac{s^* N}{2} \left(1 \pm s^{*-1/3} \right) \tag{64}$$

which completes the proof of the expression of $\Delta \mathbf{w}_{+,r}^{(t)}$.

Additionally, to show

$$\mathbf{w}_{+,r}^{(T+1)} - \mathbf{w}_{+,r}^{(0)} = \mathbf{w}_{+,r}^{(T+1)} - \mathbf{w}_{+,r}^{(0)} \tag{65}$$

we just need to note that, by the above sequence of derivations, for every $(+, r)$ $S_+^{*(0)}(\mathbf{v}_+)$, these neurons receive exactly the same update at time $t = T + 1$

$$\sum_{n=1}^N \mathbb{1}\{y_n = +\} \mathbb{1}\{|P(\mathbf{X}_n^{(T+1)}; \mathbf{v}_+)| > 0\} [1 - \text{logit}_+^{(T+1)}(\mathbf{X}_n^{(T+1)})] \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T+1)}; \mathbf{v}_+)} \left(\alpha_{n,p}^{(T+1)} \mathbf{v}_+ + \frac{1}{n,p} \right). \tag{66}$$

2. (On-diagonal finegrained-feature neuron growth)

For point 2, the proof strategy is almost identical, the only difference is that at every iteration, the expected number of patches in which subclass features appear in is

$$\begin{aligned}
& \left| \left\{ (n, p) \quad [N] \times ([P] - P(\mathbf{X}_n^{(T)}; \mathbf{v}_{+,c})) : y_n = +, |P(\mathbf{X}_n^{(T)}; \mathbf{v}_{+,c})| > 0, p \quad P(\mathbf{X}_n^{(T)}; \mathbf{v}_{+,c}) \right\} \right| \\
& = \frac{s^* N}{2k_+} \left(1 \pm s^{*-1/3} \right)
\end{aligned} \tag{67}$$

which holds with probability at least $1 - e^{- (\log^2(d))}$ for the relevant neurons. \square

Corollary D.1.1. $T_0 < O\left(\left(\eta \frac{s}{P}\right)^{-1}\right)$ $\text{poly}(d)$.

Proof. Follows from Theorem D.1. □

D.2 Lemmas

Lemma D.2. *During the time $t \in [0, T_0)$, for any $\mathbf{X}_n^{(t)}$,*

$$1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) = \frac{1}{2} \pm O(d^{-1}) \quad (68)$$

The same holds for $1 - \text{logit}_-^{(t)}(\mathbf{X}_n^{(t)})$.

Therefore, $|\psi_1| = O(d^{-1})$ for $t \in [0, T_0)$.

Proof. By definition of T_0 , for any $t \in [0, T_0]$, we have $F_c^{(t)}(\mathbf{X}_n^{(t)}) < d^{-1} + O(\eta)$ for all n , therefore, using Taylor approximation,

$$1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) = \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_+^{(t)}(\mathbf{X}_n^{(t)})) + \exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))} < \frac{\exp(d^{-1})}{1 + 1} = \frac{1}{2} + O(d^{-1}) \quad (69)$$

The lower bound can be proven due to convexity of the exponential:

$$\frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_+^{(t)}(\mathbf{X}_n^{(t)})) + \exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))} > \frac{1}{2} \exp(-d^{-1}) = \frac{1}{2} - \frac{1}{2d} \quad (70)$$

□

E Coarse-grained SGD Phase II: Loss Convergence, Large Neuron Movement

Recall that the desired probability events in Phase I happens with probability at least $1 - o(1)$.

In phase II, common-feature neurons start gaining large movement and drive the training loss down to $o(1)$. We show that the desired probability events occur with probability at least $1 - o(1)$.

We study the case of $T_1 \asymp \text{poly}(d)$, where T_1 denotes the time step at the end of training.

E.1 Main results

Theorem E.1. *With probability at least $1 - O\left(\frac{mk_+ NP T_1}{\text{poly}(d)}\right)$, the following events take place:*

1. *There exists time $T^* \asymp \text{poly}(d)$ such that for any $t \in [T^*, \text{poly}(d)]$, for any $n \in [N]$, the training loss $L(F; \mathbf{X}_n^{(t)}, y_n) = o(1)$.*

2. *(Easy sample test accuracy is nearly perfect) Given an easy test sample $(\mathbf{X}_{\text{easy}}, y)$, for $y' \in \{+1, -1\} - \{y\}$, for $t \in [T^*, \text{poly}(d)]$,*

$$\mathbb{P}\left[F_y^{(t)}(\mathbf{X}_{\text{easy}}) - F_{y'}^{(t)}(\mathbf{X}_{\text{easy}})\right] = o(1). \quad (71)$$

3. *(Hard sample test accuracy is bad) However, for all $t \in [0, \text{poly}(d)]$, given a hard test sample $(\mathbf{X}_{\text{hard}}, y)$,*

$$\mathbb{P}\left[F_y^{(t)}(\mathbf{X}_{\text{hard}}) - F_{y'}^{(t)}(\mathbf{X}_{\text{hard}})\right] = \Omega(1). \quad (72)$$

Proof. The training loss property follows from Lemma E.3 and Lemma E.4. We can set $T^* = T_{1,1}$ or any time beyond it (and upper bounded by $\text{poly}(d)$).

The test accuracy properties follow from Lemma E.8 and Lemma E.9. □

E.2 Lemmas

Lemma E.2 (Phase II, Update Expressions). *For any $T_1 \asymp \text{poly}(d)$, with probability at least $1 - O\left(\frac{mNPk_+ t}{\text{poly}(d)}\right)$, during $t \in [T_0, T_1]$, for any $(+, r) \in S_+^{*(0)}(\mathbf{v}_+)$,*

$$\begin{aligned} & \Delta w_{+,r}^{(t)} \\ &= \eta \sum_{n=1}^N \mathbb{1}\{y_n = +\} \exp\left\{-F_+^{(t)}(\mathbf{X}_n^{(t)})\right\} \\ & \quad \times \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp\left(F_-^{(t)}(\mathbf{X}_n^{(t)}) - F_+^{(t)}(\mathbf{X}_n^{(t)})\right) + 1} (1 \pm s^{*-1/3}) \frac{s^*}{NP} \left(\overline{1 \pm \iota} \mathbf{v}_+ + \binom{(t)}{n,p} \right), \end{aligned} \quad (73)$$

(where c_n^t denotes the subclass index of sample $\mathbf{X}_n^{(t)}$) and for any $(+, r) \in S_+^{*(0)}(\mathbf{v}_{+,c})$,

$$\begin{aligned} & \Delta w_{+,r}^{(t)} \\ &= \eta \exp\left\{- (1 \pm s^{*-1/3}) \overline{1 \pm \iota} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right)\right) s^* \left(A_{+,r}^{*(t)} \left|S_+^{*(0)}(\mathbf{v}_+)\right| + A_{+,c,r}^{*(t)} \left|S_+^{*(0)}(\mathbf{v}_{+,c})\right|\right)\right\} \\ & \quad \times \sum_{n=1}^N \mathbb{1}\{y_n = (+, c)\} \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp\left(F_-^{(t)}(\mathbf{X}_n^{(t)}) - F_+^{(t)}(\mathbf{X}_n^{(t)})\right) + 1} (1 \pm s^{*-1/3}) \frac{s^*}{NP} \left(\overline{1 \pm \iota} \mathbf{v}_{+,c} + \binom{(t)}{n,p} \right), \end{aligned} \quad (74)$$

In fact, for any $\mathbf{v} \in \{\mathbf{v}_+\} = \{\mathbf{v}_{+,c}\}_{c=1}^{k_+}$, every neuron in $S_+^{*(0)}(\mathbf{v})$ remain activated (on \mathbf{v} -dominated patches) and receive exactly the same updates at every iteration as shown above.

For simpler exposition, for any $(+, r^*) \in S_+^{*(0)}(\mathbf{v}_+)$, we write $A_{+,r}^{*(t)} := \mathbf{w}_{+,r}^{(t)}, \mathbf{v}_+$; similarly for $A_{+,c,r}^{*(t)} := \mathbf{w}_{+,r}, \mathbf{v}_{+,c}$ for neurons $(+, r^*) \in S_+^{*(0)}(\mathbf{v}_{+,c})$.

Moreover, on “+”-class samples, the neural network response satisfies the estimate for every $(+, r^*) \in S_+^{*(0)}(\mathbf{v}_+)$:

$$\begin{aligned} & F_+^{(t)}(\mathbf{X}_n^{(t)}) \\ &= (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) \times s^* \left(A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_{+,c_n^t}) \right| \right), \end{aligned} \quad (75)$$

The same claims hold for the “-” class neurons (with the class signs flipped).

Proof. In this proof we focus on the neurons in $S_+^{*(0)}(\mathbf{v}_+)$; the proof for the update expressions for those in $S_+^{*(0)}(\mathbf{v}_{+,c})$ are proven in virtually the same way.

Base case, $t = T_0$.

First define $A_{+,r}^{*(t)} := \mathbf{w}_{+,r}, \mathbf{v}_+, (+, r^*) \in S_+^{*(0)}(\mathbf{v}_+)$; similarly for $A_{+,c,r}^{*(t)} := \mathbf{w}_{+,r}, \mathbf{v}_{+,c}$. Note that the choice of r^* does not really matter, since we know from phase I that every neuron in $S_+^{*(0)}(\mathbf{v}_+)$ evolve at exactly the same rate, so by the end of phase I, $\mathbf{w}_{+,r}^{(T_0)} - \mathbf{w}_{+,r}^{(T_0)'} = O(\sigma_0 \log(d)) = \mathbf{w}_{+,r}^{(T_0)'} - \mathbf{w}_{+,r}^{(T_0)}$ for any $(+, r), (+, r') \in S_+^{*(0)}(\mathbf{v}_+)$.

Let $(+, r) \in S_+^{*(0)}(\mathbf{v}_+)$. Similar to phase I, consider the update equation

$$\mathbf{w}_{+,r}^{(t+1)} = \mathbf{w}_{+,r}^{(t)} + \eta \frac{1}{NP} \times \quad (76)$$

$$\sum_{n=1}^N \left(\mathbb{1}\{y_n = +\} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \right) \quad (77)$$

$$+ \mathbb{1}\{y_n = -\} [-\text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \quad (78)$$

For the on-diagonal update expression, we have

$$\begin{aligned} & \sum_{n=1}^N \mathbb{1}\{y_n = +\} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \\ &= \sum_{n=1}^N \mathbb{1}\{y_n = +\} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\ & \quad \left\{ \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}_+) > 0\} \left[\sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \mathbb{1}\left\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{1}{n,p} b_{+,r}^{(t)} \right\} \left(\alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{1}{n,p} \right) \right. \right. \\ & \quad \left. \left. + \sum_{p \notin \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \mathbb{1}\left\{ \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)} \right\} \mathbf{x}_{n,p}^{(t)} \right] \right. \\ & \quad \left. + \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}_+) = 0\} \sum_{p \in [P]} \mathbb{1}\left\{ \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)} \right\} \mathbf{x}_{n,p}^{(t)} \right\} \end{aligned} \quad (79)$$

Following from Theorem D.1 and F.1, the neurons' non-activation on the patches that do not contain \mathbf{v}_+ , and activation on the \mathbf{v}_+ -dominated patches hold with probability at least $1 - O\left(\frac{mNPk_+}{\text{poly}(d)}\right)$ at time T_0 . Therefore, the above update expression reduces to

$$\sum_{n=1}^N \mathbb{1}\{y_n = +, |P(\mathbf{X}_n^{(t)}; \mathbf{v}_+)| > 0\} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \left(\alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} \right) \quad (80)$$

Note that for samples $\mathbf{X}_n^{(t)}$ with $y_n = +$,

$$[1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] = \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)})) + \exp(F_+^{(t)}(\mathbf{X}_n^{(t)}))} \quad (81)$$

Now we need to estimate the network response $F_+^{(t)}(\mathbf{X}_n^{(t)})$. With probability at least $1 - \exp(-\Omega(s^{*1/3}))$, we have the upper bound (let $(+, c_n^t)$ denote the subclass which sample $\mathbf{X}_n^{(t)}$ belongs to):

$$\begin{aligned} & F_+^{(t)}(\mathbf{X}_n^{(t)}) \\ & \sum_{p \in \mathcal{P}(\mathbf{X}^{(t)}; \mathbf{v}_+)} \sum_{(+, r) \in S_+^{(0)}(\mathbf{v}_+)} \mathbf{w}_{+,r}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} + b_{+,r}^{(t)} \\ & + \sum_{p \in \mathcal{P}(\mathbf{X}^{(t)}; \mathbf{v}_+, c_n^t)} \sum_{(+, r) \in S_+^{(0)}(\mathbf{v}_+, c_n^t)} \mathbf{w}_{+,r}^{(t)} \mathbf{v}_+, c_n^t + \frac{(t)}{n,p} + b_{+,r}^{(t)} \\ & (1 + s^{*-1/3}) \frac{1}{1 + \iota s^*} \left(1 + O\left(\frac{1}{\log^9(d)}\right) \right) \left(A_{+,r}^{*(t)} \left| S_+^{(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(t)} \left| S_+^{(0)}(\mathbf{v}_+, c_n^t) \right| \right) \end{aligned} \quad (82)$$

The second inequality is true since $\max_r \mathbf{w}_{+,r}^{(t)} \mathbf{v}_+ + A_{+,r}^{*(t)} + O(\sigma_0 \log(d))$, and for any $(+, r) \in S_+^{(0)}(\mathbf{v}_+)$, $|\mathbf{w}_{+,r}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p}| \leq O(1/\log^9(d)) A_{+,r}^{*(t)}$. The bias value is negative (and so less than 0).

To further refine the bound, we recall $\left| S_+^{*(0)}(\mathbf{v}) \right| / \left| S_+^{(0)}(\mathbf{v}) \right|, \left| S_+^{*(0)}(\mathbf{v}) \right| / \left| S_+^{(0)}(\mathbf{v}) \right| = 1 \pm O(1/\log^5(d))$.

Therefore, we obtain the bound

$$\begin{aligned} F_+^{(t)}(\mathbf{X}_n^{(t)}) & (1 + s^{*-1/3}) \frac{1}{1 + \iota s^*} \left(1 + O\left(\frac{1}{\log^5(d)}\right) \right) \left(1 + O\left(\frac{1}{\log^5(d)}\right) \right) \\ & \times \left(A_{+,r}^{*(t)} \left| S_+^{(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(t)} \left| S_+^{(0)}(\mathbf{v}_+, c_n^t) \right| \right) \end{aligned} \quad (83)$$

Following a similar argument, we also have the lower bound

$$\begin{aligned} & F_+^{(t)}(\mathbf{X}_n^{(t)}) \\ & \sum_{p \in \mathcal{P}(\mathbf{X}^{(t)}; \mathbf{v}_+)} \sum_{(+, r) \in S_+^{(0)}(\mathbf{v}_+)} \sigma \left(\mathbf{w}_{+,r}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} + b_{+,r}^{(t)} \right) \\ & + \sum_{p \in \mathcal{P}(\mathbf{X}^{(t)}; \mathbf{v}_+, c_n^t)} \sum_{(+, r) \in S_+^{(0)}(\mathbf{v}_+, c_n^t)} \sigma \left(\mathbf{w}_{+,r}^{(t)} \mathbf{v}_+, c_n^t + \frac{(t)}{n,p} + b_{+,r}^{(t)} \right) \\ & (1 - s^{*-1/3}) \frac{1}{1 - \iota s^*} \left(1 - O\left(\frac{1}{\log^5(d)}\right) \right) \left(1 - O\left(\frac{1}{\log^5(d)}\right) \right) \\ & \times \left(A_{+,r}^{*(t)} \left| S_+^{(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(t)} \left| S_+^{(0)}(\mathbf{v}_+, c_n^t) \right| \right) \end{aligned} \quad (84)$$

The neurons in $S_+^{*(0)}(\mathbf{v}_+)$ have to activate, therefore they serve a key role in the lower bound, the bias bound for them is simply $-A_{+,r}^{*(t)} \Theta(1/\log^5(d))$; the neurons in $S_+^{(0)}(\mathbf{v}_{+,c})$ contribute at least 0 due to the ReLU activation; the rest of the neurons do not activate. The same reasoning holds for the $S_+^{*(0)}(\mathbf{v}_{+,c})$.

Knowing that neurons in $S_+^{*(0)}(\mathbf{v}_+)$ cannot activate on the patches in samples belonging to the “-” class, now we may write the update expression for every $(+, r)$ $S_+^{*(t)}(\mathbf{v}_+)$ as (their updates are identical, same as in phase I):

$$\begin{aligned}
& \Delta w_{+,r}^{(t)} \\
&= \frac{\eta}{NP} \sum_{n=1}^N \mathbb{1}\{y_n = +\} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(w_{+,r}^{(t)} \mathbf{x}_{n,p}^{(t)} + b_{+,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \\
&= \frac{\eta}{NP} \sum_{n=1}^N \mathbb{1}\{y_n = +, |P(\mathbf{X}_n^{(t)}; \mathbf{v}_+)| > 0\} \exp(-F_+^{(t)}(\mathbf{X}_n^{(t)})) \\
&\quad \times \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}) - \exp(F_+^{(t)}(\mathbf{X}_n^{(t)}))) + 1} \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \left(\alpha_{n,p}^{(t)} \mathbf{v}_+ + \beta_{n,p}^{(t)} \right) \\
&= \eta \sum_{n=1}^N \mathbb{1}\{y_n = +\} \exp \left\{ - (1 + s^{*-1/3}) \frac{1}{1 + \iota s^*} \left(1 + O\left(\frac{1}{\log^5(d)}\right) \right) \right. \\
&\quad \times \left. \left(A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_{+,c_n^t}) \right| \right) \right\} \\
&\quad \times \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}) - F_+^{(t)}(\mathbf{X}_n^{(t)})) + 1} (1 \pm s^{*-1/3}) \frac{s^*}{NP} \left(\frac{1}{1 \pm \iota} \mathbf{v}_+ + \beta_{n,p}^{(t)} \right)
\end{aligned} \tag{85}$$

This concludes the proof of the base case.

Induction step. Assume the statements hold for time period $[T_0, t]$, prove for time $t + 1$.

At step $t + 1$, based on the induction hypothesis, we know that with probability at least $1 - O\left(\frac{mNPk_+t}{\text{poly}(d)}\right)$, during time $\tau \in [T_0, t]$, for any $(+, r)$ $S_+^{*(0)}(\mathbf{v}_+)$,

$$\begin{aligned}
& \Delta w_{+,r}^{(\tau)} \\
&= \eta \sum_{n=1}^N \mathbb{1}\{y_n = +\} \exp \left\{ - (1 + s^{*-1/3}) \frac{1}{1 + \iota s^*} \left(1 + O\left(\frac{1}{\log^5(d)}\right) \right) \right. \\
&\quad \times \left. \left(A_{+,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_{+,c_n^t}) \right| \right) \right\} \\
&\quad \times \frac{\exp(F_-^{(\tau)}(\mathbf{X}_n^{(\tau)}))}{\exp(F_-^{(\tau)}(\mathbf{X}_n^{(\tau)}) - \exp(F_+^{(\tau)}(\mathbf{X}_n^{(\tau)}))) + 1} (1 \pm s^{*-1/3}) \frac{s^*}{NP} \left(\frac{1}{1 \pm \iota} \mathbf{v}_+ + \beta_{n,p}^{(\tau)} \right)
\end{aligned} \tag{86}$$

and for the bias,

$$\begin{aligned}
& \Delta b_{+,r}^{(\tau)} \\
& - \eta \frac{1}{\log^5(d)} \sum_{n=1}^N \mathbb{1}\{y_n = +\} \exp \left\{ - (1 + s^{*-1/3}) \frac{1}{1 + \iota s^*} \left(1 + O \left(\frac{1}{\log^5(d)} \right) \right) \right. \\
& \times \left. \left(A_{+,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_{+,c_n^t}) \right| \right) \right\} \\
& \times (1 - s^{*-1/3}) \frac{s^*}{NP} \left(\frac{1}{1 - \iota} - \frac{1}{\log^{10}(d)} \right) \frac{\exp(F_-^{(\tau)}(\mathbf{X}_n^{(\tau)}))}{\exp(F_-^{(\tau)}(\mathbf{X}_n^{(\tau)}) - \exp(F_+^{(\tau)}(\mathbf{X}_n^{(\tau)})) + 1}
\end{aligned} \tag{87}$$

Conditioning on the high-probability events of the induction hypothesis,

$$\begin{aligned}
& \mathbf{w}_{+,r}^{(t+1)} \\
& = \mathbf{w}_{+,r}^{(T_0)} \\
& + \eta \sum_{\tau=T_0}^t \sum_{n=1}^N \mathbb{1}\{y_n = +\} \exp \left\{ - (1 + s^{*-1/3}) \frac{1}{1 + \iota s^*} \left(1 + O \left(\frac{1}{\log^5(d)} \right) \right) \right. \\
& \times \left. \left(A_{+,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_{+,c_n^t}) \right| \right) \right\} \\
& \times \frac{\exp(F_-^{(\tau)}(\mathbf{X}_n^{(\tau)}))}{\exp(F_-^{(\tau)}(\mathbf{X}_n^{(\tau)}) - \exp(F_+^{(\tau)}(\mathbf{X}_n^{(\tau)})) + 1} (1 \pm s^{*-1/3}) \frac{s^*}{NP} \left(\frac{1}{1 \pm \iota} \mathbf{v}_+ + \frac{(\tau)}{n,p} \right)
\end{aligned} \tag{88}$$

It follows that, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, for all \mathbf{v}_+ -dominated patch $\mathbf{x}_{n,p}^{(t+1)}$,

$$\begin{aligned}
& \mathbf{w}_{+,r}^{(t+1)}, \mathbf{x}_{n,p}^{(t+1)} + b_{+,r}^{(t+1)} \\
= & \mathbf{w}_{+,r}^{(T_0)}, \overline{1 \pm \iota} \mathbf{v}_+ + \frac{(t+1)}{n,p} + b_{+,r}^{(T_0)} \\
& + \eta \sum_{\tau=T_0}^t \sum_{n=1}^N \mathbb{1}\{y_n = +\} \exp \left\{ - (1 + s^{*-1/3}) \overline{1 + \iota} s^* \left(1 + O \left(\frac{1}{\log^5(d)} \right) \right) \right. \\
& \times \left. \left(A_{+,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_+, c_n^t) \right| \right) \right\} \\
& \times \frac{\exp(F_-^{(\tau)}(\mathbf{x}_n^{(\tau)}))}{\exp(F_-^{(\tau)}(\mathbf{x}_n^{(\tau)}) - F_+^{(\tau)}(\mathbf{x}_n^{(\tau)})) + 1} (1 \pm s^{*-1/3}) \frac{s^*}{NP} \\
& \times \overline{1 \pm \iota} \mathbf{v}_+ + \frac{(\tau)}{n,p}, \overline{1 \pm \iota} \mathbf{v}_+ + \frac{(t+1)}{n,p} + \Delta b_{+,r}^{(\tau)} \\
& 0 \\
& + \eta \sum_{\tau=T_0}^t \sum_{n=1}^N \mathbb{1}\{y_n = +\} \exp \left\{ - (1 + s^{*-1/3}) \overline{1 + \iota} s^* \left(1 + O \left(\frac{1}{\log^5(d)} \right) \right) \right. \\
& \times \left. \left(A_{+,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(\tau)} \left| S_+^{*(0)}(\mathbf{v}_+, c_n^t) \right| \right) \right\} \\
& \times \frac{\exp(F_-^{(\tau)}(\mathbf{x}_n^{(\tau)}))}{\exp(F_-^{(\tau)}(\mathbf{x}_n^{(\tau)}) - F_+^{(\tau)}(\mathbf{x}_n^{(\tau)})) + 1} (1 \pm s^{*-1/3}) \frac{s^*}{NP} \\
& \times \left(1 - \iota - O \left(\frac{1}{\log^5(d)} \right) \right) \\
& > 0
\end{aligned} \tag{89}$$

Therefore the neurons $(+, r)$ $S_+^{*(0)}(\mathbf{v}_+)$ activate on the \mathbf{v}_+ -dominated patches $\mathbf{x}_{n,p}^{(t+1)}$. We also know that they cannot activate on patches that are not dominated by \mathbf{v}_+ by Theorem F.1. Following a similar derivation to the base case, we arrive at the result that, conditioning on the events of the induction hypothesis, with probability at least $1 - O\left(\frac{mNPk_+}{\text{poly}(d)}\right)$, for all $(+, r)$ $S_+^{*(0)}(\mathbf{v}_+)$,

$$\begin{aligned}
& \Delta \mathbf{w}_{+,r}^{(t+1)} \\
= & \eta \sum_{n=1}^N \mathbb{1}\{y_n = +\} \exp \left\{ - (1 + s^{*-1/3}) \overline{1 + \iota} s^* \left(1 + O \left(\frac{1}{\log^5(d)} \right) \right) \right. \\
& \times \left. \left(A_{+,r}^{*(t+1)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(t+1)} \left| S_+^{*(0)}(\mathbf{v}_+, c_n^t) \right| \right) \right\} \\
& \times (1 \pm s^{*-1/3}) \frac{s^*}{NP} \frac{\exp(F_-^{(t+1)}(\mathbf{x}_n^{(t+1)}))}{\exp(F_-^{(t+1)}(\mathbf{x}_n^{(t+1)}) - F_+^{(t+1)}(\mathbf{x}_n^{(t+1)})) + 1} \left(\overline{1 \pm \iota} \mathbf{v}_+ + \frac{(t+1)}{n,p} \right)
\end{aligned} \tag{90}$$

Consequently, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$,

$$\begin{aligned}
& \Delta b_{+,r}^{(t+1)} \\
& - \frac{1}{\log^5(d)} \sum_{n=1}^N \mathbb{1}\{y_n = +\} \eta \exp \left\{ - (1 + s^{*-1/3}) \frac{1}{1 + \iota s^*} \left(1 + O\left(\frac{1}{\log^5(d)}\right) \right) \right. \\
& \left. \times \left(A_{+,r}^{*(t+1)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(t+1)} \left| S_+^{*(0)}(\mathbf{v}_{+,c_n^t}) \right| \right) \right\} \\
& \times \frac{\exp(F_-^{(t+1)}(\mathbf{X}_n^{(t+1)}))}{\exp(F_-^{(t+1)}(\mathbf{X}_n^{(t+1)}) - F_+^{(t+1)}(\mathbf{X}_n^{(t+1)})) + 1} (1 - s^{*-1/3}) \frac{s^*}{NP} \\
& \times \left(1 - \iota - O\left(\frac{1}{\log^9(d)}\right) \right)
\end{aligned} \tag{91}$$

Utilizing the definition of conditional probability, we conclude that the expressions for $\Delta \mathbf{w}_{+,r}^{(\tau)}$ and $\Delta b_{+,r}^{(t+1)}$ are indeed as described in the theorem during time $\tau \in [T_0, t+1]$ with probability at least $\left(1 - O\left(\frac{mNPk_+ t}{\text{poly}(d)}\right)\right) \times \left(1 - O\left(\frac{mNPk_+}{\text{poly}(d)}\right)\right) = 1 - O\left(\frac{mNPk_+(t+1)}{\text{poly}(d)}\right)$.

Moreover, based on the expression of $\Delta \mathbf{w}_{+,r}^{(\tau)}$ and $\Delta b_{+,r}^{(t+1)}$, following virtually the same argument as in the base case, we can estimate the network output for any $(\mathbf{X}_n^{(t+1)}, y_n = +)$:

$$\begin{aligned}
F_+^{(t+1)}(\mathbf{X}_n^{(t+1)}) &= (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) s^* \\
& \times \left(A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c_n^t,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_{+,c_n^t}) \right| \right)
\end{aligned} \tag{92}$$

□

Lemma E.3. Define time $T_{1,1}$ to be the first point in time which the following identity holds on all $\mathbf{X}_n^{(t)}$ belonging to the “+” class:

$$\frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}) - F_+^{(t)}(\mathbf{X}_n^{(t)})) + 1} = 1 - O\left(\frac{1}{\log^5(d)}\right) \tag{93}$$

Then $T_{1,1} = \text{poly}(d)$, and for all $t \in [T_{1,1}, T_1]$, the above holds. The following also holds for this time period:

$$[1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] = O\left(\frac{1}{\log^5(d)}\right) \tag{94}$$

The same results also hold with the class signs flipped.

Proof. We first note that, the training loss $[1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})]$ on samples belonging to the “+” class at any time during $t \in [T_0, T_1]$ is, asymptotically speaking, monotonically decreasing from $\frac{1}{2} - O(d^{-1})$. This can be easily proven by observing the way $s^* \left(A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_{+,c}) \right| \right)$ monotonically increases from the proof of Lemma E.2: before $F_+^{(t)}(\mathbf{X}_n^{(t)}) = \log \log^5(d)$ on all $\mathbf{X}_n^{(t)}$ belonging to the “+” class, there

must be some samples $\mathbf{X}_n^{(t)}$ on which

$$\begin{aligned} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] &= \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)})) + \exp(F_+^{(t)}(\mathbf{X}_n^{(t)}))} \\ &= \frac{1 - O(\sigma_0 \log(d) s^* d^{c_0})}{1 + O(\sigma_0 \log(d) s^* d^{c_0}) + \log^5(d)} \\ &= \Omega\left(\frac{1}{\log^5(d)}\right). \end{aligned} \quad (95)$$

Therefore, by the update expressions in the proof of Lemma E.2, $F_+^{(t)}(\mathbf{X}_n^{(t)})$ can reach $\log \log^5(d)$ in time at most $O\left(\frac{NP \log^5(d)}{\eta s}\right)$ poly(d) (in the worst case scenario). At time $T_{1,1}$ and beyond,

$$\begin{aligned} 1 - \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}) - F_+^{(t)}(\mathbf{X}_n^{(t)})) + 1} &= 1 - \frac{\exp(1 - O(\sigma_0 d^{c_0} s^*))}{\exp(1 + O(\sigma_0 d^{c_0} s^*)) \frac{1}{\log^5(d)} + 1} \\ &= O\left(\frac{1}{\log^5(d)}\right). \end{aligned} \quad (96)$$

□

Lemma E.4. Denote $C = \eta \frac{s}{2k_+ P}$, and write (for any $c \in [k_+]$)

$$A_c(t) = s^* \left(A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_{+,c}) \right| \right) \quad (97)$$

(see Lemma E.2 for definition of $A^{*(t)}$). Define $t_{c,0} = \exp(A_c(T_{1,1}))$. We write $A(t)$ and t_0 below for cleaner notations.

Then with probability at least $1 - o(1)$, during $t \in [T_{1,1}, T_1]$,

$$A(t) = \log(C(t - T_{1,1}) + t_0) + E(t) \quad (98)$$

where $|E(t)| = O\left(\frac{1}{\log^4(d)}\right) \sum_{\tau=C^{-1}t_0}^{t-T_{1,1}+C^{-1}t_0} \frac{1}{\tau} = O\left(\frac{\log(t) - \log(C^{-1}t_0)}{\log^4(d)}\right)$.

The same results also hold with the class signs flipped.

Proof. Sidenote: To make the writing a bit cleaner, we assume in the proof below that $C^{-1}t_0$ is an integer. The general case is easy to extend to by observing that $\left| \frac{1}{t-T_{1,1}+\lceil C^{-1}t_0 \rceil} - \frac{1}{t-T_{1,1}+C^{-1}t_0} \right| \leq \frac{1}{(t-T_{1,1}+\lceil C^{-1}t_0 \rceil)(t-T_{1,1}+C^{-1}t_0)}$, which can be absorbed into the error term at every iteration since $\frac{1}{t-T_{1,1}+\lceil C^{-1}t_0 \rceil} \geq \frac{1}{\log^4(d)}$ due to $C^{-1}t_0 = \Omega(\sigma_0^{-1}/(\text{polylog}(d)d^{c_0})) = d \log^4(d)$.

Based on result from Lemmas E.2 and E.3, as long as $A(t) = O(\log(d))$, we know during time $t \in [T_{1,1}, T_1]$ the update rule for $A(t)$ is as follows:

$$\begin{aligned} A(t+1) - A(t) &= C \exp \left\{ -(1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) A(t) \right\} \\ &\quad \times \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) (1 \pm s^{*-1/3}) \left(\frac{1}{1 \pm \iota} \pm \frac{1}{\log^{10}(d)} \right) \\ &= C \exp \{-A(t)\} \exp \left\{ \pm O\left(\frac{1}{\log^4(d)}\right) \right\} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) \\ &= C \exp \{-A(t)\} \left(1 \pm \frac{C_1}{\log^4(d)} \right) \end{aligned} \quad (99)$$

where we write C_1 in place of $O(\cdot)$ for a more concrete update expression.

The base case $t = T_{1,1}$ is trivially true.

We proceed with the induction step. Assume the hypothesis true for $t \in [T_{1,1}, T]$, prove for $t + 1 = T + 1$.

Note that by Lemma E.10,

$$\begin{aligned}
A(t+1) &= \log(C(t - T_{1,1}) + t_0) + E(t) \\
&\quad + C \exp \{ -\log(C(t - T_{1,1}) + t_0) - E(t) \} \left(1 \pm \frac{C_1}{\log^4(d)} \right) \\
&= \log(C) + \log(t - T_{1,1} + C^{-1}t_0) + E(t) \\
&\quad + C \frac{1}{C(t - T_{1,1}) + t_0} (1 - E(t) \pm O(E(t)^2)) \left(1 \pm \frac{C_1}{\log^4(d)} \right) \\
&= \log(C) + \sum_{\tau=1}^{t-T_{1,1}+C^{-1}t_0-1} \frac{1}{\tau} + \frac{1}{2} \frac{1}{t - T_{1,1} + C^{-1}t_0} + \left[0, \frac{1}{8} \frac{1}{(t - T_{1,1} + C^{-1}t_0)^2} \right] \\
&\quad + \frac{1}{t - T_{1,1} + C^{-1}t_0} \pm \frac{C_1}{\log^4(d)} \frac{1}{t - T_{1,1} + C^{-1}t_0} \\
&\quad + E(t) + \frac{1}{t - T_{1,1} + C^{-1}t_0} (-E(t) \pm O(E(t)^2)) \left(1 \pm \frac{C_1}{\log^4(d)} \right) \\
&= \log(C) + \sum_{\tau=1}^{t-T_{1,1}+C^{-1}t_0} \frac{1}{\tau} + \frac{1}{2} \frac{1}{t - T_{1,1} + C^{-1}t_0} + \left[0, \frac{1}{8} \frac{1}{(t - T_{1,1} + C^{-1}t_0)^2} \right] \\
&\quad \pm \frac{C_1}{\log^4(d)} \frac{1}{t - T_{1,1} + C^{-1}t_0} \\
&\quad + E(t) + \frac{1}{t - T_{1,1} + C^{-1}t_0} (-E(t) \pm O(E(t)^2)) \left(1 \pm \frac{C_1}{\log^4(d)} \right)
\end{aligned} \tag{100}$$

Invoking Lemma E.10 again,

$$\begin{aligned}
A(t+1) &= \log(C) + \log(t+1 - T_{1,1} + C^{-1}t_0) \\
&\quad - \frac{1}{2} \frac{1}{t+1 - T_{1,1} + C^{-1}t_0} + \frac{1}{2} \frac{1}{t - T_{1,1} + C^{-1}t_0} \\
&\quad + \left[-\frac{1}{8} \frac{1}{(t+1 - T_{1,1} + C^{-1}t_0)^2}, 0 \right] + \left[0, \frac{1}{8} \frac{1}{(t - T_{1,1} + C^{-1}t_0)^2} \right] \\
&\quad \pm \frac{C_1}{\log^4(d)} \frac{1}{t - T_{1,1} + C^{-1}t_0} \\
&\quad + E(t) + \frac{1}{t - T_{1,1} + C^{-1}t_0} (-E(t) \pm O(E(t)^2)) \left(1 \pm \frac{C_1}{\log^4(d)} \right) \\
&= \log(C(t+1 - T_{1,1}) + t_0) \\
&\quad + \frac{1}{2} \frac{1}{(t+1 - T_{1,1} + C^{-1}t_0)(t - T_{1,1} + C^{-1}t_0)} \pm O \left(\frac{1}{(t+1 - T_{1,1} + C^{-1}t_0)^2} \right) \\
&\quad \pm \frac{C_1}{\log^4(d)} \frac{1}{t - T_{1,1} + C^{-1}t_0} \\
&\quad + E(t) + \frac{1}{t - T_{1,1} + C^{-1}t_0} (-E(t) \pm O(E(t)^2)) \left(1 \pm \frac{C_1}{\log^4(d)} \right)
\end{aligned} \tag{101}$$

To further refine the expression, first note that the error passed down from the previous step t does not grow in this step (in fact it slightly decreases):

$$\begin{aligned} & \left| E(t) + \frac{1}{t - T_{1,1} + C^{-1}t_0} (-E(t) \pm O(E(t)^2)) \left(1 \pm \frac{C_1}{\log^4(d)} \right) \right| \\ & < |E(t)| \\ & O\left(\frac{1}{\log^4(d)}\right) \sum_{\tau=C^{-1}t_0}^{t-T_{1,1}+C^{-1}t_0} \frac{1}{\tau}. \end{aligned} \quad (102)$$

Moreover, notice that at step $t+1$, since $\frac{1}{t+1-T_{1,1}+C^{-1}t_0} = \frac{1}{\log^4(d)}$, the error term $|E(t+1)| = |A(t+1) - \log(C(t+1 - T_{1,1}) + t_0)| = O\left(\frac{1}{\log^4(d)}\right) \sum_{\tau=C^{-1}t_0}^{t+1-T_{1,1}+C^{-1}t_0} \frac{1}{\tau}$, which finishes the inductive step. \square

Lemma E.5. *With probability at least $1 - O\left(\frac{mNPk_+T_1}{\text{poly}(d)}\right)$, for all $t \in [0, T_1]$, all $c \in [k_+]$,*

$$\begin{aligned} \frac{\Delta A_{+,c,r}^{*(t)}}{\Delta A_{+,r}^{*(t)}} &= \Theta\left(\frac{1}{k_+}\right), \\ \frac{A_{+,c,r}^{*(t)}}{A_{+,r}^{*(t)}} &= \Theta\left(\frac{1}{k_+}\right). \end{aligned} \quad (103)$$

The same identity holds for the “-”-classes.

Proof. The statements in the lemma follow trivially from Theorem D.1 for time period $[0, T_0]$. Let us focus on the phase $[T_0, T_1]$.

In this proof, we condition on the high-probability events of Lemma E.4 and Lemma E.2.

First of all, based on Lemma E.4, we know that $s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| = O(\log(d))$. We will make use of this fact later.

Base case, $t = T_0$.

The base case directly follows from our Theorem D.1.

Induction step, assume statement holds for $\tau \in [T_0, t]$, prove statement for $t+1$.

By Lemma E.2, we know that

$$\begin{aligned} & \Delta A_{+,r}^{*(t)} \\ & = \eta \sum_{c=1}^{k_+} \exp \left\{ - (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota s^*} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) \times \left(A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_{+,c}) \right| \right) \right\} \\ & \quad \times [1/3, 1] (1 \pm s^{*-1/3}) \frac{s^*}{2k_+P} \left(\frac{1}{1 \pm \iota} \pm O\left(\frac{1}{\log^9(d)}\right) \right), \end{aligned} \quad (104)$$

and for any $c \in [k_+]$,

$$\begin{aligned} & \Delta A_{+,c,r}^{*(t)} \\ & = \eta \exp \left\{ - (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota s^*} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) \times \left(A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| + A_{+,c,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_{+,c}) \right| \right) \right\} \\ & \quad \times [1/3, 1] (1 \pm s^{*-1/3}) \frac{s^*}{2k_+P} \left(\frac{1}{1 \pm \iota} \pm O\left(\frac{1}{\log^9(d)}\right) \right), \end{aligned} \quad (105)$$

Relying on the induction hypothesis, we can reduce the above expressions to

$$\begin{aligned}
& \Delta A_{+,r}^{*(t)} \\
&= \eta \sum_{c=1}^{k_+} \exp \left\{ - (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota} \left(1 \pm O \left(\frac{1}{\log^5(d)} \right) \right) \left(1 \pm O \left(\frac{1}{k_+} \right) \right) s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| \right\} \\
& \quad \times [1/3, 1] (1 \pm s^{*-1/3}) \frac{s^*}{2k_+ P} \left(\frac{1}{1 \pm \iota} \pm O \left(\frac{1}{\log^9(d)} \right) \right) \\
&= \eta \exp \left\{ - (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota} \left(1 \pm O \left(\frac{1}{\log^5(d)} \right) \right) \left(1 \pm O \left(\frac{1}{k_+} \right) \right) s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| \right\} \\
& \quad \times \Theta(1) \times \frac{s^*}{2P},
\end{aligned} \tag{106}$$

and for any $c \in [k_+]$,

$$\begin{aligned}
& \Delta A_{+,c,r}^{*(t)} \\
&= \eta \exp \left\{ - (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota} \left(1 \pm O \left(\frac{1}{\log^5(d)} \right) \right) \left(1 \pm O \left(\frac{1}{k_+} \right) \right) s^* A_{+,c,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| \right\} \\
& \quad \times \Theta(1) \times \frac{s^*}{2k_+ P}.
\end{aligned} \tag{107}$$

By invoking the property that $s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| = O(\log(d))$, we find that for all $c \in [k_+]$,

$$\begin{aligned}
\frac{\Delta A_{+,c,r}^{*(t)}}{\Delta A_{+,r}^{*(t)}} &= \exp \left\{ \pm O \left(\frac{1}{\log^5(d)} \right) s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| \right\} \times \Theta \left(\frac{1}{k_+} \right) \\
&= \left(1 \pm O \left(\frac{1}{\log^4(d)} \right) \right) \times \Theta \left(\frac{1}{k_+} \right) \\
&= \Theta \left(\frac{1}{k_+} \right).
\end{aligned} \tag{108}$$

Therefore, we can finish our induction step:

$$\frac{A_{+,c,r}^{*(t+1)}}{A_{+,r}^{*(t+1)}} = \frac{A_{+,c,r}^{*(t)} + \Delta A_{+,c,r}^{*(t)}}{A_{+,r}^{*(t)} + \Delta A_{+,r}^{*(t)}} = \frac{A_{+,c,r}^{*(t)} + \Delta A_{+,c,r}^{*(t)}}{\Theta(k_+) \times (A_{+,c,r}^{*(t)} + \Delta A_{+,c,r}^{*(t)})} = \Theta \left(\frac{1}{k_+} \right). \tag{109}$$

□

Lemma E.6. *Let $T_{(1)}$ be the first point in time such that either $s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| = \Omega(1)$ or $s^* A_{-,r}^{*(t)} \left| S_-^{*(0)}(\mathbf{v}_-) \right| = \Omega(1)$. Then for any $t < T_{(1)}$,*

$$\frac{A_{-,r}^{*(t)}}{A_{+,r}^{*(t)}} = \Theta(1) \tag{110}$$

and for any $t \in [T_{(1)}, T_1]$,

$$\frac{A_{-,r}^{*(t)}}{A_{+,r}^{*(t)}}, \frac{A_{+,r}^{*(t)}}{A_{-,r}^{*(t)}} = \Omega \left(\frac{1}{\log(d)} \right). \tag{111}$$

Proof. This lemma is a consequence of Theorem D.1, Lemma E.2 and Lemma E.4.

Due to Theorem D.1, we already know that $\frac{A_{+,r}^{(t)}}{A_{+,r}^{(0)}} = \Theta(1)$ up to time T_0 . In addition, with Lemma E.2 we know that before $s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| \Omega(1)$, the loss term (on a +-class sample) $1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) = \Theta(1)$ (the same holds with the class signs flipped), in which case it is also easy to derive $\frac{A_{+,r}^{(t)}}{A_{+,r}^{(0)}} = \Theta(1)$ by noting that the update expressions $\Delta A_{-,r}^{*(t)} / \Delta A_{+,r}^{*(t)} = \Theta(1)$.

Beyond time $T_{(1)}$, by Lemma E.4, we know that $s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right|, s^* A_{-,r}^{*(t)} \left| S_-^{*(0)}(\mathbf{v}_-) \right| = O(\log(d))$. With the understanding that $s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right|, s^* A_{-,r}^{*(t)} \left| S_-^{*(0)}(\mathbf{v}_-) \right| \Omega(1)$ beyond $T_{(1)}$ due to the monotonicity of these functions, and the property $\left| \frac{|S_-^{*(0)}(\mathbf{v}_-)|}{|S_+^{*(0)}(\mathbf{v}_+)|} - 1 \right| = O\left(\frac{1}{\log^5(d)}\right)$ from Proposition 1, the rest of the lemma follows. \square

Lemma E.7. *With probability at least $1 - O\left(\frac{mNPk_+ t}{\text{poly}(d)}\right)$, for all $t \in [0, T_1]$ and all $(+, r) \in S_+^{*(0)}(\mathbf{v}_+)$,*

$$\frac{\Delta b_{+,r}^{(t)}}{\Delta A_{+,r}^{*(t)}} = -\Theta\left(\frac{1}{\log^5(d)}\right). \quad (112)$$

The same holds with the +-class signs replaced by the --class signs.

Proof. Choose any $(+, r) \in S_+^{*(0)}(\mathbf{v}_+)$.

The statement in this lemma for time period $t \in [0, T_0]$ follows easily from Theorem D.1 and its proof. Let us examine the period $t \in [T_0, T_1]$.

Based on Lemma E.2 and its proof and Lemma E.5, we know that for $t \in [T_0, T_1]$, with probability at least $1 - O\left(\frac{mNPk_+ t}{\text{poly}(d)}\right)$,

$$\begin{aligned} & \Delta A_{+,r}^{*(t)} \\ &= \eta \exp \left\{ - (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) \left(1 \pm O\left(\frac{1}{k_+}\right) \right) s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| \right\} \\ & \quad \times (1 \pm s^{*-1/3}) \frac{s^*}{NP} \left(\frac{1}{1 \pm \iota} \pm O\left(\frac{1}{\log^9(d)}\right) \right) \sum_{n=1}^N \mathbb{1}\{y_n = +\} \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}) - F_+^{(t)}(\mathbf{X}_n^{(t)})) + 1} \end{aligned} \quad (113)$$

Furthermore,

$$\begin{aligned} & \Delta b_{+,r}^{(t)} \\ &= - \frac{\Delta \mathbf{w}_{+,r}^{(t)}}{\log^5(d)} \\ &= - \eta \frac{1}{\log^5(d)} \exp \left\{ - (1 \pm s^{*-1/3}) \frac{1}{1 \pm \iota} \left(1 \pm O\left(\frac{1}{\log^5(d)}\right) \right) \left(1 \pm O\left(\frac{1}{k_+}\right) \right) s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| \right\} \\ & \quad \times (1 \pm s^{*-1/3}) \frac{s^*}{NP} \left(1 \pm \iota \pm \frac{1}{\log^9(d)} \right) \sum_{n=1}^N \mathbb{1}\{y_n = +\} \frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}) - F_+^{(t)}(\mathbf{X}_n^{(t)})) + 1} \end{aligned} \quad (114)$$

With the understanding that $s^* A_{+,r}^{*(t)} \left| S_+^{*(0)}(\mathbf{v}_+) \right| = O(\log(d))$ from Lemma E.4 and the fact that $\frac{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}))}{\exp(F_-^{(t)}(\mathbf{X}_n^{(t)}) - F_+^{(t)}(\mathbf{X}_n^{(t)})) + 1} = \Theta(1)$, we have

$$\begin{aligned}
\frac{\Delta b_{+,r}^{(t)}}{\Delta A_{+,r}^{(t)}} &= -\Theta\left(\frac{1}{\log^5(d)}\right) \exp\left\{-\left(1 \pm O\left(\frac{1}{\log^5(d)}\right)\right) s^* A_{+,r}^{*(t)} \left|S_+^{*(0)}(\mathbf{v}_+)\right|\right\} \\
&= -\Theta\left(\frac{1}{\log^5(d)}\right) \left(1 \pm O\left(\frac{1}{\log^4(d)}\right)\right) \\
&= -\Theta\left(\frac{1}{\log^5(d)}\right).
\end{aligned} \tag{115}$$

□

Lemma E.8 (Probability of mistake on hard samples is high). *For all $t \in [0, T_1]$, given a hard test sample (\mathbf{X}_{hard}, y) , $y' = y$,*

$$\mathbb{P}\left[F_y^{(T)}(\mathbf{X}_{hard}) = F_y^{(T)}(\mathbf{X}_{hard})\right] = \Omega(1). \tag{116}$$

Proof. We first show that at time $t = 0$, the probability of the network making a mistake on hard test samples is $\Omega(1)$, then prove that for the rest of the time, i.e. $t \in (0, T_1]$, the model still makes mistake on hard test samples with probability $\Omega(1)$.

At time $t = 0$, by Lemma H.3, we know that for any $r \in [m]$, with probability $\Omega(1)$,

$$\mathbf{w}_{+,r}^{(0)} \cdot \mathbf{s}^* = \Omega(\sigma_0 \sigma_\zeta \sqrt{d}) \pm \Omega(\sigma_0 \text{polylog}(d)) \pm \Omega\left(\sigma_0 \sqrt{\log(d)}\right). \tag{117}$$

Relying on concentration of the binomial random variable, with probability at least $1 - e^{-\Omega(\text{polylog}(d))}$,

$$\sum_{r=1}^m \sigma \left(\mathbf{w}_{+,r}^{(0)} \cdot \mathbf{s}^* + b_{+,r}^{(0)} \right) = \Omega(m \sigma_0 \sigma_\zeta \sqrt{d}), \tag{118}$$

which is asymptotically larger than the activation from the features, which, following from Proposition 1, is upper bounded by $O\left(\sigma_0 \sqrt{\log(d)} s^* d^{c_0}\right)$. The same can be said for the “−” class. In other words,

$$\begin{aligned}
&F_-^{(0)}(\mathbf{X}_{hard}) - F_+^{(0)}(\mathbf{X}_{hard}) > 0 \\
&\left\{ \sum_{r=1}^m \mathbb{1}\{\mathbf{w}_{-,r}^{(0)} \cdot \mathbf{s}^* + b_{-,r}^{(0)} > 0\} \mathbf{w}_{-,r}^{(0)} \cdot \mathbf{s}^* \right. \\
&\quad \left. - \sum_{r=1}^m \mathbb{1}\{\mathbf{w}_{+,r}^{(0)} \cdot \mathbf{s}^* + b_{+,r}^{(0)} > 0\} \mathbf{w}_{+,r}^{(0)} \cdot \mathbf{s}^* \right\} (1 \pm o(1)) > 0
\end{aligned} \tag{119}$$

which clearly holds with probability $\Omega(1)$.

Now consider $t \in (0, T_1]$.

During this period of time, by Theorem D.1 and Lemma E.2, we note that for any $c \in [k_+]$ and $(+, r)$ $S_+^{*(0)}(\mathbf{v}_{+,c}), \Delta_{+,r}^{(t)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\zeta_{+,r}}^{(t)2} I_d)$, with $\sigma_{\zeta_{+,r}}^{(t)} = \Theta\left(\Delta A_{+,c,r}^{(t)} \sqrt{\frac{2k_+}{sN}} \sigma_\zeta\right)$. The same can be said for $(+, r)$ $S_+^{*(0)}(\mathbf{v}_+)$, although with the $\Delta A_{+,c,r}^{(t)} \sqrt{\frac{2k_+}{sN}}$ factor replaced by $\Delta A_{+,r}^{(t)} \sqrt{\frac{2}{sN}}$. Also from the proofs of Theorem D.1 and Lemma E.2, and using the property $|U_{+,r}^{(0)}| = O(1)$ from Proposition 1, we know that for all neurons, the updates to the neurons also take the feature-plus-Gaussian-noise form of $\sum_{\mathbf{v}' \in \mathcal{U}_{+,r}^{(0)}} c^{(t)}(\mathbf{v}') \mathbf{v}' + \Delta_{+,r}^{(t)}$, with $c^{(t)}(\mathbf{v}') = \left(1 + O\left(\frac{1}{\log^5(d)}\right)\right) \Delta A_{+,c,r}^{(t)}$ if $\mathbf{v}' = \mathbf{v}_{+,c}$ for some $c \in [k_+]$, or $c^{(t)}(\mathbf{v}') = \left(1 + O\left(\frac{1}{\log^5(d)}\right)\right) \Delta A_{+,r}^{(t)}$ if $\mathbf{v}' = \mathbf{v}_+$ (because the \mathbf{v}' component of a \mathbf{v}' -singleton neuron’s update is already the maximum possible).

Moreover, if $\mathbf{v}_+ \in U_{+,r}^{(0)}$, then $\sigma_{\zeta_{+,r}}^{(t)} = O\left(\Delta A_{+,r}^{(t)} \sqrt{\frac{2}{sN}} \sigma_\zeta\right) + O\left(\Delta A_{+,c,r}^{(t)} \sqrt{\frac{2k_+}{sN}} \sigma_\zeta\right) = O\left(\Delta A_{+,r}^{(t)} \sqrt{\frac{2}{sN}} \sigma_\zeta\right)$, otherwise, if $U_{+,r}^{(0)}$ only contains the fine-grained features, then $\sigma_{\zeta_{+,r}}^{(t)} = O\left(\Delta A_{+,c,r}^{(t)} \sqrt{\frac{2k_+}{sN}} \sigma_\zeta\right)$.

With the understanding that only neurons in $S_y^{(0)}(\mathbf{v}_y)$ and $S_y^{(0)}(\mathbf{v}_{y,c})$ can possibly activate on the feature patches of a sample when $t = T_1$ (coming from Theorem F.1), we have

$$\begin{aligned} F_+^{(t)}(\mathbf{X}_{\text{hard}}) &= \sum_{(+,r) \in S_+^{(0)}(\mathbf{v}_{+,c})} \sum_{p \in \mathcal{P}(\mathbf{X}_{\text{hard}}; \mathbf{v}_{+,c})} \sigma\left(\mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, \overline{1 \pm \iota} \mathbf{v}_{+,c} + \mathbf{p} + b_{+,r}^{(t)}\right) \\ &+ \sum_{r \in [m]} \sigma\left(\mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, * + b_{+,r}^{(t)}\right) \\ &+ \sum_{(+,r) \in S_+^{(0)}(\mathbf{v}_-)} \sum_{p \in \mathcal{P}(\mathbf{X}_{\text{hard}}; \mathbf{v}_-)} \sigma\left(\mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, \alpha_p^\dagger \mathbf{v}_- + \mathbf{p} + b_{+,r}^{(t)}\right) \end{aligned} \quad (120)$$

To further refine this upper bound, we first note that with probability at least $1 - O\left(\frac{mNPk_+t}{\text{poly}(d)}\right)$, the following holds with arbitrary choice of $(+, r^*) \in S_+^{(0)}(\mathbf{v}_{+,c})$:

$$\sum_{(+,r) \in S_+^{(0)}(\mathbf{v}_{+,c})} \sum_{p \in \mathcal{P}(\mathbf{X}_{\text{hard}}; \mathbf{v}_{+,c})} \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, \overline{1 \pm \iota} \mathbf{v}_{+,c} + \mathbf{p} = O\left(s^* \left|S_+^{(0)}(\mathbf{v}_{+,c})\right| \sum_{\tau=0}^{t-1} \Delta A_{+,c,r}^{(\tau)}\right) \quad (121)$$

Invoking Lemma E.5, we obtain (for arbitrary $(+, r^*) \in S_+^{(0)}(\mathbf{v}_+)$):

$$\sum_{(+,r) \in S_+^{(0)}(\mathbf{v}_{+,c})} \sum_{p \in \mathcal{P}(\mathbf{X}_{\text{hard}}; \mathbf{v}_{+,c})} \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, \overline{1 \pm \iota} \mathbf{v}_{+,c} + \mathbf{p} = O\left(\frac{1}{k_+} s^* \left|S_+^{(0)}(\mathbf{v}_{+,c})\right| \sum_{\tau=0}^{t-1} \Delta A_{+,r}^{(\tau)}\right) \quad (122)$$

Let us examine the term $\sum_{r \in [m]} \sigma\left(\mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, * + b_{+,r}^{(t)}\right)$ more carefully. First of all, denoting $S_+^{(0)} = \sum_{c=1}^{k_+} S_+^{(0)}(\mathbf{v}_{+,c}) + \sum_{c=1}^{k_-} S_+^{(0)}(\mathbf{v}_{-,c}) + S_+^{(0)}(\mathbf{v}_+) + S_+^{(0)}(\mathbf{v}_-)$, neurons $(+, r) \in S_+^{(0)}$ cannot receive any update at all during training due to Theorem F.1. Therefore we can rewrite the term

$$\begin{aligned} &\sum_{r \in [m]} \sigma\left(\mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, * + b_{+,r}^{(t)}\right) \\ &= \sum_{(+,r) \in S_+^{(0)}} \sigma\left(\mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, * + b_{+,r}^{(t)}\right) + \sum_{(+,r) \notin S_+^{(0)}} \sigma\left(\mathbf{w}_{+,r}^{(0)}, * + b_{+,r}^{(0)}\right) \end{aligned} \quad (123)$$

Relying on Corollary F.1.1, we know

$$\sum_{\tau=0}^{t-1} \Delta b_{+,r}^{(\tau)} < \sum_{\tau=0}^{t-1} -\Omega\left(\frac{\text{polylog}(d)}{\log^5(d)}\right) \left|\Delta \mathbf{w}_{+,r}^{(\tau)}, *\right|. \quad (124)$$

Therefore, we know that for $r \in [m]$,

$$\sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, * + \Delta b_{+,r}^{(\tau)} = 0 \quad (125)$$

As a consequence, we can write the naive upper bound

$$\begin{aligned}
& \sum_{r \in [m]} \sigma \left(\mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, * + b_{+,r}^{(t)} \right) \\
& \sum_{(+,r) \in S_+^{(0)}} \sigma \left(\mathbf{w}_{+,r}^{(0)}, * + b_{+,r}^{(0)} \right) + \sum_{(+,r) \notin S_+^{(0)}} \sigma \left(\mathbf{w}_{+,r}^{(0)}, * + b_{+,r}^{(0)} \right) \\
& = \sum_{r \in [m]} \sigma \left(\mathbf{w}_{+,r}^{(0)}, * + b_{+,r}^{(0)} \right)
\end{aligned} \tag{126}$$

Additionally, due to Theorem F.1 (and its proof), we know that

$$\begin{aligned}
& \sum_{(+,r) \in S_+^{(0)}(\mathbf{v}_-)} \sum_{p \in \mathcal{P}(X_{\text{hard}}; \mathbf{v}_-)} \sigma \left(\mathbf{w}_{+,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{+,r}^{(\tau)}, \alpha_p^\dagger \mathbf{v}_- + p + b_{+,r}^{(t)} \right) \\
& \sum_{(+,r) \in S_+^{(0)}(\mathbf{v}_-)} \sum_{p \in \mathcal{P}(X_{\text{hard}}; \mathbf{v}_-)} \sigma \left(\mathbf{w}_{+,r}^{(0)}, \alpha_p^\dagger \mathbf{v}_- + p + b_{+,r}^{(0)} \right)
\end{aligned} \tag{127}$$

It follows that

$$\begin{aligned}
& F_+^{(t)}(\mathbf{X}_{\text{hard}}) \\
& O \left(\frac{1}{k_+} s^* \left| S_+^{(0)}(\mathbf{v}_{+,c}) \right| \sum_{\tau=0}^{t-1} \Delta A_{+,r}^{(\tau)} \right) + \sum_{(+,r) \in S_+^{(0)}(\mathbf{v}_{+,c})} \sum_{p \in \mathcal{P}(X_{\text{hard}}; \mathbf{v}_{+,c})} \left| \mathbf{w}_{+,r}^{(0)}, \overline{1 \pm \iota} \mathbf{v}_{+,c} + p \right| \\
& + \sum_{r \in [m]} \sigma \left(\mathbf{w}_{+,r}^{(0)}, * + b_{+,r}^{(0)} \right) + \sum_{(+,r) \in S_+^{(0)}(\mathbf{v}_-)} \sum_{p \in \mathcal{P}(X_{\text{hard}}; \mathbf{v}_-)} \sigma \left(\mathbf{w}_{+,r}^{(0)}, \alpha_p^\dagger \mathbf{v}_- + p + b_{+,r}^{(0)} \right)
\end{aligned} \tag{128}$$

On the other hand, for the “−” neurons, denoting $S_-^{(0)} = \sum_{c=1}^{k_+} S_-^{(0)}(\mathbf{v}_{+,c}) \sum_{c=1}^{k_-} S_-^{(0)}(\mathbf{v}_{-,c}) S_-^{(0)}(\mathbf{v}_+) S_-^{(0)}(\mathbf{v}_-)$,

$$\begin{aligned}
& F_-^{(t)}(\mathbf{X}_{\text{hard}}) \\
& \sum_{(+,r) \in S_-^{(0)}(\mathbf{v}_-)} \sum_{p \in \mathcal{P}(X_{\text{hard}}; \mathbf{v}_-)} \sigma \left(\mathbf{w}_{-,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{-,r}^{(\tau)}, \alpha_p^\dagger \mathbf{v}_- + p + b_{+,r}^{(t)} \right) \\
& + \sum_{(+,r) \notin S_-^{(0)}} \sigma \left(\mathbf{w}_{-,r}^{(0)}, * + b_{+,r}^{(0)} \right),
\end{aligned} \tag{129}$$

note that the last line is true because neurons outside the set $S_-^{(0)}$ cannot receive any update during training with probability at least $1 - O\left(\frac{mNPk_+t}{\text{poly}(d)}\right)$ due to Theorem F.1. Estimating the activation value of the neurons from $S_-^{(0)}(\mathbf{v}_-)$ on the feature noise patches requires some care. We define time t_- to be the first point in time such that any $(-, r^*) \in S_-^{(0)}(\mathbf{v}_-)$ satisfies $\sum_{\tau=0}^{t_-} \Delta A_{-,r}^{(\tau)} \geq \sigma_0 \log^5(d)$, and beyond this point in time, i.e. for $t \in [t_-, T_1]$, the neurons in $S_-^{(0)}(\mathbf{v}_-)$ have to activate with high probability, since

$$\begin{aligned}
& \mathbf{w}_{-,r}^{(0)} + \sum_{\tau=0}^{t-1} \Delta \mathbf{w}_{-,r}^{(\tau)}, \alpha_p^\dagger \mathbf{v}_- + p + b_{+,r}^{(t)} \geq \left(1 - O\left(\frac{1}{\log^5(d)}\right)\right) \sigma_0 \log^5(d) / \log^4(d) - O(\sigma_0 \sqrt{\log(d)}) \\
& > 0.
\end{aligned} \tag{130}$$

Now we can proceed to prove the lemma for $t \in (0, T_1]$ by combining the above estimates for $F_+^{(t)}(\mathbf{X}_{\text{hard}})$ and $F_-^{(t)}(\mathbf{X}_{\text{hard}})$.

For $t \in (0, t_-]$, relying argument similar to the situation of $t = 0$ and the fact that $m - |S_-^{(0)}| = (1 - o(1))m$,

$$\begin{aligned} & \left\{ \sum_{(+, r) \notin S_-^{(0)}} \mathbb{1}\{ \mathbf{w}_{-, r}^{(0), * } + b_{-, r}^{(0)} > 0 \} \mathbf{w}_{-, r}^{(0), * } \right. \\ & \quad \left. - \sum_{r=1}^m \mathbb{1}\{ \mathbf{w}_{+, r}^{(0), * } + b_{+, r}^{(0)} > 0 \} \mathbf{w}_{+, r}^{(0), * } \right\} (1 \pm o(1)) > 0 \\ & = F_-^{(t)}(\mathbf{X}_{\text{hard}}) - F_+^{(t)}(\mathbf{X}_{\text{hard}}) > 0 \end{aligned} \quad (131)$$

which has to be true with probability $\Omega(1)$.

On the other hand, with $t \in (t_-, T_1]$, we have

$$\begin{aligned} & F_-^{(t)}(\mathbf{X}_{\text{hard}}) - F_+^{(t)}(\mathbf{X}_{\text{hard}}) \\ & \left\{ \sum_{\tau=0}^{t-1} \left(1 - O\left(\frac{1}{\log^5(d)}\right) \right) s^\dagger / S_-^{*(0)}(\mathbf{v}_-) / \Delta A_{-, r}^{(\tau)} - O(\sigma_0 \sqrt{\log(d)}) \right. \\ & \quad \left. - O\left(\frac{1}{k_+} s^* \left| S_+^{(0)}(\mathbf{v}_{+, c}) \right| \sum_{\tau=0}^{t-1} \Delta A_{+, r}^{(\tau)} \right) \right\} \\ & + \left\{ \sum_{(+, r) \notin S_-^{(0)}} \sigma\left(\mathbf{w}_{-, r}^{(0), * } + b_{+, r}^{(0)}\right) - \sum_{(+, r) \in S_+^{(0)}(\mathbf{v}_{+, c})} \sum_{p \in \mathcal{P}(\mathbf{X}_{\text{hard}}; \mathbf{v}_{+, c})} \left| \mathbf{w}_{+, r}^{(0), \overline{1 \pm \iota} \mathbf{v}_{+, c} + p} \right| \right. \\ & \quad \left. - \sum_{r \in [m]} \sigma\left(\mathbf{w}_{+, r}^{(0), * } + b_{+, r}^{(0)}\right) - \sum_{(+, r) \in S_+^{(0)}(\mathbf{v}_-)} \sum_{p \in \mathcal{P}(\mathbf{X}_{\text{hard}}; \mathbf{v}_-)} \sigma\left(\mathbf{w}_{+, r}^{(0), \alpha_p^\dagger \mathbf{v}_- + p} + b_{+, r}^{(0)}\right) \right\} \end{aligned} \quad (132)$$

Let us begin analyzing the first $\{\cdot\}$ bracket.

By Proposition 1 we know that $\left| S_-^{*(0)}(\mathbf{v}_-) \right| = (1 \pm O(1/\log^5(d))) \left| S_+^{(0)}(\mathbf{v}_{+, c}) \right|$, and by Lemma E.5, we know that $\Delta A_{+, r}^{(\tau)} = O(\log(d) \Delta A_{-, r}^{(\tau)})$, therefore,

$$\begin{aligned} O\left(\frac{1}{k_+} s^* \left| S_+^{(0)}(\mathbf{v}_{+, c}) \right| \sum_{\tau=0}^{t-1} \Delta A_{+, r}^{(\tau)}\right) &= O\left(\frac{\log(d)}{k_+} s^* \left| S_-^{*(0)}(\mathbf{v}_-) \right| \sum_{\tau=0}^{t-1} \Delta A_{-, r}^{(\tau)}\right) \\ &= \sum_{\tau=0}^{t-1} \left(1 - O\left(\frac{1}{\log^5(d)}\right) \right) s^\dagger / S_-^{*(0)}(\mathbf{v}_-) / \Delta A_{-, r}^{(\tau)} - O(\sigma_0 \sqrt{\log(d)}) \end{aligned} \quad (133)$$

Therefore, we obtained the simpler lower bound

$$\begin{aligned} & F_-^{(t)}(\mathbf{X}_{\text{hard}}) - F_+^{(t)}(\mathbf{X}_{\text{hard}}) \\ & \left\{ \sum_{(+, r) \notin S_-^{(0)}} \sigma\left(\mathbf{w}_{-, r}^{(0), * } + b_{+, r}^{(0)}\right) - \sum_{(+, r) \in S_+^{(0)}(\mathbf{v}_{+, c})} \sum_{p \in \mathcal{P}(\mathbf{X}_{\text{hard}}; \mathbf{v}_{+, c})} \left| \mathbf{w}_{+, r}^{(0), \overline{1 \pm \iota} \mathbf{v}_{+, c} + p} \right| \right. \\ & \quad \left. - \sum_{r \in [m]} \sigma\left(\mathbf{w}_{+, r}^{(0), * } + b_{+, r}^{(0)}\right) - \sum_{(+, r) \in S_+^{(0)}(\mathbf{v}_-)} \sum_{p \in \mathcal{P}(\mathbf{X}_{\text{hard}}; \mathbf{v}_-)} \sigma\left(\mathbf{w}_{+, r}^{(0), \alpha_p^\dagger \mathbf{v}_- + p} + b_{+, r}^{(0)}\right) \right\} \end{aligned} \quad (134)$$

which is greater than 0 with probability $\Omega(1)$ (by relying on an argument almost identical to the $t = 0$ case again, and noting that $m - |S_-^{(0)}| = (1 - o(1))m$). This concludes the proof. \square

Lemma E.9 (Probability of mistake on easy samples is low after training). For $t \in [T_{1,1}, T_1]$, given an easy test sample $(\mathbf{X}_{\text{easy}}, y)$,

$$\mathbb{P} \left[F_y^{(T)}(\mathbf{X}_{\text{easy}}) - F_y^{(T)}(\mathbf{X}_{\text{easy}}) \right] = o(1). \quad (135)$$

Proof. Without loss of generality, assume the true label of \mathbf{X}_{easy} is $+1$. Assume $t \in [T_{1,1}, T_1]$.

Firstly, conditioning on the events of Theorem F.1, the following upper bound on $F_-^{(t)}(\mathbf{X}_{\text{easy}})$ holds with probability at least $1 - O\left(\frac{m}{\text{poly}(d)}\right)$:

$$\begin{aligned} F_-^{(t)}(\mathbf{X}_{\text{easy}}) &= \sum_{(-,r) \in S_-^{(0)}(V_+)} \sum_{p \in \mathcal{P}(X_{\text{easy}}; V_+)} \sigma \left(\mathbf{w}_{-,r}^{(t)}, \overline{1 \pm \iota V_+} + p + b_{-,r}^{(t)} \right) \\ &+ \sum_{(-,r) \in S_-^{(0)}(V_{+,c})} \sum_{p \in \mathcal{P}(X_{\text{easy}}; V_{+,c})} \sigma \left(\mathbf{w}_{-,r}^{(t)}, \overline{1 \pm \iota V_{+,c}} + p + b_{-,r}^{(t)} \right) \\ &+ \sum_{(-,r) \in S_-^{(0)}(V_+)} \sum_{p \in \mathcal{P}(X_{\text{easy}}; V_+)} \sigma \left(\mathbf{w}_{-,r}^{(0)}, \overline{1 \pm \iota V_+} + p + b_{-,r}^{(0)} \right) \\ &+ \sum_{(-,r) \in S_-^{(0)}(V_{+,c})} \sum_{p \in \mathcal{P}(X_{\text{easy}}; V_{+,c})} \sigma \left(\mathbf{w}_{-,r}^{(0)}, \overline{1 \pm \iota V_{+,c}} + p + b_{-,r}^{(0)} \right) \\ &< O(s^* d^{c_0} \sigma_0) \\ &= o(1), \end{aligned} \quad (136)$$

and on the other hand,

$$\begin{aligned} F_+^{(t)}(\mathbf{X}_{\text{easy}}) &= \sum_{(+,r) \in S_+^{(0)}(V_+)} \sum_{p \in \mathcal{P}(X_{\text{easy}}; V_+)} \sigma \left(\mathbf{w}_{+,r}^{(t)}, \overline{1 \pm \iota V_+} + p + b_{+,r}^{(t)} \right) \\ &+ \sum_{(+,r) \in S_+^{(0)}(V_{+,c})} \sum_{p \in \mathcal{P}(X_{\text{easy}}; V_{+,c})} \sigma \left(\mathbf{w}_{+,r}^{(t)}, \overline{1 \pm \iota V_{+,c}} + p + b_{+,r}^{(t)} \right) \\ &> \Omega(1). \end{aligned} \quad (137)$$

Therefore, $F_+^{(t)}(\mathbf{X}_{\text{easy}}) - F_-^{(t)}(\mathbf{X}_{\text{easy}}) > \Omega(1)$, which completes the proof. \square

Lemma E.10 (Jr. & John W. Wrench (1971)). The partial sum of harmonic series satisfies the following identity:

$$\sum_{k=1}^{n-1} \frac{1}{k} = \log(n) + E - \frac{1}{2n} - \epsilon_n \quad (138)$$

where E is the Euler–Mascheroni constant (approximately 0.58), and $\epsilon_n \in [0, 1/8n^2]$.

F Coarse-grained SGD, Poly-time properties

In this section, set $T_e = \text{poly}(d)$.

Please note that we are performing stochastic gradient descent on easy samples only.

Theorem F.1. Fix any $t \in [0, T_e]$.

1. (Non-activation invariance) For any $\tau \leq t$, with probability at least $1 - O\left(\frac{mk_+ NPt}{\text{poly}(d)}\right)$, any feature $\mathbf{v} = \{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\}$, any $t' \leq t$, $(+, r) \in S_+^{(0)}(\mathbf{v})$ and \mathbf{v} -dominated patch sample $\mathbf{x}_{n,p}^{(\tau)} = \alpha_{n,p}^{(\tau)} \mathbf{v} + \beta_{n,p}^{(\tau)}$, the following holds:

$$\sigma\left(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(t)}\right) = 0 \quad (139)$$

2. (Non-activation on noise patches) For any $\tau \leq t$, with probability at least $1 - O\left(\frac{mNPt}{\text{poly}(d)}\right)$, for every $t' \leq t$, $r \in [m]$ and noise patch $\mathbf{x}_{n,p}^{(\tau)} = \beta_{n,p}^{(\tau)}$, the following holds:

$$\sigma\left(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(t)}\right) = 0 \quad (140)$$

3. (Off-diagonal nonpositive growth) For any $\tau \leq t$, with probability at least $1 - O\left(\frac{mk_+ NPt}{\text{poly}(d)}\right)$, for any $t' \leq t$, any feature $\mathbf{v} = \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_-\}$, any $(+, r) \in S_+^{(0)}(\mathbf{v})$ and \mathbf{v} -dominated patch $\mathbf{x}_{n,p}^{(\tau)} = \alpha_{n,p}^{(\tau)} \mathbf{v} + \beta_{n,p}^{(\tau)}$, $\sigma\left(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(t)}\right) \leq \sigma\left(\mathbf{w}_{+,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(0)}\right)$.

Proof. **Base case $t = 0$.**

1. (Nonactivation invariance)

Choose any $\tau \leq 0$, \mathbf{v}^* from the set $\{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\}$. We will work with neuron sets in the “+” class in this proof; the “−”-class case can be handled in the same way.

First, we need to show that, for every n such that $|P(\mathbf{X}_n^{(\tau)}; \mathbf{v}^*)| > 0$ and $p \in P(\mathbf{X}_n^{(\tau)}; \mathbf{v}^*)$, for every $(+, r)$ neuron index,

$$\mathbf{w}_{+,r}^{(0)}, \mathbf{v}^* \leq \sigma_0 \sqrt{4 + 2c_0} \sqrt{\log(d) - \frac{1}{\log^5(d)}} = \sigma\left(\mathbf{w}_{+,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(0)}\right) = 0 \quad (141)$$

This is indeed true. The following holds with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$ for all $(+, r) \in S_+^{(0)}(\mathbf{v})$ and all such $\mathbf{x}_{n,p}^{(\tau)}$:

$$\begin{aligned} \mathbf{w}_{+,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(0)} &\leq \sigma_0 \sqrt{4 + 2c_0} \sqrt{\log(d) - 1/\log^5(d)} + O\left(\frac{\sigma_0}{\log^9(d)}\right) - \sqrt{4 + 2c_0} \sqrt{\log(d)} \sigma_0 \\ &= \sigma_0 \left(\frac{(4 + 2c_0)(1 + \iota)(\log(d) - 1/\log^5(d)) - (4 + 2c_0) \log(d)}{\sqrt{(4 + 2c_0)(\log(d) - 1/\log^5(d)) + \sqrt{4 + 2c_0} \sqrt{\log(d)}}} + O\left(\frac{1}{\log^9(d)}\right) \right) \\ &= \sigma_0 \left(\frac{(4 + 2c_0)\iota \log(d) - (1 + \iota)/\log^5(d)}{\sqrt{(4 + 2c_0)(\log(d) - 1/\log^5(d)) + \sqrt{4 + 2c_0} \sqrt{\log(d)}}} + O\left(\frac{1}{\log^9(d)}\right) \right) \\ &< 0, \end{aligned} \quad (142)$$

The first equality holds by utilizing the identity $a - b = \frac{a^2 - b^2}{a + b}$. As a consequence, $\sigma\left(\mathbf{w}_{+,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(0)}\right) = 0$.

2. (*Non-activation on noise patches*) Invoking Lemma H.3, for any $\tau \geq 0$, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, we have for all possible choices of $r \in [m]$ and the noise patches $\mathbf{x}_{n,p}^{(\tau)} = \mathbf{x}_{n,p}^{(\tau)}$:

$$\left| \mathbf{w}_{+,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} \right| \leq O(\sigma_0 \sigma_\zeta \sqrt{d \log(d)}) + O\left(\frac{\sigma_0}{\log^9(d)}\right) b_{+,r}^{(0)}. \quad (143)$$

Therefore, no neuron can activate on the noise patches at time $t = 0$.

3. (*Off-diagonal nonpositive growth*) This point is trivially true at $t = 0$.

Inductive step: we assume the induction hypothesis for $t \in [0, T]$ (with $T < T_e$ of course), and prove the statements for $t = T + 1$.

1. (*Nonactivation invariance*)

Choose any \mathbf{v}^* from the set $\{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\}$. We will work with neuron sets in the “+” class in this proof; the “-”-class case can be handled in the same way.

We need to prove that given $\tau = T + 1$, with probability at least $1 - O\left(\frac{mNP(T+1)}{\text{poly}(d)}\right)$, for every $t' = T + 1$, $(+, r)$ neuron index and \mathbf{v}^* -dominated patch $\mathbf{x}_{n,p}^{(\tau)}$,

$$(+, r) / S_+^{(0)}(\mathbf{v}^*) = \sigma\left(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(t)}\right) = 0. \quad (144)$$

Conditioning on the (high-probability) event of the induction hypothesis of point 1., the following is already true on all the \mathbf{v}^* -dominated patches at time $t' = T$:

$$(+, r) / S_+^{(0)}(\mathbf{v}^*) = \sigma\left(\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(T)} + b_{+,r}^{(t)}\right) = 0. \quad (145)$$

In particular, $\sigma\left(\mathbf{w}_{+,r}^{(T)}, \mathbf{x}_{n,p}^{(T)} + b_{+,r}^{(T)}\right) = 0$.

In other words, no $(+, r) / S_+^{(0)}(\mathbf{v}^*)$ can be updated on the \mathbf{v}^* -dominated patches at time $t = T$. Furthermore, the induction hypothesis of point 2. also states that the network cannot activate on any noise patch $\mathbf{x}_{n,p}^{(T)} = \mathbf{x}_{n,p}^{(T)}$ with probability at least $1 - O\left(\frac{mNPT}{\text{poly}(d)}\right)$. Therefore, the neuron update for those $(+, r) / S_+^{(0)}(\mathbf{v}^*)$ takes the form

$$\begin{aligned} \Delta \mathbf{w}_{+,r}^{(T)} &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v}^*)} \sum_{n=1}^N \mathbb{1}\{\ell(\mathbf{X}_n^{(T)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \mathbf{w}_{n,p}^{(T)} + b_{c,r}^{(T)} > 0\} \left(\alpha_{n,p}^{(T)} \mathbf{v} + \mathbf{w}_{n,p}^{(T)}\right) \end{aligned} \quad (146)$$

Now we can invoke Lemma F.2 and obtain that, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, the following holds for all relevant neurons and \mathbf{v}^* -dominated patches:

$$\Delta \mathbf{w}_{+,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta b_{+,r}^{(T)} < 0. \quad (147)$$

In conclusion, with $\tau = T + 1$, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, for every $(+, r) / S_+^{(0)}(\mathbf{v}^*)$ and relevant (n, p) 's,

$$\mathbf{w}_{+,r}^{(T)} + \Delta \mathbf{w}_{+,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(T)} + \Delta b_{+,r}^{(T)} = \mathbf{w}_{+,r}^{(T+1)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(T+1)} < 0, \quad (148)$$

which leads to $\mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(t)} < 0$ for all $t' = T + 1$ with probability at least $1 - O\left(\frac{mk_+ NP(T+1)}{\text{poly}(d)}\right)$ (also taking union bound over all the possible choices of \mathbf{v}^*). This finishes the inductive step for point 1.

2. (Non-activation on noise patches)

Relying on the event of the induction hypothesis, for any $\tau = T$, the following holds for every $r = [m]$ and noise patch $\mathbf{x}_{n,p}^{(\tau)} = \begin{pmatrix} \tau \\ n,p \end{pmatrix}$,

$$\mathbf{w}_{+,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(T)} < 0. \quad (149)$$

Conditioning on this high-probability event, this means no neuron $\mathbf{w}_{+,r}^{(T)}$ can be updated on the noise patches. Denoting the set of features $\mathcal{M} = \{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\}$, for every $r = [m]$, its update is reduced to

$$\begin{aligned} \Delta \mathbf{w}_{+,r}^{(T)} &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(T)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \begin{pmatrix} T \\ n,p \end{pmatrix} + b_{+,r}^{(T)} > 0\} \left(\alpha_{n,p}^{(T)} \mathbf{v} + \begin{pmatrix} T \\ n,p \end{pmatrix} \right), \end{aligned} \quad (150)$$

Invoking Lemma F.3, we have that, for any $\tau = T + 1$, the following inequality holds with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$ for every $r = [m]$ and noise patches,

$$\Delta \mathbf{w}_{+,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta b_{+,r}^{(T)} < 0. \quad (151)$$

Consequently, for any $\tau = T + 1$, the following inequality holds with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$ for every $r = [m]$ and noise patches $\mathbf{x}_{n,p}^{(\tau)} = \begin{pmatrix} \tau \\ n,p \end{pmatrix}$:

$$\mathbf{w}_{+,r}^{(T)} + \Delta \mathbf{w}_{+,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(T)} + \Delta b_{+,r}^{(T)} = \mathbf{w}_{+,r}^{(T+1)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,r}^{(T+1)} < 0. \quad (152)$$

This finishes the inductive step for point 2.

3. (Off-diagonal nonpositive growth) Choose any $\mathbf{v}^* = \{\mathbf{v}_-\} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-}$.

Choose any neuron with index $(+, r)$. Similar to our proof for point 2., we know that its update, when taken inner product with a \mathbf{v}^* -dominated patch $\mathbf{x}_{n,p}^{(\tau)} = \overline{1 \pm \iota \mathbf{v}^*} + \begin{pmatrix} \tau \\ n,p \end{pmatrix}$, has to take the form

$$\begin{aligned} &\Delta \mathbf{w}_{+,r}^{(T)}, \overline{1 \pm \iota \mathbf{v}^*} + \begin{pmatrix} \tau \\ n,p \end{pmatrix} \\ &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(T)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \begin{pmatrix} T \\ n,p \end{pmatrix} + b_{+,r}^{(T)} > 0\} \alpha_{n,p}^{(T)} \mathbf{v} + \begin{pmatrix} T \\ n,p \end{pmatrix}, \overline{1 \pm \iota \mathbf{v}^*} + \begin{pmatrix} \tau \\ n,p \end{pmatrix} \\ &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M} - \{\mathbf{v}^*\}} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(T)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \begin{pmatrix} T \\ n,p \end{pmatrix} + b_{+,r}^{(T)} > 0\} \left(\begin{pmatrix} T \\ n,p \end{pmatrix}, \overline{1 \pm \iota \mathbf{v}^*} + \alpha_{n,p}^{(T)} \mathbf{v} + \begin{pmatrix} T \\ n,p \end{pmatrix}, \begin{pmatrix} \tau \\ n,p \end{pmatrix} \right) \\ &\quad - \frac{\eta}{NP} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(T)}; \mathbf{v}^*)| > 0\} [\text{logit}_+^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v}^*)} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \begin{pmatrix} T \\ n,p \end{pmatrix} + b_{+,r}^{(T)} > 0\} \alpha_{n,p}^{(T)} \mathbf{v}^* + \begin{pmatrix} T \\ n,p \end{pmatrix}, \overline{1 \pm \iota \mathbf{v}^*} + \begin{pmatrix} \tau \\ n,p \end{pmatrix} \end{aligned} \quad (153)$$

With probability at least $1 - O\left(\frac{NP}{\text{poly}(d)}\right)$, $\alpha_{n,p}^{(T)} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p}$, $\overline{1 \pm \iota} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p} > 0$, and $\frac{(\tau)}{n,p}$, $\overline{1 \pm \iota} \mathbf{v}_+^{*} + \alpha_{n,p}^{(T)} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p} < O(1/\log^9(d))$. Therefore,

$$\begin{aligned} \Delta \mathbf{w}_{+,r}^{(T)}, \mathbf{v}_+^* &< \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M} - \{\mathbf{v}\}} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(T)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(T)}(\mathbf{X}_n^{(T)})] \\ &\times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v}_+ + \frac{(\tau)}{n,p} + b_{+,r}^{(T)} > 0\} O\left(\frac{1}{\log^9(d)}\right) \end{aligned} \quad (154)$$

Invoking Lemma F.3, we know that

$$\begin{aligned} &\Delta b_{+,r}^{(T)} \\ &- \frac{1}{\log^5(d)} \frac{\eta}{NP} \left(\overline{1 - \iota} - \frac{1}{\log^9(d)} \right) \\ &\times \left(\sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(T)}; \mathbf{v})| > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(T)}(\mathbf{X}_n^{(T)}) \right| \right. \\ &\times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v}_+ + \frac{(\tau)}{n,p} + b_{+,r}^{(T)} > 0\} \right). \end{aligned} \quad (155)$$

It follows that

$$\begin{aligned} &\Delta \mathbf{w}_{+,r}^{(T)}, \overline{1 \pm \iota} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p} + \Delta b_{+,r}^{(T)} \\ &< O\left(\frac{1}{\log^9(d)}\right) \frac{\eta}{NP} \left(\sum_{\mathbf{v} \in \mathcal{M} - \{\mathbf{v}\}} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(T)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(T)}(\mathbf{X}_n^{(T)})] \right. \\ &\times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v}_+ + \frac{(\tau)}{n,p} + b_{c,r}^{(T)} > 0\} \right) \\ &- \Omega\left(\frac{1}{\log^5(d)}\right) \frac{\eta}{NP} \left(\sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(T)}; \mathbf{v})| > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(T)}(\mathbf{X}_n^{(T)}) \right| \right. \\ &\times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v}_+ + \frac{(\tau)}{n,p} + b_{c,r}^{(T)} > 0\} \right) \\ &< 0. \end{aligned} \quad (156)$$

Consequently,

$$\begin{aligned} &\sigma\left(\mathbf{w}_{+,r}^{(T+1)}, \overline{1 \pm \iota} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p} + b_{+,r}^{(T+1)}\right) \\ &= \sigma\left(\mathbf{w}_{+,r}^{(T)}, \overline{1 \pm \iota} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p} + b_{+,r}^{(T)} + \Delta \mathbf{w}_{+,r}^{(T)}, \overline{1 \pm \iota} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p} + \Delta b_{+,r}^{(T)}\right) \\ &\sigma\left(\mathbf{w}_{+,r}^{(T)}, \overline{1 \pm \iota} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p} + b_{+,r}^{(T)}\right) \\ &\sigma\left(\mathbf{w}_{+,r}^{(0)}, \overline{1 \pm \iota} \mathbf{v}_+^{*} + \frac{(\tau)}{n,p} + b_{+,r}^{(0)}\right). \end{aligned} \quad (157)$$

□

Corollary F.1.1 (Bias update upper bound). *Choose any $T_e = \text{poly}(d)$. With probability at least $1 - O\left(\frac{mk + NP T_e}{\text{poly}(d)}\right)$, for all $t \in [0, T_e]$, any neuron $\mathbf{w}_{+,r}$, and any $\mathbf{v} \in U_{+,r}^{(0)}$,*

$$\Delta b_{+,r}^{(t)} < -\Omega\left(\frac{\text{polylog}(d)}{\log^5(d)}\right) \left| \Delta \mathbf{w}_{+,r}^{(t)}, \mathbf{v}_+^* \right|. \quad (158)$$

Proof. Conditioning on the high-probability events of Theorem F.1 above, we know that for any neuron indexed $(+, r)$, at any time $t \in T_e$, its update takes the form

$$\begin{aligned} \Delta \mathbf{w}_{+,r}^{(t)} &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{U}_{+,r}^{(0)}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} \left(\alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} \right), \end{aligned} \quad (159)$$

It follows that, with probability at least $1 - O\left(\frac{1}{\text{poly}(d)}\right)$,

$$\begin{aligned} \left| \Delta \mathbf{w}_{+,r}^{(t)}, * \right| &= \left| \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{U}_{+,r}^{(0)}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ &\quad \times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p}, * \right| \\ &\leq \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{U}_{+,r}^{(0)}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} O\left(\frac{1}{\text{polylog}(d)}\right) \end{aligned} \quad (160)$$

On the other hand,

$$\begin{aligned} \left\| \Delta \mathbf{w}_{+,r}^{(t)} \right\|_2 &\leq \left\| \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{U}_{+,r}^{(0)}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ &\quad \times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} \alpha_{n,p}^{(t)} \mathbf{v} \right\|_2 \\ &\leq \left\| \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{U}_{+,r}^{(0)}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ &\quad \times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} \frac{(t)}{n,p} \right\|_2 \\ &\leq \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{U}_{+,r}^{(0)}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} \left(\frac{1}{1-\iota} - O\left(\frac{1}{\log^9(d)}\right) \right) \end{aligned} \quad (161)$$

Clearly,

$$\left\| \Delta \mathbf{w}_{+,r}^{(t)} \right\|_2 \leq \Omega\left(\text{polylog}(d) \left| \Delta \mathbf{w}_{+,r}^{(t)}, * \right| \right). \quad (162)$$

The conclusion follows. \square

Lemma F.2 (Nonactivation invariance). *Let the assumptions in Theorem D.1 hold.*

Denote the set of features $\mathcal{C}(\mathbf{V}^*) = \{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\} - \{\mathbf{v}^*\}$. If the update term for neuron $\mathbf{w}_{+,r}^{(t)}$ can be written as follows

$$\begin{aligned} \Delta \mathbf{w}_{+,r}^{(t)} &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{V})} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{v}_{n,p}^{(t)} + b_{c,r}^{(t)} > 0\} \left(\alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{v}_{n,p}^{(t)} \right), \end{aligned} \quad (163)$$

then given any $\tau > t$, the following inequality holds with probability at least $1 - O\left(\frac{NP}{\text{poly}(d)}\right)$ for all \mathbf{v}^* -dominated patch $\mathbf{x}_{n,p}^{(\tau)}$:

$$\Delta \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta b_{+,r}^{(t)} < 0 \quad (164)$$

Proof. Let us fix a neuron $\mathbf{w}_{+,r}$ satisfying the update expression in the Lemma statement, and fix some $\tau > t$. Firstly, the bias update for this neuron can be upper bounded via the reverse triangle inequality:

$$\begin{aligned} \Delta b_{+,r}^{(t)} &= - \frac{\left\| \Delta \mathbf{w}_{+,r}^{(t)} \right\|_2}{\log^5(d)} \\ &\quad - \frac{1}{\log^5(d)} \frac{\eta}{NP} \left\| \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{V})} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ &\quad \times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{v}_{n,p}^{(t)} + b_{c,r}^{(t)} > 0\} \alpha_{n,p}^{(t)} \mathbf{v} \right\|_2 \\ &\quad + \frac{1}{\log^5(d)} \frac{\eta}{NP} \left\| \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{V})} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ &\quad \times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{v}_{n,p}^{(t)} + b_{c,r}^{(t)} > 0\} \mathbf{v}_{n,p}^{(t)} \right\|_2 \end{aligned} \quad (165)$$

Let us further upper bound the two ℓ_2 terms separately. Firstly,

$$\begin{aligned}
& \left\| \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \right. \\
& \quad \times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \mathbf{v} \right\|_2 \\
&= \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \left\| \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \right. \\
& \quad \times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \mathbf{v} \right\|_2 \\
&= \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \\
& \quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \|\mathbf{v}\|_2 \\
& \quad \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \\
& \quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \frac{1}{1 - \iota}
\end{aligned} \tag{166}$$

Secondly, with probability at least $1 - O\left(\frac{NP}{\text{poly}(d)}\right)$,

$$\begin{aligned}
& \left\| \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \right. \\
& \quad \times \left. \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \mathbf{v} \right\|_2 \\
& \quad \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \\
& \quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \left\| \alpha_{n,p}^{(t)} \mathbf{v} \right\|_2 \\
& \quad \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \\
& \quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \frac{1}{\log^9(d)}
\end{aligned} \tag{167}$$

Therefore, with probability at least $1 - O\left(\frac{NP}{\text{poly}(d)}\right)$, we can bound the update to the bias as follows:

$$\begin{aligned} & \Delta b_{+,r}^{(t)} \\ & - \frac{1}{\log^5(d)} \frac{\eta}{NP} \left(\frac{1}{1-\iota} - \frac{1}{\log^9(d)} \right) \\ & \times \left(\sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \right. \\ & \left. \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \right) \end{aligned} \quad (168)$$

Furthermore, with probability at least $1 - e^{-(d)+O(\log(d))} > 1 - O\left(\frac{NP}{\text{poly}(d)}\right)$, the following holds for all n, p :

$$\alpha_{n,p}^{(t)} \mathbf{v}, \frac{(t)}{n,p}, \frac{(t)}{n,p}, \alpha_{n,p}^{(\tau)} \mathbf{v}^*, \frac{(t)}{n,p}, \frac{(t)}{n,p} < O\left(\frac{1}{\log^9(d)}\right). \quad (169)$$

Combining the above derivations, they imply that with probability at least $1 - O\left(\frac{NP}{\text{poly}(d)}\right)$, for any $\mathbf{x}_{n,p}^{(\tau)}$ dominated by \mathbf{v}^* ,

$$\begin{aligned} & \Delta \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta b_{+,r}^{(t)} \\ & = \Delta \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(\tau)} \mathbf{v}^* + \frac{(t)}{n,p} + \Delta b_{+,r}^{(t)} \\ & = \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\ & \quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p}, \alpha_{n,p}^{(\tau)} \mathbf{v}^* + \frac{(t)}{n,p} + \Delta b_{+,r}^{(t)} \\ & = \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\ & \quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \left(\alpha_{n,p}^{(t)} \mathbf{v}, \frac{(t)}{n,p} + \frac{(t)}{n,p}, \alpha_{n,p}^{(\tau)} \mathbf{v}^* + \frac{(t)}{n,p}, \frac{(t)}{n,p} \right) \\ & \quad + \Delta b_{+,r}^{(t)} \\ & = \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \\ & \quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \times O\left(\frac{1}{\log^9(d)}\right) + \Delta b_{+,r}^{(t)} \\ & = \frac{\eta}{NP} \left(O\left(\frac{1}{\log^9(d)}\right) - \frac{1}{\log^5(d)} \left(\frac{1}{1-\iota} - \frac{1}{\log^9(d)} \right) \right) \\ & \quad \times \left(\sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \right. \\ & \quad \left. \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0 \} \right) \\ & < 0. \end{aligned} \quad (170)$$

This completes the proof. \square

Lemma F.3 (Nonactivation on noise patches). *Let the assumptions in Theorem D.1 hold.*

Denote the set of features $\mathcal{M} = \{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\}$. If the update term for neuron $\mathbf{w}_{+,r}^{(t)}$ can be written as follows

$$\begin{aligned} \Delta \mathbf{w}_{+,r}^{(t)} &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} \left(\alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} \right), \end{aligned} \quad (171)$$

then

$$\begin{aligned} &\Delta b_{+,r}^{(t)} \\ &\quad - \frac{1}{\log^5(d)} \frac{\eta}{NP} \left(\frac{1}{1-\iota} - \frac{1}{\log^9(d)} \right) \\ &\quad \times \left(\sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \right. \\ &\quad \left. \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} \right). \end{aligned} \quad (172)$$

Moreover, for any $\tau > t$, the following inequality holds with probability at least $1 - O\left(\frac{NP}{\text{poly}(d)}\right)$ for all noise patches $\mathbf{x}_{n,p}^{(\tau)} = \frac{(\tau)}{n,p}$:

$$\Delta \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta b_{+,r}^{(t)} < 0 \quad (173)$$

Proof. Similar to the proof of Lemma F.2, we can estimate the update to the bias term

$$\begin{aligned} &\Delta b_{+,r}^{(t)} \\ &\quad - \frac{1}{\log^5(d)} \frac{\eta}{NP} \left(\frac{1}{1-\iota} - \frac{1}{\log^9(d)} \right) \\ &\quad \times \left(\sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \right. \\ &\quad \left. \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{c,r}^{(t)} > 0\} \right) \end{aligned} \quad (174)$$

Then for any $\mathbf{x}_{n,p}^{(\tau)} = \mathbf{x}_{n,p}^{(t)}$ with $\tau > t$, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$,

$$\begin{aligned}
& \Delta \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta b_{+,r}^{(t)} \\
&= \Delta \mathbf{w}_{+,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + \Delta b_{+,r}^{(t)} \\
&= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\
&\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{b}_{c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{b}_{n,p}^{(t)}, \mathbf{b}_{n,p}^{(\tau)} + \Delta b_{+,r}^{(t)} \\
&= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\
&\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{b}_{n,p}^{(t)} + b_{c,r}^{(t)} > 0 \} \left(\alpha_{n,p}^{(t)} \mathbf{v}, \mathbf{b}_{n,p}^{(\tau)} + \mathbf{b}_{n,p}^{(t)}, \mathbf{b}_{n,p}^{(\tau)} \right) + \Delta b_{+,r}^{(t)} \\
&\quad \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \tag{175} \\
&\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{b}_{n,p}^{(t)} + b_{c,r}^{(t)} > 0 \} \times O\left(\frac{1}{\log^9(d)}\right) + \Delta b_{+,r}^{(t)} \\
&\quad \frac{\eta}{NP} \left(O\left(\frac{1}{\log^9(d)}\right) - \frac{1}{\log^5(d)} \left(\frac{1}{1-\iota} - \frac{1}{\log^9(d)} \right) \right) \\
&\quad \times \left(\sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = +\} - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)}) \right| \right) \\
&\quad \times \left(\sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \mathbf{b}_{n,p}^{(t)} + b_{c,r}^{(t)} > 0 \} \right) \\
&< 0.
\end{aligned}$$

□

G Fine-grained Learning

This section treats the learning dynamics of using fine-grained labels to train the NN; the analysis will be much simpler since the technical analysis overlaps significantly with that in the previous sections.

The training procedure is exactly the same as in the coarse-grained training setting. We explicitly write them out here to avoid any possible confusion.

The learner for fine-grained classification is written as follows for $c \in [k_+]$:

$$F_{+,c}(\mathbf{X}) = \sum_{r=1}^{m_{+,c}} a_{+,c,r} \sum_{p=1}^P \sigma(\mathbf{w}_{+,c,r}, \mathbf{x}_p + b_{+,c,r}), \quad c \in [k_+] \quad (176)$$

with frozen linear classifier weights $a_{+,c,r} = 1$. Same definition applies to the $-$ classes.

The SGD dynamics induced by the training loss is now

$$\begin{aligned} \mathbf{w}_{+,c,r}^{(t+1)} = \mathbf{w}_{+,c,r}^{(t)} + \eta \frac{1}{NP} \sum_{n=1}^N \left(\mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{c,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} + \right. \\ \left. \mathbb{1}\{y_n = (+, c)\} [-\text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \sum_{p \in [P]} \sigma'(\mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(t)} + b_{c,r}^{(t)}) \mathbf{x}_{n,p}^{(t)} \right) \end{aligned} \quad (177)$$

The bias is manually tuned according to the update rule

$$b_{+,c,r}^{(t+1)} = b_{+,c,r}^{(t)} - \frac{\Delta \mathbf{w}_{+,c,r}^{(t)}}{\log^5(d)} \quad (178)$$

We assign $m_{+,c} = \Theta(d^{1+2c_0})$ neurons to each subclass $(+, c)$. For convenience, we write $m = dm_{+,c}$.

The initialization scheme is identical to the coarse-training case, except we choose a slightly less negative $b_{c,r}^{(0)} = -\sigma_0 \sqrt{2 + 2c_0} \sqrt{\log(d)}$.

The parameter choices remain the same as before.

G.1 Initialization geometry

Definition G.1. Define the following sets of interest of the hidden neurons:

1. $U_{+,c,r}^{(0)} = \{\mathbf{v} \in \mathcal{V} : \|\mathbf{w}_{+,c,r}^{(0)}, \mathbf{v}\| \leq \sigma_0 \sqrt{2 + 2c_0} \sqrt{\log(d) - \frac{1}{\log^5(d)}}\}$
2. Given $\mathcal{V} \subseteq \mathcal{V}$, $S_{+,c}^{*(0)}(\mathcal{V}) \subseteq (+, c) \times [m_{+,c}]$ satisfies:
 - (a) $\|\mathbf{w}_{+,c,r}^{(0)}, \mathbf{v}\| \leq \sigma_0 \sqrt{2 + 2c_0} \sqrt{\log(d) + \frac{1}{\log^5(d)}}$
 - (b) $\mathbf{v}' \in \mathcal{V}$ s.t. $\|\mathbf{v}', \mathbf{v}\| \leq \sigma_0 \sqrt{2 + 2c_0} \sqrt{\log(d) - \frac{1}{\log^5(d)}}$
3. Given $\mathcal{V} \subseteq \mathcal{V}$, $S_{+,c}^{(0)}(\mathcal{V}) \subseteq (+, c) \times [m_{+,c}]$ satisfies:
 - (a) $\|\mathbf{w}_{+,c,r}^{(0)}, \mathbf{v}\| \leq \sigma_0 \sqrt{2 + 2c_0} \sqrt{\log(d) - \frac{1}{\log^5(d)}}$
4. For any $(+, c, r) \in S_{+,c,reg}^{*(0)} \subseteq (+, c) \times [m_{+,c}]$:
 - (a) $\|\mathbf{w}_{+,c,r}^{(0)}, \mathbf{v}\| \leq \sigma_0 \sqrt{10} \sqrt{\log(d)}$ $\forall \mathbf{v} \in \mathcal{V}$
 - (b) $|U_{+,c,r}^{(0)}| = O(1)$

The same definitions apply to the $-$ -class neurons.

Proposition 2. *At $t = 0$, for all $\mathbf{v} \in D$, the following properties are true with probability at least $1 - d^{-2}$ over the randomness of the initialized kernels:*

1. $|S_{+,c}^{*(0)}(\mathbf{v})|, |S_{+,c}^{(0)}(\mathbf{v})| = \Theta\left(\frac{1}{\log(d)}\right) d^{c_0}$
2. In particular, $\left|\frac{|S_{y'}^{(0)}(\mathbf{v})|}{|S_y^{(0)}(\mathbf{v})|} - 1\right| = O\left(\frac{1}{\log^5(d)}\right)$ and $\left|\frac{|S_{y'}^{(0)}(\mathbf{v})|}{|S_y^{(0)}(\mathbf{v})|} - 1\right| = O\left(\frac{1}{\log^5(d)}\right)$ for any y, y'
 $\{(+, c)\}_{c=1}^{k_+}$ $\{(-, c)\}_{c=1}^{k_-}$ and common or fine-grained features \mathbf{v}, \mathbf{v}' .
3. $S_{+,c,reg}^{(0)} = [m_{+,c}]$

The same properties apply to the $-$ -class neurons.

Proof. This proof proceeds in virtually the same way as in the proof of Proposition 1, so we omit it here. \square

G.2 Poly-time properties

Theorem G.1. *Fix any $t \in [0, T_e]$, assuming $T_e = \text{poly}(d)$.*

1. (Non-activation invariance) *For any $\tau \leq t$, with probability at least $1 - O\left(\frac{mk_+ N P t}{\text{poly}(d)}\right)$, for any feature $\mathbf{v} \in \{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\}$, for every $t' \leq t$, $(+, c, r) \in S_{+,c}^{(0)}(\mathbf{v})$ and \mathbf{v} -dominated patch sample $\mathbf{x}_{n,p}^{(\tau)} = \alpha_{n,p}^{(\tau)} \mathbf{v} + \beta_{n,p}^{(\tau)}$, the following holds:*

$$\sigma\left(\mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(t)}\right) = 0 \quad (179)$$

2. (Non-activation on noise patches) *For any $\tau \leq t$, with probability at least $1 - O\left(\frac{m N P t}{\text{poly}(d)}\right)$, for every $c \in [k_+]$, $r \in [m]$ and noise patch $\mathbf{x}_{n,p}^{(\tau)} = \beta_{n,p}^{(\tau)}$, the following holds:*

$$\sigma\left(\mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(t)}\right) = 0 \quad (180)$$

3. (Off-diagonal nonpositive growth) *Given fine-grained class $(+, c)$ and any $\tau \leq t$, with probability at least $1 - O\left(\frac{mk_+ N P t}{\text{poly}(d)}\right)$, for any $t' \leq t$, any feature $\mathbf{v} \in \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_-\} \cup \{\mathbf{v}_{+,c'}\}_{c' \neq c}$, any neuron $\mathbf{w}_{+,c,r} \in S_{+,c}^{(0)}(\mathbf{v})$ and any \mathbf{v} -dominated patch $\mathbf{x}_{n,p}^{(\tau)} = \alpha_{n,p}^{(\tau)} \mathbf{v} + \beta_{n,p}^{(\tau)}$, $\sigma\left(\mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(t)}\right) \leq \sigma\left(\mathbf{w}_{+,c,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(0)}\right)$.*

Proof. The proof of this theorem is similar to that of Theorem F.1, but with some subtle differences.

Base case $t = 0$.

1. (Nonactivation invariance)

Choose any \mathbf{v}^* from the set $\{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\}$. We will work with neuron sets in the “+” class in this proof; the “-”-class case can be handled in the same way.

First, given $\tau = 0$, we need to show that, for every n such that $|P(\mathbf{X}_n^{(0)}; \mathbf{v}^*)| > 0$ and $p \in P(\mathbf{X}_n^{(0)}; \mathbf{v}^*)$, for every $(+, c, r)$ neuron index,

$$\mathbf{w}_{+,c,r}^{(0)}, \mathbf{v}^* \leq \sigma_0 \sqrt{\frac{1}{2 + 2c_0} \log(d) - \frac{1}{\log^5(d)}} = \sigma\left(\mathbf{w}_{+,c,r}^{(0)}, \mathbf{x}_{n,p}^{(0)} + b_{+,c,r}^{(0)}\right) = 0 \quad (181)$$

This is indeed true. The following holds with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$ for all $(+, r) / S_+^{(0)}(\mathbf{v})$ and all such $\mathbf{x}_{n,p}^{(\tau)}$:

$$\begin{aligned}
& \mathbf{w}_{+,c,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(0)} \\
& \sigma_0 \frac{1 + \iota \sqrt{(2 + 2c_0)(\log(d) - 1/\log^5(d))}}{1 + \iota \sqrt{(2 + 2c_0)(\log(d) - 1/\log^5(d))}} + O\left(\frac{\sigma_0}{\log^9(d)}\right) - \frac{1}{2 + 2c_0} \sqrt{\log(d)} \sigma_0 \\
& = \sigma_0 \left(\frac{(2 + 2c_0)(1 + \iota)(\log(d) - 1/\log^5(d)) - (2 + 2c_0) \log(d)}{\sqrt{(2 + 2c_0)(\log(d) - 1/\log^5(d))} + \frac{1}{4 + 2c_0} \sqrt{\log(d)}} + O\left(\frac{1}{\log^9(d)}\right) \right) \\
& = \sigma_0 \left(\frac{(2 + 2c_0)(\iota \log(d) - (1 + \iota)/\log^5(d))}{\sqrt{(2 + 2c_0)(\log(d) - 1/\log^5(d))} + \frac{1}{2 + 2c_0} \sqrt{\log(d)}} + O\left(\frac{1}{\log^9(d)}\right) \right) \\
& < 0,
\end{aligned} \tag{182}$$

The first equality holds by utilizing the identity $a - b = \frac{a^2 - b^2}{a + b}$. As a consequence, $\sigma(\mathbf{w}_{+,c,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(0)}) = 0$.

2. (Non-activation on noise patches) Invoking Lemma H.3, for any $\tau = 0$, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, we have for all possible choices of $r \in [m]$ and the noise patches $\mathbf{x}_{n,p}^{(\tau)} = \mathbf{x}_{n,p}^{(\tau)}$:

$$\left| \mathbf{w}_{+,c,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} \right| \leq O(\sigma_0 \sigma_\zeta \sqrt{d \log(d)}) + O\left(\frac{\sigma_0}{\log^9(d)}\right) + b_{+,c,r}^{(0)}. \tag{183}$$

Therefore, no neuron can activate on the noise patches at time $t = 0$.

3. (Off-diagonal nonpositive growth) This point is trivially true at $t = 0$.

Inductive step: we assume the induction hypothesis for $t \in [0, T]$ (with $T < T_e$ of course), and prove the statements for $t = T + 1$.

1. (Nonactivation invariance)

Again, choose any \mathbf{v}^* from the set $\{\mathbf{v}_{+,c}\}_{c=1}^{k_+} \cup \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_+, \mathbf{v}_-\}$.

We need to prove that given $\tau = T + 1$, with probability at least $1 - O\left(\frac{mk_+ NP(T+1)}{\text{poly}(d)}\right)$, for every $t' = T + 1$, $(+, c, r)$ neuron index and \mathbf{v}^* -dominated patch $\mathbf{x}_{n,p}^{(\tau)}$,

$$(+, c, r) / S_{+,c}^{(0)}(\mathbf{v}^*) = \sigma\left(\mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(t)}\right) = 0. \tag{184}$$

By the induction hypothesis of point 1., with probability at least $1 - O\left(\frac{mk_+ NPT}{\text{poly}(d)}\right)$, the following is already true on all the \mathbf{v}^* -dominated patches at time $t' = T$:

$$(+, c, r) / S_{+,c}^{(0)}(\mathbf{v}^*) = \sigma\left(\mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(T)} + b_{+,c,r}^{(t)}\right) = 0. \tag{185}$$

In particular, $\sigma\left(\mathbf{w}_{+,c,r}^{(T)}, \mathbf{x}_{n,p}^{(T)} + b_{+,c,r}^{(T)}\right) = 0$.

In other words, no $(+, c, r) / S_{+,c}^{(0)}(\mathbf{v}^*)$ can be updated on the \mathbf{v}^* -dominated patches at time $t = T$. Furthermore, the induction hypothesis of point 2. also states that the network cannot activate on any noise patch $\mathbf{x}_{n,p}^{(T)} = \mathbf{x}_{n,p}^{(T)}$ with probability at least $1 - O\left(\frac{mNPT}{\text{poly}(d)}\right)$. Therefore, the neuron update for those

$(+, c, r) / S_{+,c}^{(0)}(\mathbf{v}^*)$ takes the form

$$\begin{aligned} \Delta \mathbf{w}_{+,c,r}^{(T)} &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(T)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(T)}; \mathcal{V})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \frac{(T)}{n,p} + b_{+,c,r}^{(T)} > 0 \} \left(\alpha_{n,p}^{(T)} \mathbf{v} + \frac{(T)}{n,p} \right) \end{aligned} \quad (186)$$

Conditioning on this high-probability event, we have

$$\begin{aligned} \Delta \mathbf{b}_{+,c,r}^{(t)} &= - \frac{\|\Delta \mathbf{w}_{+,c,r}^{(t)}\|_2}{\log^5(d)} \\ &\quad - \frac{1}{\log^5(d)} \frac{\eta}{NP} \left\| \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ &\quad \times \left. \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \mathbf{v} \right\|_2 \\ &\quad + \frac{1}{\log^5(d)} \frac{\eta}{NP} \left\| \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ &\quad \times \left. \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0 \} \frac{(t)}{n,p} \right\|_2 \end{aligned} \quad (187)$$

Let us further upper bound the two \cdot_2 terms separately. Firstly,

$$\begin{aligned} &\left\| \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ &\quad \times \left. \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \mathbf{v} \right\|_2 \\ &= \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \right| \\ &\quad \times \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0 \} \alpha_{n,p}^{(t)} \|\mathbf{v}\|_2 \\ &\quad \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \right| \\ &\quad \times \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0 \} \frac{1}{1-\iota} \end{aligned} \quad (188)$$

For the second \cdot_2 term consisting purely of noise, note that since all the $\frac{(t)}{n,p}$'s are independent Gaussian random vectors, the standard deviation of the sum is in fact

$$\begin{aligned} &\left\{ \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0 \} \right. \\ &\quad \times \left. [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})]^2 \right\}^{1/2} \sigma_\zeta. \end{aligned} \quad (189)$$

With the basic property that $\sqrt{\sum_j c_j^2} = \sum_j |c_j|$ for any sequence of real numbers c_1, c_2, \dots , we know this standard deviation can be upper bounded by

$$\begin{aligned} & \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v})| > 0\} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0\} \\ & \times \left| \mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \right| \sigma_\zeta \end{aligned} \quad (190)$$

It follows that with probability at least $1 - O\left(\frac{1}{\text{poly}(d)}\right)$,

$$\begin{aligned} & \left\| \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v})| > 0\} [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \right. \\ & \times \left. \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0\} \frac{(t)}{n,p} \right\|_2 \\ & \sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v})| > 0\} \left| \mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \right| \\ & \times \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0\} \frac{1}{\log^9(d)} \end{aligned} \quad (191)$$

Therefore, we can upper bound the bias update as follows:

$$\begin{aligned} \Delta b_{+,c,r}^{(t)} & - \frac{1}{\log^5(d)} \frac{\eta}{NP} \left(\frac{1}{1-\iota} - \frac{1}{\log^9(d)} \right) \\ & \times \left(\sum_{\mathbf{v} \in \mathcal{C}(\mathcal{V})} \sum_{n=1}^N \mathbb{1}\{|P(\mathbf{X}_n^{(t)}; \mathbf{v})| > 0\} \left| \mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \right| \right. \\ & \times \left. \sum_{p \in \mathcal{P}(\mathcal{X}_n^{(t)}; \mathcal{V})} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0\} \right) \end{aligned} \quad (192)$$

Furthermore, with probability at least $1 - O\left(\frac{NP}{\text{poly}(d)}\right)$, the following holds for all n, p :

$$\alpha_{n,p}^{(t)} \mathbf{v}, \frac{(t)}{n,p}, \frac{(t)}{n,p}, \alpha_{n,p}^{(\tau)} \mathbf{v}^*, \frac{(t)}{n,p}, \frac{(\tau)}{n,p} < O\left(\frac{1}{\log^9(d)}\right). \quad (193)$$

Combining the above derivations, they imply that with probability at least $1 - O\left(\frac{NP}{\text{poly}(d)}\right)$, for any $\mathbf{x}_{n,p}^{(\tau)}$ dominated by \mathbf{v}^* ,

$$\begin{aligned}
& \Delta \mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta b_{+,c,r}^{(t)} \\
&= \Delta \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(\tau)} \mathbf{v}^* + \alpha_{n,p}^{(\tau)} + \Delta b_{+,c,r}^{(t)} \\
&= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \\
&\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \alpha_{n,p}^{(t)} + b_{+,c,r}^{(t)} > 0\} \alpha_{n,p}^{(t)} \mathbf{v} + \alpha_{n,p}^{(t)} + b_{+,c,r}^{(t)} + \Delta b_{+,c,r}^{(t)} \\
&= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \\
&\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \alpha_{n,p}^{(t)} + b_{+,c,r}^{(t)} > 0\} \left(\alpha_{n,p}^{(t)} \mathbf{v}, \alpha_{n,p}^{(t)} + \alpha_{n,p}^{(\tau)} \mathbf{v}^* + \alpha_{n,p}^{(\tau)} + b_{+,c,r}^{(t)} \right) \\
&\quad + \Delta b_{+,c,r}^{(t)} \tag{194} \\
&\frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \right| \\
&\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \alpha_{n,p}^{(t)} + b_{+,c,r}^{(t)} > 0\} \times O\left(\frac{1}{\log^9(d)}\right) + \Delta b_{+,c,r}^{(t)} \\
&\frac{\eta}{NP} \left(O\left(\frac{1}{\log^9(d)}\right) - \frac{1}{\log^5(d)} \left(\frac{1}{1-\iota} - \frac{1}{\log^9(d)} \right) \right) \\
&\quad \times \left(\sum_{\mathbf{v} \in \mathcal{C}(\mathbf{v})} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \right| \right) \\
&\quad \times \left(\sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \alpha_{n,p}^{(t)} + b_{+,c,r}^{(t)} > 0\} \right) \\
&< 0.
\end{aligned}$$

Therefore, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, the following holds for the relevant neurons and \mathbf{v}^* -dominated patches:

$$\Delta \mathbf{w}_{+,c,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta b_{+,c,r}^{(T)} < 0. \tag{195}$$

In conclusion, with $\tau = T + 1$, with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, for every $(+, c, r) \in S_{+,c}^{(0)}(\mathbf{v}^*)$ and relevant (n, p) 's,

$$\mathbf{w}_{+,c,r}^{(T)} + \Delta \mathbf{w}_{+,c,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(T)} + \Delta b_{+,c,r}^{(T)} = \mathbf{w}_{+,c,r}^{(T+1)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(T+1)} < 0, \tag{196}$$

which leads to $\mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + b_{+,c,r}^{(t)} < 0$ for all $t' = T + 1$ with probability at least $1 - O\left(\frac{mk_+ NP(T+1)}{\text{poly}(d)}\right)$ (also by taking union bound over all the possible choices of \mathbf{v}^* at time $T + 1$). This finishes the inductive step for point 1.

2. (Non-activation on noise patches)

The inductive step for this part is very similar to (and even simpler than) the inductive step of point 1, so we omit the calculations here.

3. (*Off-diagonal nonpositive growth*) By the induction hypothesis's high-probability event, we already have that, given any fine-grained class $(+, c)$, $\tau = T + 1$, for any feature $\mathbf{v}^* \in \{\mathbf{v}_{-,c}\}_{c=1}^{k_-} \cup \{\mathbf{v}_-\} \cup \{\mathbf{v}_{+,c}\}_{c \neq +}$ and any neuron $\mathbf{w}_{+,c,r}$, $\sigma \left(\mathbf{w}_{+,c,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + \mathbf{b}_{+,c,r}^{(T)} \right) = \sigma \left(\mathbf{w}_{+,c,r}^{(0)}, \mathbf{x}_{n,p}^{(\tau)} + \mathbf{b}_{+,c,r}^{(0)} \right)$. We just need to show that $\Delta \mathbf{w}_{+,c,r}^{(t)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta \mathbf{b}_{+,c,r}^{(t)} = 0$ to finish the proof; the rest proceeds in a similar fashion to the induction step of point 3 in the proof of Theorem F.1.

Similar to the induction step of point 1, denoting \mathcal{M} to be the set of all common and fine-grained features, the update expression of any neuron $(+, c, r)$ has to be

$$\begin{aligned} \Delta \mathbf{w}_{+,c,r}^{(T)} &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(T)}; \mathbf{v}) > 0\} [\mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \frac{(T)}{n,p} + \mathbf{b}_{+,c,r}^{(T)} > 0 \} \left(\alpha_{n,p}^{(T)} \mathbf{v} + \frac{(T)}{n,p} \right) \end{aligned} \quad (197)$$

Written more explicitly,

$$\begin{aligned} \Delta \mathbf{w}_{+,c,r}^{(T)} &= \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M} - \{\mathbf{v}^*\}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(T)}; \mathbf{v}) > 0\} \mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_{+,c}^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \frac{(T)}{n,p} + \mathbf{b}_{+,c,r}^{(T)} > 0 \} \left(\alpha_{n,p}^{(T)} \mathbf{v} + \frac{(T)}{n,p} \right) \\ &\quad - \frac{\eta}{NP} \sum_{n=1}^N \mathbb{1}\{y_n = (+, c)\} \mathbb{1}\{P(\mathbf{X}_n^{(T)}; \mathbf{v}^*) > 0\} [\text{logit}_{+,c}^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v}^*)} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v}^* + \frac{(T)}{n,p} + \mathbf{b}_{+,c,r}^{(T)} > 0 \} \left(\alpha_{n,p}^{(T)} \mathbf{v}^* + \frac{(T)}{n,p} \right) \end{aligned} \quad (198)$$

It follows that with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, for relevant n, p, r , we have

$$\begin{aligned} &\Delta \mathbf{w}_{+,c,r}^{(T)}, \alpha_{n,p}^{(\tau)} \mathbf{v}^* + \frac{(\tau)}{n,p} \\ &< \frac{\eta}{NP} \sum_{\mathbf{v} \in \mathcal{M} - \{\mathbf{v}^*\}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(T)}; \mathbf{v}) > 0\} \mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_{+,c}^{(T)}(\mathbf{X}_n^{(T)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(T)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(T)}, \alpha_{n,p}^{(T)} \mathbf{v} + \frac{(T)}{n,p} + \mathbf{b}_{+,c,r}^{(T)} > 0 \} O\left(\frac{1}{\log^9(d)}\right) \end{aligned} \quad (199)$$

Furthermore, similar to the induction step of point 1, we can estimate the bias update as follows:

$$\begin{aligned} &\Delta \mathbf{b}_{+,c,r}^{(t)} \\ &- \Omega\left(\frac{1}{\log^5(d)}\right) \frac{\eta}{NP} \left(\sum_{\mathbf{v} \in \mathcal{M}} \sum_{n=1}^N \mathbb{1}\{P(\mathbf{X}_n^{(t)}; \mathbf{v}) > 0\} \left| \mathbb{1}\{y_n = (+, c)\} - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \right| \right. \\ &\quad \left. \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v} + \frac{(t)}{n,p} + \mathbf{b}_{+,c,r}^{(t)} > 0 \} \right) \end{aligned} \quad (200)$$

It follows that, indeed, $\Delta \mathbf{w}_{+,c,r}^{(T)}, \mathbf{x}_{n,p}^{(\tau)} + \Delta \mathbf{b}_{+,c,r}^{(T)} = 0$, which completes the induction step of point 3. \square

G.3 Training

Choose an arbitrary constant $B \in [\Omega(1), \log(3/2)]$.

Definition G.2. Let $T_0(B) > 0$ be the first time that there exists some $\mathbf{X}_n^{(t)}$ and c such that $F_y^{(T_0(B))}(\mathbf{X}_n^{(T_0(B))}) \leq B$ for any $n \in [N]$ and $y \in \{(+, c)\}_{c=1}^{k_+} \cup \{(-, c)\}_{c=1}^{k_-}$.

We write $T_0(B)$ as T_0 for simplicity of notation when the context is clear.

Lemma G.2. *With probability at least $1 - O\left(\frac{mk_+ N P T_0}{\text{poly}(d)}\right)$, the following holds for all $t \in [0, T_0]$:*

1. (On-diagonal common-feature neuron growth) For every $c \in [k_+]$, every $(+, c, r), (+, c, r') \in S_{+,c}^{*(0)}(\mathbf{v}_+)$,

$$\mathbf{w}_{+,c,r}^{(t)} - \mathbf{w}_{+,c,r}^{(0)} = \mathbf{w}_{+,c,r}^{(t)} - \mathbf{w}_{+,c,r}^{(0)} \quad (201)$$

Moreover,

$$\Delta \mathbf{w}_{+,c,r}^{(t)} = [1/4, 2/3] \frac{1}{1 \pm \iota} \left(1 \pm s^{*-1/3}\right) \eta \frac{s^*}{2k_+ P} \mathbf{v}_+ + \Delta_{+,r}^{(t)} \quad (202)$$

where $\Delta_{+,c,r}^{(t)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\zeta_{+,c,r}}^{(t)2} \mathbf{I})$, $\sigma_{\zeta_{+,c,r}}^{(t)} = \Theta(1) \times \eta \sigma_\zeta \frac{\sqrt{s}}{P\sqrt{2N}}$.

The bias updates satisfy

$$\Delta b_{+,c,r}^{(t)} = -\Theta\left(\frac{\eta s^*}{k_+ P \log^5(d)}\right). \quad (203)$$

Furthermore, every $(+, r) \in S_+^{*(0)}(\mathbf{v}_+)$ activates on all the \mathbf{v}_+ -dominated patches at time t .

2. (On-diagonal finegrained-feature neuron growth) For every $c \in [k_+]$ and every $(+, c, r), (+, c, r') \in S_{+,c}^{*(0)}(\mathbf{v}_{+,c})$,

$$\mathbf{w}_{+,c,r}^{(t)} - \mathbf{w}_{+,c,r}^{(0)} = \mathbf{w}_{+,c,r}^{(t)} - \mathbf{w}_{+,c,r}^{(0)} \quad (204)$$

Moreover,

$$\Delta \mathbf{w}_{+,c,r}^{(t)} = \left(1 \pm O\left(\frac{1}{k_+}\right)\right) \frac{1}{1 \pm \iota} \left(1 \pm s^{*-1/3}\right) \eta \frac{s^*}{2k_+ P} \mathbf{v}_{+,c} + \Delta_{+,r}^{(t)} \quad (205)$$

where $\Delta_{+,c,r}^{(t)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\zeta_{+,c,r}}^{(t)2} \mathbf{I})$, and $\sigma_{\zeta_{+,c,r}}^{(t)} = \left(1 \pm O\left(\frac{1}{k_+}\right)\right) (1 \pm s^{*-1/3}) \eta \sigma_\zeta \frac{\sqrt{s}}{P \cdot 2Nk_+}$.

The bias updates satisfy

$$\Delta b_{+,c,r}^{(t)} = -\Theta\left(\frac{\eta s^*}{k_+ P \log^5(d)}\right). \quad (206)$$

Furthermore, every $(+, c, r) \in S_{+,c}^{*(0)}(\mathbf{v}_{+,c})$ activates on all the $\mathbf{v}_{+,c}$ -dominated patches at time t .

3. The above results also hold with the “+” and “-” class signs flipped.

Proof. The proof of this theorem proceeds in a similar fashion to Theorem D.1, with some variations for the common-feature neurons.

We shall prove the statements in this theorem via induction. We focus on the +-class neurons; --class neurons' proofs are done in the same fashion.

First of all, relying on the (high-probability) event of Theorem G.1, we know that we can simplify the update expressions for the neurons in $S_{+,c}^{*(0)}(\mathbf{v}_{+,c})$ to the form

$$\begin{aligned} \Delta \mathbf{w}_{+,c,r}^{(t)} &= \frac{\eta}{NP} \sum_{n=1}^N \mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \\ &\quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_{+,c})} \mathbb{1}\{\mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v}_{+,c} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0\} \left(\alpha_{n,p}^{(t)} \mathbf{v}_{+,c} + \frac{(t)}{n,p}\right), \end{aligned} \quad (207)$$

and for the neurons in $S_{+,c}^{*(0)}(\mathbf{v}_+)$, the updates take the form

$$\begin{aligned} & \Delta \mathbf{w}_{+,c,r}^{(t)} \\ &= \frac{\eta}{NP} \sum_{n=1}^N \left\{ \mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] + \sum_{c \in [k_+] - \{c\}} \mathbb{1}\{y_n = (+, c')\} [-\text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \right\} \\ & \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v}_+ + \alpha_{n,p}^{(t)} + b_{+,c,r}^{(t)} > 0 \} \left(\alpha_{n,p}^{(t)} \mathbf{v}_+ + \alpha_{n,p}^{(t)} \right). \end{aligned} \quad (208)$$

By definition of T_0 and the fact that $B = \log(3/2)$, for any $n \in [N]$ and $t < T_0$, we can write down a simple upper bound of $\text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})$:

$$\begin{aligned} \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) &= \frac{\exp(F_{+,c}(\mathbf{X}_n^{(t)}))}{\sum_{c=1}^{k_+} \exp(F_{+,c}(\mathbf{X}_n^{(t)})) + \sum_{c=1}^{k_-} \exp(F_{-,c}(\mathbf{X}_n^{(t)}))} \\ & \leq \frac{3}{2k_+} = \frac{3}{4k_+}, \end{aligned} \quad (209)$$

and we can lower bound it as follows

$$\text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)}) \geq \frac{1}{2k_+ \times \frac{3}{2}} = \frac{1}{3k_+}, \quad (210)$$

The inductive proof for the fine-grained neurons $S_{+,c}^{*(0)}(\mathbf{v}_{+,c})$ is almost identical to that in the proof of Theorem D.1. The only notable difference here is that $[1 - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})]$ has the estimate $(1 \pm O(\frac{1}{k_+}))$.

The inductive proof of the common-feature neurons $S_{+,c}^{*(0)}(\mathbf{v}_+)$ requires more care as its update expression equation 208 is qualitatively different from the coarse-grained training case in Theorem D.1, so we present the full proof here.

Base case, $t = 0$.

With probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$, for every $c \in [k_+]$ and every $(+, c, r) \in S_{+,c}^{*(0)}(\mathbf{v}_+)$,

$$\begin{aligned} & \mathbf{w}_{+,c,r}^{(0)}, \alpha_{n,p}^{(0)} \mathbf{v}_+ + \alpha_{n,p}^{(0)} + b_{+,c,r}^{(0)} \\ & \geq \sigma_0 \left(\sqrt{(1-\iota)(2+2c_0)(\log(d) + 1/\log^5(d))} - \sqrt{(2+2c_0)\log(d)} - O\left(\frac{1}{\log^9(d)}\right) \right) \\ &= \sigma_0 \left(\frac{(1-\iota)(2+2c_0)(\log(d) + 1/\log^5(d)) - (2+2c_0)\log(d)}{\sqrt{(1-\iota)(2+2c_0)(\log(d) + 1/\log^5(d))} + \sqrt{(2+2c_0)\log(d)}} - O\left(\frac{1}{\log^9(d)}\right) \right) \\ &= \sigma_0 \left(\frac{(2+2c_0)(-\iota\log(d) + (1-\iota)/\log^5(d))}{\sqrt{(1-\iota)(2+2c_0)(\log(d) + 1/\log^5(d))} + \sqrt{(2+2c_0)\log(d)}} - O\left(\frac{1}{\log^9(d)}\right) \right) \\ & > 0. \end{aligned} \quad (211)$$

This means all the \mathbf{v}_+ -singleton neurons will be updated on all the \mathbf{v}_+ -dominated patches at time $t = 0$. Therefore, we can write update expression equation 208 as follows

$$\begin{aligned} & \Delta \mathbf{w}_{+,c,r}^{(0)} \\ &= \frac{\eta}{NP} \sum_{n=1}^N \left\{ \mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_{+,c}^{(0)}(\mathbf{X}_n^{(0)})] + \sum_{c \in [k_+] - \{c\}} \mathbb{1}\{y_n = (+, c')\} [-\text{logit}_{+,c}^{(0)}(\mathbf{X}_n^{(0)})] \right\} \\ & \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(0)}; \mathbf{v}_+)} \left(\alpha_{n,p}^{(0)} \mathbf{v}_+ + \frac{(\cdot)}{n,p} \right). \end{aligned} \quad (212)$$

By concentration of the binomial random variable, we know that with probability at least $1 - e^{-\log^2(d)}$, for all n ,

$$\left| P(\mathbf{X}_n^{(0)}; \mathbf{v}_+) \right| = \left(1 \pm s^{*-1/3} \right) s^*. \quad (213)$$

Now, with the estimates we derived for $\text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})$ at the beginning of the proof and the independence of all the noise vectors $\frac{(\cdot)}{n,p}^{(0)}$'s, we arrive at

$$\Delta \mathbf{w}_{+,r}^{(0)} = [1/4, 2/3] \frac{1}{1 \pm \iota} \left(1 \pm s^{*-1/3} \right) \eta \frac{s^*}{2k_+P} \mathbf{v}_+ + \Delta_{+,r}^{(0)} \quad (214)$$

where $\sigma_{\zeta_{+,c,r}}^{(0)} = \Theta(1) \times \eta \sigma_\zeta \frac{\sqrt{s}}{P\sqrt{2N}}$.

Additionally, a byproduct of the above proof steps is that all the $S_{+,c}^*(\mathbf{v}_+)$ neurons indeed activate on all the \mathbf{v}_+ -dominated patches at $t = 0$ with high probability.

Now we examine the bias update. We first estimate $\left\| \Delta \mathbf{w}_{+,c,r}^{(0)} \right\|_2$. With probability at least $1 - O\left(\frac{m}{\text{poly}(d)}\right)$ the following upper bound holds for all neurons in $S_{+,c}^*(\mathbf{v}_+)$:

$$\begin{aligned} \left\| \Delta \mathbf{w}_{+,c,r}^{(0)} \right\|_2 & \leq O\left(\eta \frac{s^*}{k_+P}\right) \|\mathbf{v}_+\|_2 + \left\| \Delta_{+,r}^{(0)} \right\|_2 \\ & \leq O\left(\eta \frac{s^*}{k_+P}\right) + O\left(\eta \sigma_\zeta \frac{\sqrt{s}}{P} \frac{1}{N} \bar{d}\right) \\ & \leq O\left(\eta \frac{s^*}{k_+P}\right), \end{aligned} \quad (215)$$

and the following lower bound holds (via the reverse triangle inequality):

$$\begin{aligned} \left\| \Delta \mathbf{w}_{+,c,r}^{(0)} \right\|_2 & \geq \Omega\left(\eta \frac{s^*}{k_+P}\right) \|\mathbf{v}_+\|_2 - \left\| \Delta_{+,r}^{(0)} \right\|_2 \\ & \geq \Omega\left(\eta \frac{s^*}{k_+P}\right) - O\left(\eta \sigma_\zeta \frac{\sqrt{s}}{P} \frac{1}{N} \bar{d}\right) \\ & \geq \Omega\left(\eta \frac{s^*}{k_+P}\right), \end{aligned} \quad (216)$$

It follows that $\left\| \Delta \mathbf{w}_{+,c,r}^{(0)} \right\|_2 = \Theta\left(\eta \frac{s^*}{k_+P}\right)$, which means

$$\Delta b_{+,c,r}^{(0)} = -\frac{\left\| \Delta \mathbf{w}_{+,c,r}^{(0)} \right\|_2}{\log^5(d)} = -\Theta\left(\frac{\eta s^*}{k_+P \log^5(d)}\right). \quad (217)$$

This completes the proof of the base case.

Induction step. Assume statements for time $[0, t]$, prove for $t + 1$.

First, by the induction hypothesis, we know that neurons in $S_{+,c}^{*(0)}(\mathbf{v}_+)$ must activate on all the \mathbf{v}_+ -dominated patches at time t . Therefore, we can write the update expression equation 208 as follows:

$$\begin{aligned} & \Delta \mathbf{w}_{+,c,r}^{(t)} \\ &= \frac{\eta}{NP} \sum_{n=1}^N \left\{ \mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] + \sum_{c \in [k_+] - \{c\}} \mathbb{1}\{y_n = (+, c')\} [-\text{logit}_{+,c}^{(t)}(\mathbf{X}_n^{(t)})] \right\} \\ & \quad \times \sum_{p \in \mathcal{P}(\mathbf{X}_n^{(t)}; \mathbf{v}_+)} \left(\alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} \right). \end{aligned} \quad (218)$$

Following the same argument as in the base case, we have that with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$,

$$\Delta \mathbf{w}_{+,c,r}^{(t)} = [1/4, 2/3] \overline{1 \pm \iota} \left(1 \pm s^{*-1/3}\right) \eta \frac{s^*}{2k_+ P} \mathbf{v}_+ + \Delta_{+,c,r}^{(t)}, \quad (219)$$

and $\sigma_{\zeta_{+,c,r}}^{(t)} = \Theta(1) \times \eta \sigma_{\zeta} \frac{\sqrt{s}}{P\sqrt{2N}}$.

Now we need to show that $\mathbf{w}_{+,c,r}^{(t+1)}$ indeed activate on all the \mathbf{v}_+ -dominated patches at time $t + 1$ with high probability.

So far, we know that for $\tau \in [0, t + 1]$,

$$\Delta \mathbf{w}_{+,c,r}^{(\tau)} = [1/4, 2/3] \overline{1 \pm \iota} \left(1 \pm s^{*-1/3}\right) \eta \frac{s^*}{2k_+ P} \mathbf{v}_+ + \Delta_{+,c,r}^{(\tau)}, \quad (220)$$

and $\sigma_{\zeta_{+,c,r}}^{(\tau)} = \Theta(1) \times \eta \sigma_{\zeta} \frac{\sqrt{s}}{P\sqrt{2N}}$. It follows that

$$\mathbf{w}_{+,r}^{(t+1)} = \mathbf{w}_{+,c,r}^{(0)} + (t + 1) [1/4, 2/3] \overline{1 \pm \iota} \left(1 \pm s^{*-1/3}\right) \eta \frac{s^*}{2k_+ P} \mathbf{v}_+ + \Delta_{+,c,r}^{(t+1)}, \quad (221)$$

where $\sigma_{\zeta_{+,c,r}}^{(t+1)} = \Theta(1) \times \overline{t + 1} \eta \sigma_{\zeta} \frac{\sqrt{s}}{P\sqrt{2N}}$.

The following holds with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$ over all the \mathbf{v}_+ -dominated patches $\mathbf{X}_{n,p}^{(t+1)} = \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p}$ (which are independent of $\mathbf{w}_{+,r}^{(t+1)}$) and the \mathbf{v}_+ -singleton neurons:

$$\begin{aligned} & \mathbf{w}_{+,c,r}^{(t+1)}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p} \\ &= \mathbf{w}_{+,c,r}^{(0)}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p} + (t + 1) [1/4, 2/3] \overline{1 \pm \iota} \left(1 \pm s^{*-1/3}\right) \left(1 \pm O\left(\frac{1}{\log^9(d)}\right)\right) \eta \frac{s^*}{2k_+ P} \\ & \quad + \frac{(t+1)}{+,c,r}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p} \end{aligned} \quad (222)$$

Note that with probability at least $1 - O\left(\frac{1}{\text{poly}(d)}\right)$,

$$\frac{(t+1)}{+,c,r}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ = O(1) \times \overline{T} \eta \sigma_{\zeta} \frac{s^*}{P \cdot 2N} \sqrt{d \log(d)}, \quad (223)$$

and since $\overline{t + 1} \leq t + 1$, $\overline{s^*} < s^*$, $\sigma_{\zeta} \sqrt{d \log(d)} < \frac{1}{\log^9(d)}$, and $N > dk_+$, we know that

$$\frac{(t+1)}{+,c,r}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ = O\left(\frac{1}{d}\right) \times (t + 1) \eta \frac{s^*}{2k_+ P}. \quad (224)$$

Similarly, with probability at least $1 - O\left(\frac{1}{\text{poly}(d)}\right)$,

$$\mathcal{W}_{+,c,r}^{(t+1)}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p} \quad O(1) \times \frac{\bar{T} \eta \sigma_\zeta^2}{P} \frac{\bar{s}^*}{2N} \sqrt{d \log(d)} \quad O\left(\frac{1}{d}\right) \times (t+1) \eta \frac{s^*}{2k_+ P}. \quad (225)$$

It follows that with probability at least $1 - O\left(\frac{mNP}{\text{poly}(d)}\right)$,

$$\begin{aligned} & \mathcal{W}_{+,c,r}^{(t+1)}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p} \\ & \mathcal{W}_{+,c,r}^{(0)}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p} + \frac{1}{4}(t+1)(1-\iota) \left(1 - s^{*-1/3}\right) \left(1 - O\left(\frac{1}{\log^9(d)}\right)\right) \eta \frac{s^*}{2k_+ P}. \end{aligned} \quad (226)$$

Next, let us estimate the bias updates for $\tau \in [0, t+1]$.

Estimating $\Delta b_{+,c,r}^{(t)}$ follows an almost identical argument as in the base case (with the only main difference being relying on Theorem G.1 for non-activation on non- \mathbf{v}_+ -dominated patches), so we skip its calculations.

Therefore, $b_{+,c,r}^{(t+1)} = b_{+,c,r}^{(0)} - \Theta\left(\frac{\eta s}{k_+ P \log^5(d)}\right)$. This means

$$\begin{aligned} & \mathcal{W}_{+,c,r}^{(t+1)}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p} + b_{+,c,r}^{(t+1)} \\ & \mathcal{W}_{+,c,r}^{(0)}, \alpha_{n,p}^{(t+1)} \mathbf{v}_+ + \frac{(t+1)}{n,p} + b_{+,c,r}^{(0)} \\ & + \frac{1}{4}(t+1)(1-\iota) \left(1 - s^{*-1/3}\right) \left(1 - O\left(\frac{1}{\log^9(d)}\right)\right) \eta \frac{s^*}{2k_+ P} - O\left(\frac{\eta s^*(t+1)}{k_+ P \log^5(d)}\right) \end{aligned} \quad (227)$$

> 0 .

This completes the inductive step. □

Corollary G.2.1. *At time $t = T_0$, $\frac{\eta s}{k_+ P} \times s^* \left|S_{+,c}^{*(0)}(\mathbf{v}_+)\right|, \frac{\eta s}{k_+ P} \times s^* \left|S_{+,c}^{*(0)}(\mathbf{v}_{+,c})\right| = \Theta(1)$.*

Proof. Directly follows from Lemma G.2 and Theorem G.1. □

G.4 Model error after training

In this subsection, we show the model's error after fine-grained training. We also discuss that finetuning the model further increases its feature extractor's response to the true features, so it is even more robust/generalizing in downstream classification tasks.

Theorem G.3. *Define $\widehat{F}_+(\mathbf{X}) = \max_{c \in [k_+]} F_{+,c}(\mathbf{X})$, $\widehat{F}_-(\mathbf{X}) = \max_{c \in [k_-]} F_{-,c}(\mathbf{X})$.*

With probability at least $1 - O\left(\frac{mk_+^2 NPT_0}{\text{poly}(d)}\right)$, the following events take place:

1. *(Fine-grained easy & hard sample test accuracies are nearly perfect) Given an easy or hard fine-grained test sample (\mathbf{X}, y) where $y \in \{(+, c)\}_{c=1}^{k_+} \cup \{(-, c)\}_{c=1}^{k_-}$, $\mathbb{P}\left[F_y^{(T_0)}(\mathbf{X}) - \max_{y \neq y'} F_y^{(T_0)}(\mathbf{X})\right] = o(1)$.*
2. *(Coarse-grained easy & hard sample test accuracy are nearly perfect) Given an easy or hard coarse-grained test sample (\mathbf{X}, y) where $y \in \{+1, -1\}$, $\mathbb{P}\left[\widehat{F}_y^{(T_0)}(\mathbf{X}) - \widehat{F}_y^{(T_0)}(\mathbf{X})\right] = o(1)$.*

Proof. **Probability of mistake on easy samples.**

Without loss of generality, assume \mathbf{X} is a $(+, c)$ -class easy sample.

Conditioning on the events of Theorem G.1 and Lemma G.2, we know that for all $c' \in [k_-]$,

$$F_{-,c'}^{(T_0)} = O(m_{+,c} \sigma_0 \sqrt{\log(d)}) = o(1), \quad (228)$$

and for all $c' \in [k_+] - \{c\}$,

$$F_{+,c}^{(T_0)} = \sum_{p \in \mathcal{P}(X; \mathbf{v}_+)} \sum_{(+,r) \in S_{+,c}^{(0)}(\mathbf{v}_+)} \sigma \left(\mathbf{w}_{+,r}^{(T_0)}, \alpha_{n,p} \mathbf{v}_+ + \mathbf{n}_{+,p} + b_{+,c,r}^{(T_0)} \right) + O(m_{+,c} \sigma_0 \sqrt{\log(d)})$$

$$s^* \left| S_{+,c}^{(0)}(\mathbf{v}_+) \right| \frac{2}{3} (1 + \iota) \left(1 + s^{*-1/3} \right) \left(1 + \left(\frac{1}{\log^9(d)} \right) \right) \eta T_0 \frac{s^*}{2k_+ P}$$
(229)

moreover,

$$F_{+,c}^{(T_0)} = \sum_{p \in \mathcal{P}(X; \mathbf{v}_+)} \sum_{(+,r) \in S_{+,c}^{(0)}(\mathbf{v}_+)} \sigma \left(\mathbf{w}_{+,c,r}^{(T_0)}, \alpha_{n,p} \mathbf{v}_+ + \mathbf{n}_{+,p} + b_{+,c,r}^{(T_0)} \right)$$

$$+ \sum_{p \in \mathcal{P}(X; \mathbf{v}_{+,c})} \sum_{(+,r) \in S_{+,c}^{(0)}(\mathbf{v}_{+,c})} \sigma \left(\mathbf{w}_{+,c,r}^{(T_0)}, \alpha_{n,p} \mathbf{v}_{+,c} + \mathbf{n}_{+,p} + b_{+,c,r}^{(T_0)} \right)$$

$$s^* \left| S_{+,c}^{(0)}(\mathbf{v}_+) \right| \frac{1}{4} (1 - \iota) \left(1 - s^{*-1/3} \right) \left(1 - \left(\frac{1}{\log^5(d)} \right) \right) \eta T_0 \frac{s^*}{2k_+ P}$$

$$+ s^* \left| S_{+,c}^{(0)}(\mathbf{v}_{+,c}) \right| \left(1 - O \left(\frac{1}{k_+} \right) \right) (1 - \iota) \left(1 - s^{*-1/3} \right) \left(1 - \left(\frac{1}{\log^5(d)} \right) \right) \eta T_0 \frac{s^*}{2k_+ P}$$
(230)

Relying on Proposition 2, we know $\left| S_{+,c}^{(0)}(\mathbf{v}_+) \right| = \left(1 \pm \left(\frac{1}{\log^5(d)} \right) \right) \left| S_{+,c}^{*(0)}(\mathbf{v}_+) \right|$ and $\left| S_{+,c}^{*(0)}(\mathbf{v}_{+,c}) \right| = \left(1 \pm \left(\frac{1}{\log^5(d)} \right) \right) \left| S_{+,c}^{*(0)}(\mathbf{v}_+) \right|$, therefore $F_{+,c}^{(T_0)}(\mathbf{X}) > \max_{c \neq c'} F_{+,c'}^{(T_0)}(\mathbf{X})$ has to be true. With Corollary G.2.1, we also have $F_{+,c}^{(T_0)}(\mathbf{X}) \geq \Omega(1) > o(1) \max_{c \in [k_+]} F_{-,c}^{(T_0)}(\mathbf{X})$. It follows that the probability of mistake on an easy test sample is indeed at most $o(1)$.

Probability of mistake on hard samples. Without loss of generality, assume \mathbf{X} is a $(+,c)$ -class hard sample.

By Theorem G.1 (and its proof) and Lemma G.2, we know that for any $c' \in [k_+]$, the neurons $\mathbf{w}_{+,c',r}$ can only possibly receive update on \mathbf{v} -dominated patches for $\mathbf{v} \in \mathcal{U}_{+,c',r}^{(0)}$, and the updates to the neurons take the feature-plus-Gaussian-noise form of $\sum_{\mathbf{v} \in \mathcal{U}_{+,c',r}^{(0)}} c(\mathbf{v}) \mathbf{v} + \Delta_{+,c',r}^{(t)}$, with $c(\mathbf{v}) \in \left[\frac{2}{3}, \frac{1}{1+\iota} (1 + s^{*-1/3}) \eta \frac{s}{2k_+ P} \right]$ if \mathbf{v} is a fine-grained feature, or $c(\mathbf{v}) \in \left[\frac{2}{3}, \frac{1}{1+\iota} (1 + s^{*-1/3}) \eta \frac{s}{2k_+ P} \right]$ if $\mathbf{v} = \mathbf{v}_+$ (because the \mathbf{v} component of a \mathbf{v} -singleton neuron's update is already the maximum possible). Moreover, $\sigma_{\zeta_{+,c',r}}^{(t)} = O \left(\eta \sigma_{\zeta} \frac{\sqrt{s}}{P\sqrt{2N}} \right)$.

Relying on Theorem G.1, Lemma G.2, Corollary G.2.1 and previous observations, we have

$$F_{+,c}^{(T_0)}(\mathbf{X}) = \sum_{p \in \mathcal{P}(X; \mathbf{v}_{+,c})} \sum_{(+,c,r) \in S_{+,c}^{(0)}(\mathbf{v}_{+,c})} \sigma \left(\mathbf{w}_{+,c,r}^{(T_0)}, \alpha_{n,p} \mathbf{v}_{+,c} + \mathbf{n}_{+,p} + b_{+,c,r}^{(T_0)} \right)$$

$$s^* \left| S_{+,c}^{(0)}(\mathbf{v}_{+,c}) \right| \left(1 - O \left(\frac{1}{k_+} \right) \right) (1 - \iota) \left(1 - s^{*-1/3} \right) \left(1 - O \left(\frac{1}{\log^5(d)} \right) \right) \eta T_0 \frac{s^*}{2k_+ P}$$

$$\Omega(1),$$
(231)

and for $c' = c$,

$$\begin{aligned}
F_{+,c}^{(T_0)}(\mathbf{X}) &= \sum_{r=1}^{m_{+,c}} \sigma \left(\mathbf{w}_{+,c,r}^{(T_0)} * + b_{+,c,r}^{(T_0)} \right) \\
&+ \sum_{p \in \mathcal{P}(X; \mathbf{v}_{+,c})} \sum_{(+,c,r) \in S_{+,c}^{(0)}(\mathbf{v}_{+,c})} \sigma \left(\mathbf{w}_{+,c,r}^{(T_0)}, \alpha_{n,p} \mathbf{v}_{+,c} + \mathbf{v}_{+,c} + b_{+,c,r}^{(T_0)} \right) \\
&+ \sum_{p \in \mathcal{P}(X; \mathbf{v}_{-,c})} \sum_{(+,c,r) \in S_{+,c}^{(0)}(\mathbf{v}_{-,c})} \sigma \left(\mathbf{w}_{+,c,r}^{(T_0)}, \alpha_{n,p}^\dagger \mathbf{v}_{-,c} + \mathbf{v}_{+,c} + b_{+,c,r}^{(T_0)} \right) \\
&O(1) \times \left(\sum_{(+,c,r) \in \mathcal{U}_{+,c,r}^{(0)}} \sum_{\tau=0}^{T_0-1} \Delta \mathbf{w}_{+,c,r}^{(\tau)} * + \sum_{r \in [m_{+,c}]} \mathbf{w}_{+,c,r}^{(0)} * \right) \\
&+ \sum_{p \in \mathcal{P}(X; \mathbf{v}_{+,c})} \sum_{(+,c,r) \in S_{+,c}^{(0)}(\mathbf{v}_{+,c})} \sigma \left(\mathbf{w}_{+,c,r}^{(0)}, \alpha_{n,p} \mathbf{v}_{+,c} + \mathbf{v}_{+,c} + b_{+,c,r}^{(0)} \right) \\
&+ \sum_{p \in \mathcal{P}(X; \mathbf{v}_{-,c})} \sum_{(+,c,r) \in S_{+,c}^{(0)}(\mathbf{v}_{-,c})} \sigma \left(\mathbf{w}_{+,c,r}^{(0)}, \alpha_{n,p}^\dagger \mathbf{v}_{-,c} + \mathbf{v}_{+,c} + b_{+,c,r}^{(0)} \right) \\
&O\left(\frac{1}{\text{polylog}(d)}\right).
\end{aligned} \tag{232}$$

Moreover, for any $c' \in [k_-]$, similar to before,

$$\begin{aligned}
F_{-,c}^{(T_0)}(\mathbf{X}) &= \sum_{r=1}^{m_{-,c}} \sigma \left(\mathbf{w}_{-,c,r}^{(T_0)} * + b_{-,c,r}^{(T_0)} \right) \\
&+ \sum_{p \in \mathcal{P}(X; \mathbf{v}_{+,c})} \sum_{(-,c,r) \in S_{-,c}^{(0)}(\mathbf{v}_{+,c})} \sigma \left(\mathbf{w}_{-,c,r}^{(T_0)}, \alpha_{n,p} \mathbf{v}_{+,c} + \mathbf{v}_{-,c} + b_{-,c,r}^{(T_0)} \right) \\
&+ \sum_{p \in \mathcal{P}(X; \mathbf{v}_{-,c})} \sum_{(-,c,r) \in S_{-,c}^{(0)}(\mathbf{v}_{-,c})} \sigma \left(\mathbf{w}_{-,c,r}^{(T_0)}, \alpha_{n,p}^\dagger \mathbf{v}_{-,c} + \mathbf{v}_{-,c} + b_{-,c,r}^{(T_0)} \right) \\
&O(1) \times \left(\sum_{(-,c,r) \in \mathcal{U}_{-,c,r}^{(0)}} \mathbf{w}_{-,c,r}^{(T_0)} * + \sum_{r \in [m_{-,c}]} \mathbf{w}_{-,c,r}^{(0)} * \right) \\
&+ \sum_{p \in \mathcal{P}(X; \mathbf{v}_{+,c})} \sum_{(-,c,r) \in S_{-,c}^{(0)}(\mathbf{v}_{+,c})} \sigma \left(\mathbf{w}_{-,c,r}^{(0)}, \alpha_{n,p} \mathbf{v}_{+,c} + \mathbf{v}_{-,c} + b_{-,c,r}^{(0)} \right) \\
&+ O(1) \times s^\dagger \left| S_{-,c}^{(0)}(\mathbf{v}_{-,c}) \right| \times (\iota_{upper}^\dagger + O(\sigma_0 \log(d))) \\
&O\left(\frac{1}{\text{polylog}(d)}\right) + O\left(\sigma_0 \sqrt{\log(d)}\right) + O\left(\frac{1}{\log(d)}\right) \\
&o(1).
\end{aligned} \tag{233}$$

Therefore, $F_{+,c}^{(T_0)}(\mathbf{X}) > \max_{y \neq (+,c)} F_y^{(T_0)}(\mathbf{X})$, which means $\widehat{F}_+^{(T_0)}(\mathbf{X}) > \widehat{F}_-^{(T_0)}(\mathbf{X})$ indeed. \square

Remark. First of all, note that the feature extractor, after fine-grained training, is already well-performing, as it responds strongly ($\Omega(1)$ strength) to the true features, and very weakly ($o(1)$ strength) to any off-diagonal features and noise. In other words, we stop training when the margin is at least $\Omega(1)$, i.e. when we have $F_{y_n}^{(T)}(X_n^{(T)}) - \max_{y \neq y_n} F_y^{(T)}(X_n^{(T)}) \geq \Omega(1)$ for all n at some $T = \text{poly}(d)$, and with high probability, we just

need T_0 time to reach it. This can already help us explain the linear-probing result we saw on ImageNet21k in Appendix A.2, since linear probing does not alter the feature extractor after fine-grained pretraining (on ImageNet21k), it only retrains a new linear classifier on top of the feature extractor for classifying on the target ImageNet1k dataset.

At a high level, *finetuning* \widehat{F} can only further enhance the feature extractor’s response to the features, therefore making the model even more robust for challenging downstream classification problems; it will not degrade the feature extractor’s response to any true feature. A rigorous proof of this statement is almost a repetition of the proofs for fine-grained training, so we do not repeat them here. Intuitively speaking, we just need to note that the properties stated in Theorem G.1 will continue to hold during finetuning (as long as we stay in polynomial time), and with similar argument to those in the proof of Lemma G.2, we note that the neurons responsible for detecting fine-grained features, i.e. the $S_{+,c}^{*(0)}(\mathbf{v}_{+,c})$, will continue to only receive (positive) updates on the $\mathbf{v}_{+,c}$ -dominated patches of the following form:

$$\begin{aligned} \Delta \mathbf{w}_{+,c,r}^{(t)} &= \frac{\eta}{NP} \sum_{n=1}^N \mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\ &\quad \times \sum_{p \in \mathcal{P}(X_n^{(t)}; \mathbf{v}_{+,c})} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v}_{+,c} + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0 \} \left(\alpha_{n,p}^{(t)} \mathbf{v}_{+,c} + \frac{(t)}{n,p} \right), \end{aligned} \quad (234)$$

and similar update expression can be stated for the $S_{+,c}^{*(0)}(\mathbf{v}_+)$ neurons:

$$\begin{aligned} &\Delta \mathbf{w}_{+,c,r}^{(t)} \\ &= \frac{\eta}{NP} \sum_{n=1}^N \mathbb{1}\{y_n = (+, c)\} [1 - \text{logit}_+^{(t)}(\mathbf{X}_n^{(t)})] \\ &\quad \times \sum_{p \in \mathcal{P}(X_n^{(t)}; \mathbf{v}_+)} \mathbb{1}\{ \mathbf{w}_{+,c,r}^{(t)}, \alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} + b_{+,c,r}^{(t)} > 0 \} \left(\alpha_{n,p}^{(t)} \mathbf{v}_+ + \frac{(t)}{n,p} \right). \end{aligned} \quad (235)$$

Indeed, these feature-detector neurons will continue growing in the direction of the features they are responsible for detecting instead of degrade in strength.

H Probability Lemmas

Lemma H.1 (Laurent-Massart χ^2 Concentration (Laurent & Massart (2000) Lemma 1)). *Let $\mathbf{g} \sim N(\mathbf{0}, \mathbf{I}_d)$. For any vector $\mathbf{a} \in \mathbb{R}_{\geq 0}^d$, any $t > 0$, the following concentration inequality holds:*

$$\mathbb{P} \left[\sum_{i=1}^d a_i g_i^2 \geq \|\mathbf{a}\|_1 + 2 \|\mathbf{a}\|_2 \sqrt{t} + 2 \|\mathbf{a}\|_\infty t \right] \leq e^{-t} \quad (236)$$

Lemma H.2. *Let $\mathbf{g} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_d)$. Then,*

$$\mathbb{P} \left[\|\mathbf{g}\|_2^2 \geq 5\sigma^2 d \right] \leq e^{-d} \quad (237)$$

Proof. By Lemma H.1, setting $a_i = 1$ for all i and $t = d$ yields

$$\mathbb{P} \left[\|\mathbf{g}\|_2^2 \geq \sigma^2 d + 2\sigma^2 d + 2\sigma^2 d \right] \leq e^{-d} \quad (238)$$

□

Lemma H.3 (Shen et al. (2022a)). *Let $\mathbf{g}_1 \sim N(\mathbf{0}, \sigma_1^2 \mathbf{I}_d)$ and $\mathbf{g}_2 \sim N(\mathbf{0}, \sigma_2^2 \mathbf{I}_d)$ be independent. Then, for any $\delta \in (0, 1)$ and sufficiently large d , there exist constants c_1, c_2 such that*

$$\mathbb{P} \left[\|\mathbf{g}_1, \mathbf{g}_2\| \leq c_1 \sigma_1 \sigma_2 \sqrt{d \log(1/\delta)} \right] \geq 1 - \delta \quad (239)$$

$$\mathbb{P} \left[\|\mathbf{g}_1, \mathbf{g}_2\| \leq c_2 \sigma_1 \sigma_2 \sqrt{d} \right] \geq \frac{1}{4} \quad (240)$$