

---

# Unsupervised Transfer Learning via Adversarial Contrastive Training

---

Anonymous Authors<sup>1</sup>

## Abstract

Learning transferable data representations from abundant unlabeled data remains a critical yet challenging task in machine learning. While numerous self-supervised contrastive learning methods have emerged to address this challenge, a notable class of these approaches focuses on aligning the covariance or correlation matrix with the identity matrix. Despite their impressive performance across various downstream tasks, these methods often suffer from biased sample risk. This bias not only leads to significant optimization offsets, especially in mini-batch scenarios, but also complicates the development of theoretical frameworks. In this paper, we introduce Adversarial Contrastive Training (ACT), a novel unbiased self-supervised transfer learning approach. This method allows us to develop a comprehensive end-to-end theoretical analysis for self-supervised contrastive learning. Our theoretical results reveal that minimizing the loss function of ACT can lead to the downstream data distribution being clustered in the representation space, provided that the upstream unlabeled sample size is sufficient. As a result, even with a few downstream samples, ACT can achieve outstanding classification performance, offering valuable insights for few-shot learning. Furthermore, ACT demonstrates state-of-the-art classification performance across multiple benchmark datasets.

## 1. Introduction

Collecting unlabeled data is far more convenient and cost-effective than gathering labeled data in real-world applications. Consequently, learning representations from abundant unlabeled data presents a highly valuable yet challenging problem. The learned representations can be transferred to downstream tasks to enhance model performance or reduce the sample size required for those tasks.

Recently, self-supervised contrastive learning has emerged as a leading approach for learning representations from unlabeled data. This method aims to learn representations that are invariant to data augmentation. However, solely min-

imizing the distance between similar pairs leads to trivial solutions, known as model collapse. To address this issue, researchers have developed various strategies, broadly categorized into three types.

The first strategy treats augmented views of different images as negative pairs, ensuring their representations remain dissimilar (Ye et al., 2019; He et al., 2020; Chen et al., 2020a,b; HaoChen et al., 2021; Zhang et al., 2023). However, these methods require large batch sizes to ensure sufficient negative samples, leading to substantial computational and memory demands that may be prohibitive in many applications. Additionally, by treating augmented views of different images as negative pairs, these approaches fail to account for semantic similarities between distinct images, potentially forcing apart representations of conceptually related content. As pointed out by Chuang et al. (2020; 2022), this design can hurt the representation performance.

The second strategy prevents model collapse through asymmetric network architectures (Grill et al., 2020; Chen & He, 2021; Caron et al., 2020; 2021). Although eliminating the need for negative pairs, they exhibit significant sensitivity to architectural design choices, where minor modifications can lead to collapsed solutions (Grill et al., 2020; Chen & He, 2021). The specific architectural constraints may also limit the neural network’s approximation capabilities. Besides, these methods simultaneously introduce significant challenges for explanation.

The third strategy prevents model collapse by imposing a regularization term to align the covariance or correlation matrix with the identity matrix (Zbontar et al., 2021; Ermolov et al., 2021; Bardes et al., 2022; HaoChen et al., 2022; HaoChen & Ma, 2023; Huang et al., 2023), encouraging the separation of category centers. These methods do not require negative samples and also facilitate a clear theoretical understanding. Among them, a typical regularization term takes the form (Zbontar et al., 2021; HaoChen & Ma, 2023; Huang et al., 2023) as:

$$\mathcal{R}(f) = \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1)f(\mathbf{x}_2)^\top\} - I_{d^*} \right\|_F^2, \quad (1)$$

where  $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d^*}$  denotes the representation mapping from the original image space to the representation space,  $\|\cdot\|_F$  denotes the Frobenius norm,  $\mathbf{x}$  represents an original image, and  $\mathcal{A}(\mathbf{x})$  denotes the collection of all augmented

views of  $\mathbf{x}$ . The terms  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})$  indicate two augmented views independently and uniformly sampled from  $\mathcal{A}(\mathbf{x})$ , while  $I_{d^*}$  is the identity matrix with the same dimension  $d^*$  as the representation space.

The population risk defined in equation (1) is typically intractable. The following sample-level risk is used to estimate it (HaoChen et al., 2022; HaoChen & Ma, 2023; Zbontar et al., 2021):

$$\widehat{\mathcal{R}}(f) = \left\| \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_1^{(i)}) f(\mathbf{x}_2^{(i)})^\top - I_{d^*} \right\|_F^2, \quad (2)$$

where  $\{\mathbf{x}^{(i)}\}_{i \in [n]}$  denotes the original dataset, and  $\widetilde{D}_s = \{(\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)}) : \mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)} \in \mathcal{A}(\mathbf{x}^{(i)})\}_{i \in [n]}$  represents the augmented dataset for learning representations. Unfortunately, it is evident that  $\widehat{\mathcal{R}}(f)$  is a biased estimator of  $\mathcal{R}(f)$ , i.e.,  $\mathbb{E}_{\widetilde{D}_s} \{\widehat{\mathcal{R}}(f)\} \neq \mathcal{R}(f)$  due to the non-commutativity between the expectation and the Frobenius norm. Two significant challenges emerge due to this bias nature.

Firstly, the biased estimator (2) used in HaoChen et al. (2022); HaoChen & Ma (2023); Zbontar et al. (2021) introduces significant optimization deviations during the training procedure. Although theoretically  $\mathbb{E}_{\widetilde{D}_s} \{\widehat{\mathcal{R}}(f)\}$  converges to  $\mathcal{R}(f)$  as  $n$  approaches infinity, practical constraints necessitate the use of mini-batch samples to estimate  $\mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_2)^\top\}$ . In this regard, the bias leads to an offset in the optimization direction. Furthermore, this offset would compound across successive training iterations, as each gradient direction strongly depends on the previous one, ultimately resulting in a learned representation that may diverge significantly from the intended minimizer of the population risk in equation (1), as shown in Table 1.

Secondly, this inherent bias presents significant obstacles in establishing end-to-end theoretical guarantees. The development of such guarantees requires addressing three crucial aspects: *how does the downstream task error converge with respect to both the number of unlabeled samples in the source domain and labeled samples in the target domain, how does the abundance of unlabeled samples in self-supervised learning benefit downstream tasks, and why do self-supervised learning methods maintain their effectiveness even with limited downstream labeled data?*

Recent theoretical studies have significantly advanced our understanding of self-supervised learning. These studies can be categorized into two main lines of research. The first line (Garrido et al., 2022; HaoChen et al., 2022; Awasthi et al., 2022; Huang et al., 2023) focuses on analyzing the population risk of self-supervised learning methods. Consequently, these fundamental questions remain incompletely addressed due to the lack of discussion at the sample level. A comprehensive theoretical analysis requires bridging the

gap between population-level and sample-level risks, which is a challenging task due to the biasedness of methods such as Zbontar et al. (2021); HaoChen et al. (2022); HaoChen & Ma (2023).

The second line of theoretical research (Saunshi et al., 2019; HaoChen et al., 2021; Ash et al., 2022; Lei et al., 2023; HaoChen & Ma, 2023) studies the generalization error through Rademacher complexity while overlooking the approximation error. Since the learning performance is decided by the overall error, which is the summation of generalization error (evaluated by the Rademacher complexity) and approximation error, the error analysis yielded from the second line may be invalid.

In this study, we introduce **A**dversarial **C**ontrastive **T**raining (ACT), a novel unbiased approach to self-supervised learning. ACT implements an innovative iteration format that eliminates the bias between the population risk (1) and its sample-level counterpart. This advancement effectively addresses two critical challenges: the training deviation and the theoretical obstacle introduced by bias. Through comprehensive end-to-end analysis of ACT, we demonstrate how the number of unlabeled data in the self-supervised pre-training phase enhances downstream task performance. Specifically, we demonstrate that through representation learning using ACT, the downstream data can be clustered in the representation space, provided that the upstream unlabeled sample size is sufficient. As a result, even with a few downstream samples, ACT can achieve outstanding classification performance, offering valuable insights for few-shot learning.

## 1.1. Related Work

**Self-Supervised Loss** The loss function proposed by HaoChen et al. (2022) can be regarded as a special version of ACT with the constraint  $\mathbf{x}_1 = \mathbf{x}_2$ . The main difference between ACT and the approach by HaoChen et al. (2022) lies in the iteration format. As stated in Section 1, optimization deviation can accumulate with each iteration, particularly in the mini-batch scenario, while ACT employs adversarial training to mitigate this issue. The same problem is encountered by Zbontar et al. (2021), which can be loosely regarded as a biased sample version of (1).

**Self-Supervised Theory** Recent theoretical studies can be categorized into two main lines of research. The first line (Garrido et al., 2022; HaoChen et al., 2022; Awasthi et al., 2022; Huang et al., 2023) focuses on analyzing the population risk of self-supervised learning methods, which can not characterize how the error in downstream tasks diminishes with increasing sample size. The second line of research (Saunshi et al., 2019; HaoChen et al., 2021; Ash et al., 2022; Lei et al., 2023; HaoChen & Ma, 2023) studies

the generalization error through Rademacher complexity without the consideration of approximation error. However, the scarcity of approximation error makes the resulting error analysis ineffective. Specifically, ignoring the approximation error by simply supposing  $f$  belonging to a deep neural network class, the Rademacher complexity can be significantly reduced by controlling the scale of the network class, leading to impressive upper bounds. However, this controlled neural network class intuitively limits its approximation capacity. The increasing approximation error results in a larger overall error. Therefore, these studies cannot provide theoretical guidance for hypothesis class selection nor fully characterize the total error of self-supervised learning methods. In contrast, our work provides a comprehensive convergence analysis that characterizes how the downstream task error converges with respect to both the number of unlabeled samples in the source domain and labeled samples in the target domain.

## 1.2. Contributions

Our main contributions can be summarized as follows:

- We introduce Adversarial Contrastive Training (ACT), a novel self-supervised transfer learning method. This approach learns representations from unlabeled data by solving a min-max optimization problem that corrects the bias inherent in existing methods (HaoChen et al., 2022; Zbontar et al., 2021).
- Through extensive experiments, we demonstrate that ACT significantly outperforms traditional biased iterative methods (Table 1). Our empirical evaluation shows that ACT achieves state-of-the-art classification performance across multiple benchmark datasets using both fine-tuned linear probes and  $k$ -nearest neighbor ( $k$ -nn) protocols (Table 2).
- We establish comprehensive end-to-end theoretical guarantees for ACT in transfer learning scenarios under misspecified and overparameterized settings (Theorem 3.9). Our theoretical analysis demonstrates that ACT-learned representations that minimaxing the loss function of ACT can lead to the downstream data distribution being clustered by category in the representation space, provided that the upstream unlabeled sample size is sufficient. Hence, even with a few downstream samples, ACT can achieve outstanding classification performance, offering valuable insights for few-shot learning.

## 1.3. Preliminaries

Given an integer  $n \in \mathbb{N}$ , we use  $[n]$  to represent the integer set  $\{1, 2, \dots, n\}$ . For any vector  $v$ , we denote

$\|v\|_2$  and  $\|v\|_\infty$  as the 2-norm and  $\infty$ -norm of  $v$  respectively. Let  $A, B \in \mathbb{R}^{d_1 \times d_2}$  be two matrices, we denote their Frobenius inner product by  $\langle A, B \rangle_F = \text{tr}(A^T B)$ . Moreover, we denote  $\|A\|_F$  as the Frobenius norm of  $A$ , which is the norm induced by Frobenius inner product, and  $\|A\|_\infty = \sup_{\|x\|_\infty \leq 1} \|Ax\|_\infty$  as the  $\infty$ -norm of  $A$ , which is the maximum 1-norm of the rows of  $A$ . For a given map  $f$  and  $0 \leq a_1 \leq a_2$ , we use  $a_1 \leq \|f\|_2 \leq a_2$  to denote  $b_1 \leq \inf_v \|f(v)\|_2 \leq \sup_v \|f(v)\|_2 \leq b_2$ . Besides that, the Lipschitz norm of  $f$  is given by  $\|f\|_{\text{Lip}} = \sup_{u \neq v} \frac{\|f(u) - f(v)\|_2}{\|u - v\|_2}$ . Furthermore, for a given function  $f : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$ , we use  $f \in \text{Lip}(L)$  to represent  $\|f\|_{\text{Lip}} \leq L$ . For ease of presentation, throughout this paper, we use  $X \lesssim Y$  or  $Y \gtrsim X$  to denote the statement that  $X \leq CY$  for two quantities  $X$  and  $Y$ , where  $C > 0$  can be arbitrary constant.

We will adopt the following ReLU neural network class as the hypothesis space in the subsequent content.

**Definition 1.1** (ReLU neural network class). Given  $0 < d_1, d_2; L, N_1, \dots, N_L \in \mathbb{N}; 0 < \mathcal{K}$  and  $0 < B_1 \leq B_2$ , define  $W = \max\{N_1, \dots, N_L\}$ , a deep ReLU network class with parameter  $(W, L, \mathcal{K}, B_1, B_2)$ ,  $\mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}, B_1, B_2)$ , is defined as the collection of all maps of the form

$$f_\theta(x) = A_L \sigma(A_{L-1} \sigma(\dots \sigma(A_0 x + b_0))) + b_{L-1}$$

such that  $B_1 \leq \|f_\theta\|_2 \leq B_2$  and  $\kappa(\theta) \leq \mathcal{K}$ , where  $\sigma(x) = \max\{0, x\}$  is the ReLU activate function,  $N_0 = d_1, N_{L+1} = d_2, A_i \in \mathbb{R}^{N_{i+1} \times N_i}$  and  $b_i \in \mathbb{R}^{N_{i+1}}$ . The integers  $W$  and  $L$  are called the width and depth of the neural network respectively. The parameters set of the neural network is defined as  $\theta := ((A_0, b_0), \dots, (A_{L-1}, b_{L-1}), A_L)$ . Further,  $\kappa(\theta)$  is defined as  $\kappa(\theta) = \|A_L\|_\infty \prod_{l=0}^{L-1} \max\{\|(A_l, b_l)\|_\infty, 1\}$ .

For any  $f_\theta \in \mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}, B_1, B_2)$ , we can justify  $\|f_\theta\|_{\text{Lip}} \leq \mathcal{K}$ . The proof details are deferred to Appendix A.1.

Besides that, for any two measures  $\mu$  and  $\nu$ , we define the 1-Wasserstein distance as  $\mathcal{W}(\mu, \nu) = \max_{g \in \text{Lip}(1)} \mathbb{E}_{X \sim \mu}\{g(X)\} - \mathbb{E}_{Y \sim \nu}\{g(Y)\}$ .

## 1.4. Organization

This paper is structured as follows: Section 2 introduces the core concept of ACT and presents our alternating optimization algorithm. In Section 3, we develop a comprehensive end-to-end theoretical guarantee for ACT. Section 4 demonstrates ACT's effectiveness through extensive experimental evaluations across diverse datasets and metrics. Section 5 concludes with a summary of our findings. All detailed proofs are provided in Section A.

## 2. Adversarial Contrastive Training

### 2.1. Notations for Unsupervised Transfer Learning

Throughout this paper, we use  $d$  and  $d^*$  to represent the dimensions of the original image and the representation dimension, respectively. We denote image instances from the source domain  $\mathcal{X}_s \subseteq [0, 1]^d$  with distribution  $\mathbb{P}_s$  using the letter  $\mathbf{x}$  and its subscripted or superscripted variants. In contrast, we use the letter  $\mathbf{z}$  and its subscripted or superscripted variants for image instances from the target domain  $\mathcal{X}_t \subseteq [0, 1]^d$  with distribution  $\mathbb{P}_t$ . In this context, we can independently and identically sample a total of  $n_s$  source image instances from  $\mathbb{P}_s$  and  $n_t$  downstream samples from  $\mathbb{P}_t$ . Notably, the label for each  $\mathbf{z}^{(i)} \sim \mathbb{P}_t$  is observable. We refer to these two datasets as  $D_s = \{\mathbf{x}^{(i)}\}_{i \in [n_s]}$  and  $D_t = \{(\mathbf{z}^{(i)}, y_i)\}_{i \in [n_t]}$ , respectively.

Since the primary objective of contrastive learning is to learn a representation that is invariant to augmentations, data augmentation plays a crucial role in this area. A data augmentation  $A : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is essentially a predefined transformation applied to original images. Common augmentations include a composition of random transformations, such as Random-Crop, HorizontalFlip, and Color Distortion (Chen et al., 2020a). We refer to the collection of used data augmentations as  $\mathcal{A} = \{A_i(\cdot)\}_{i \in [m]}$  as the collection of used data augmentations, where  $m$  is the total number of data augmentation under consideration. Theoretically,  $m$  could be infinite. But we might consider only a finite but sufficiently large  $m$  for convenient theoretical treatment. In fact, as long as  $m$  is sufficiently large, essentially any type of data augmentation might be well approximated by some  $A \in \mathcal{A}$ . Base on  $\mathcal{A}$ , we can construct an augmented dataset  $\tilde{D}_s = \{\tilde{\mathbf{x}}^{(i)}\}_{i \in [n_s]}$ , where  $\tilde{\mathbf{x}}^{(i)} = (\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)}) = (A_{i,1}(\mathbf{x}^{(i)}), A_{i,2}(\mathbf{x}^{(i)}))$ , and  $A_{i,1}$  and  $A_{i,2}$  are independently drawn from the uniform distribution on  $\mathcal{A}$ .

### 2.2. Adversarial Contrastive Training

We begin by recalling  $\mathcal{R}(f)$  defined in (1), which is the regularization term adopted by various studies (HaoChen et al., 2022; HaoChen & Ma, 2023; Huang et al., 2023) to prevent model collapse. Its empirical version at the sample level is given by  $\widehat{\mathcal{R}}(f)$ . However, as stated in Section 1,  $\widehat{\mathcal{R}}(f)$  is a biased counterpart of  $\mathcal{R}(f)$ , i.e.,  $\mathbb{E}_{\tilde{D}_s} \{\widehat{\mathcal{R}}(f)\} \neq \mathcal{R}(f)$ , which hinders establishing a theoretical foundation at the sample level and introduces optimization deviation.

To address these two issues, we then propose a novel sample-level estimator for the population risk (1). A key observation to motivate ACT is that we can rewrite  $\mathcal{R}(f)$  as

$$\mathcal{R}(f) = \sup_{G \in \mathcal{G}(f)} \mathcal{R}(f, G), \quad (3)$$

where  $G \in \mathbb{R}^{d^* \times d^*}$  is a matrix variable, and

$$\begin{aligned} \mathcal{R}(f, G) &= \langle \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1)f(\mathbf{x}_2)^\top\} - I_{d^*}, G \rangle_F, \\ \mathcal{G}(f) &= \{G \in \mathbb{R}^{d^* \times d^*} : \|G\|_F \leq \sqrt{\mathcal{R}(f)}\}. \end{aligned}$$

The equation (3) holds because of the fact that  $\langle A, B \rangle_F \leq \|A\|_F \|B\|_F$  for any matrices  $A, B$  of same dimension, with equality holding if and only if  $A = B$ . Correspondingly, the sample-level counterpart associated with (3) is given by

$$\widehat{\mathcal{R}}(f) = \sup_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{R}}(f, G),$$

where

$$\begin{aligned} \widehat{\mathcal{R}}(f, G) &= \left\langle \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{x}_1^{(i)})f(\mathbf{x}_2^{(i)})^\top - I_{d^*}, G \right\rangle_F, \\ \widehat{\mathcal{G}}(f) &= \left\{ G \in \mathbb{R}^{d^* \times d^*} : \|G\|_F \leq \sqrt{\widehat{\mathcal{R}}(f)} \right\}. \end{aligned}$$

It can be shown,

$$\mathcal{R}(f, G) = \mathbb{E}_{\tilde{D}_s} \{\widehat{\mathcal{R}}(f, G)\}.$$

Hence, the equivalent transformation (3) help us avoid the issue of biasedness. Specifically, with the equivalent transformation (3) and its empirical version, we learn the contrastive representation through the Adversarial Contrastive Training (ACT) at the sample level, which can be formulated as a mini-max problem as follows:

$$\begin{aligned} \hat{f}_{n_s} &\in \arg \min_{f \in \mathcal{F}} \max_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{L}}(f, G), \quad (4) \\ \widehat{\mathcal{L}}(f, G) &= \widehat{\mathcal{L}}_{\text{align}}(f) + \lambda \widehat{\mathcal{R}}(f, G), \\ \widehat{\mathcal{L}}_{\text{align}}(f) &= \frac{1}{n_s} \sum_{i=1}^{n_s} \|f(\mathbf{x}_1^{(i)}) - f(\mathbf{x}_2^{(i)})\|_2^2, \end{aligned}$$

where  $\mathcal{F}$  is defined as  $\mathcal{N}\mathcal{N}(W, L, \mathcal{K}, B_1, B_2)$ . We will specify the appropriate parameters  $(W, L, \mathcal{K}, B_1, B_2)$  to satisfy the theoretical requirements in Section 3. The term  $\widehat{\mathcal{L}}_{\text{align}}(f)$  embodies the core idea of contrastive learning: learning a representation that is invariant to augmentations. Additionally,  $\lambda > 0$  serves as the regularization hyperparameter.

This mini-max problem naturally leads to an alternative optimization algorithm for solving it, where  $G$  is fixed during the optimization of the encoder  $f$  and  $f$  is fixed when optimizing  $G$ . We present this algorithm in Algorithm 1. It is important to note that  $G_t$  has been detached from the computational graph when updating the encoder parameters  $\theta$ . This detachment implies that the gradient with respect to  $\theta$  is as given by the seventh line of Algorithm 1, rather than  $\nabla_{\theta} \left\| \frac{1}{N} \sum_{i=1}^N f_{\theta}(\mathbf{x}_1^{(n_i^t)}) f_{\theta}(\mathbf{x}_2^{(n_i^t)})^\top - I_{d^*} \right\|_F^2$ , which is the mini-batch gradient of  $\widehat{\mathcal{R}}(f)$ . In this regard, such a

**Algorithm 1** Alternative Optimization Algorithm

**Require:** Augmented dataset  $D_s = \{\tilde{\mathbf{x}}^{(i)}\}_{i \in [n]}$ , initial encoder parameter  $\theta_0$ , iteration horizon  $T$ , mini-batch size  $N$ , learning rate  $\eta$ .

- 1: **for**  $t \in \{0\} \cup [T - 1]$  **do**
- 2:   Sample a mini-batch  $\mathcal{B}_t = \{\mathbf{x}^{(n_i^t)}\}_{i \in N} \subseteq D_s$  of size  $N$ , where  $n_i^t$  represents the index of the  $i$ -th sample in the mini-batch  $\mathcal{B}_t$  within  $D_s$ .
- 3:   **if**  $t = 0$  **then**
- 4:      $G_0 = \sum_{i=1}^N f_{\theta_0}(\mathbf{x}_1^{(n_i^t)}) f_{\theta_0}(\mathbf{x}_2^{(n_i^t)})^\top - I_{d^*}$ .
- 5:     Detach:  $G_0 \leftarrow G_0.\text{detach}()$ .
- 6:   **end if**
- 7:   Update encoder  $\theta_{t+1} = \theta_t - \eta \Delta\theta$ , where  $\Delta\theta = \nabla_{\theta} \frac{1}{N} \sum_{i=1}^N \|f_{\theta}(\mathbf{x}_1^{(n_i^t)}) - f_{\theta}(\mathbf{x}_2^{(n_i^t)})\|_2^2 + \langle \nabla_{\theta} \frac{1}{N} \sum_{i=1}^N f_{\theta}(\mathbf{x}_1^{(n_i^t)}) f_{\theta}(\mathbf{x}_2^{(n_i^t)})^\top - I_{d^*}, G_t \rangle_F$ .
- 8:    $G_{t+1} = \sum_{i=1}^N f_{\theta_{t+1}}(\mathbf{x}_1^{(n_i^t)}) f_{\theta_{t+1}}(\mathbf{x}_2^{(n_i^t)})^\top - I_{d^*}$ .
- 9:   Detach:  $G_{t+1} \leftarrow G_{t+1}.\text{detach}()$ .
- 10: **end for**

**output** The learned encoder  $f_{\theta_T}$ .

mini-max iteration format will yield a distinctly different encoder in the mini-batch scenario compared to previous studies (Zbontar et al., 2021; HaoChen et al., 2022).

We compare ACT against two biased self-supervised learning methods: Barlow Twins (Zbontar et al., 2021) and the approach proposed by HaoChen et al. (2022), across multiple benchmark datasets. The experimental results, summarized in Table 1, demonstrate that ACT significantly improves downstream classification accuracy compared to both baseline methods, which are implemented using our repository, with a total training of 1000 epochs and a representation dimension of 512. While, ACT employs representation dimensions of 64, 64, and 128, which are significantly lower than those of Zbontar et al. (2021); HaoChen et al. (2022); yet, it achieved the most outstanding performance.

Table 1. Classification accuracy (top 1) of a linear classifier and a 5-nearest neighbors classifier for different methods and datasets. Here BT indicates Barlow Twins (Zbontar et al., 2021) while BS refers to the method proposed by HaoChen et al. (2022).

Method	CIFAR-10		CIFAR-100		Tiny ImageNet	
	Linear	$k$ -nn	Linear	$k$ -nn	Linear	$k$ -nn
BT	83.96	81.18	56.75	47.91	34.08	19.40
BS	86.95	82.83	53.75	48.40	35.80	20.36
ACT	<b>92.11</b>	<b>90.01</b>	<b>68.24</b>	<b>58.35</b>	<b>49.72</b>	<b>36.40</b>

### 3. End-to-End Theoretical Guarantee

#### 3.1. Problem Formulation

We first define the ideal  $f^*$  for  $f$  as the minimizer at the population level, which represents the ideal objective of ACT.

$$\begin{aligned}
 f^* &\in \arg \min_{f: B_1 \leq \|f\|_2 \leq B_2} \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G), \\
 \mathcal{L}(f, G) &= \mathcal{L}_{\text{align}}(f) + \lambda \mathcal{R}(f, G), \\
 \mathcal{L}_{\text{align}}(f) &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \left\{ \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2 \right\}.
 \end{aligned}$$

Apart from that, we further denote

$$\mathcal{L}(f) = \mathcal{L}_{\text{align}}(f) + \lambda \mathcal{R}(f).$$

Intuitively, in data representation, the most critical aspect is the differentiation between various features, rather than the specific ranges of their values. Therefore, the constraint  $B_1 \leq \|f\|_2 \leq B_2$  will not diminish the performance of the encoder; instead, it facilitates the establishment of theoretical foundations for ACT.

Moreover, following a similar process to that used for obtaining  $\tilde{D}_s$ , we can construct the downstream augmented dataset  $\tilde{D}_t = \{(\tilde{\mathbf{z}}^{(i)}, y_i)\}_{i \in [n_t]}$ , where  $\tilde{\mathbf{z}}^{(i)} = \{(\mathbf{z}_1^{(i)}, \mathbf{z}_2^{(i)})\}_{i \in [n_t]}$  with  $\mathbf{z}_1^{(i)} = A_{i,1}(\mathbf{z}^{(i)})$ ,  $\mathbf{z}_2^{(i)} = A_{i,2}(\mathbf{z}^{(i)})$ . Therein,  $A_{i,1}$ ,  $A_{i,2}$  are independently and identically distributed samples drawn from the uniform distribution defined on  $\mathcal{A}$ . In this context, we construct the following linear probe as a classifier:

$$Q_{\hat{f}_{n_s}}(\mathbf{z}) = \arg \max_{k \in [K]} (\widehat{W} \hat{f}_{n_s}(\mathbf{z}))_k, \quad (5)$$

where the  $k$ -th row of  $\widehat{W}$  is given as  $\hat{\mu}_t(k) = \frac{1}{2n_t(k)} \sum_{i=1}^{n_t} (\hat{f}_{n_s}(\mathbf{z}_1^{(i)}) + \hat{f}_{n_s}(\mathbf{z}_2^{(i)})) \mathbb{1}\{y_i = k\}$ , therein,  $n_t(k) = \sum_{i=1}^{n_t} \mathbb{1}\{y_i = k\}$ . The classifier defined in (5) indicates that by calculating the average representations for each class, we build a template for each downstream class individually. Whenever a new sample needs to be classified, it is assigned to the category of the template that it most closely resembles. Furthermore, we use the following misclassification rate to evaluate the representation learned by ACT.

$$\text{Err}(Q_{\hat{f}_{n_s}}) = \sum_{k=1}^K \mathbb{P}_t \{Q_{\hat{f}_{n_s}}(\mathbf{z}) \neq k, \mathbf{z} \in C_t(k)\}, \quad (6)$$

where  $C_t(k)$  is a set such that  $\mathbf{z} \in C_t(k)$  if and only if  $\mathbf{z}$  belongs to the  $k$ -th class. Correspondingly, similar to Huang et al. (2023), we assume that any upstream instance  $\mathbf{x}$  can be categorized into one or more latent classes  $\{C_s(k)\}_{k \in [K]}$ . For ease of presentation, let  $p_s(k) = \mathbb{P}_s\{\mathbf{x} \in C_s(k)\}$  and  $\mathbb{P}_s(k)(\cdot) = \mathbb{P}_s\{\cdot | \mathbf{x} \in C_s(k)\}$ . Similarly, let  $p_t(k) =$

$\mathbb{P}_t\{\mathbf{z} \in C_t(k)\}$  and  $\mathbb{P}_t(k)(\cdot) = \mathbb{P}_t(\cdot|\mathbf{z} \in C_t(k))$ . In this context, we use the quantities

$$\begin{aligned} \epsilon_1 &= \max_{k \in [K]} \mathcal{W}(\mathbb{P}_s(k), \mathbb{P}_t(k)), \\ \epsilon_2 &= \max_{k \in [K]} |p_s(k) - p_t(k)|, \end{aligned} \quad (7)$$

to measure the divergence between the source and the target domains, where  $\mathcal{W}$  denotes 1-Wasserstein distance.

### 3.2. Theoretical Limitation of Bias

In this section, we aim to elucidate the limitations imposed by bias from a theoretical perspective. We first assert that  $\mathbb{E}_{\tilde{D}_s}\{\text{Err}(Q_{\hat{f}_{n_s}})\} \lesssim \mathbb{E}_{\tilde{D}_s}\{\mathcal{L}(\hat{f}_{n_s})\}$  under specific conditions, the details of which can be found in Section A.2. Consequently, analyzing the sample complexity of  $\mathbb{E}_{\tilde{D}_s}\{\mathcal{L}(\hat{f}_{n_s})\}$  is essential to establish an end-to-end theory for ACT. However, this analysis poses a significant challenge due to the presence of bias.

In fact, in the field of learning theory, the condition  $\mathbb{E}_{\tilde{D}_s}\{\hat{\mathcal{L}}(f)\} = \mathcal{L}(f)$  is quite important to establish the upper bound of  $\mathbb{E}_{\tilde{D}_s}\{\mathcal{L}(\hat{f}_{n_s})\}$ . Specifically, let  $\bar{f}$  satisfy  $\mathcal{L}(\bar{f}) - \mathcal{L}(f^*) = \inf_{f \in \mathcal{F}}\{\mathcal{L}(f) - \mathcal{L}(f^*)\}$ , then

$$\begin{aligned} \mathcal{L}(\hat{f}_{n_s}) &= \mathcal{L}(\hat{f}_{n_s}) - \hat{\mathcal{L}}(\hat{f}_{n_s}) + \hat{\mathcal{L}}(\hat{f}_{n_s}) - \hat{\mathcal{L}}(f^*) + \hat{\mathcal{L}}(f^*) \\ &\quad - \mathcal{L}(f^*) + \mathcal{L}(f^*) \\ &\leq \mathcal{L}(f^*) + 2 \sup_{f \in \mathcal{F}} |\mathcal{L}(f) - \hat{\mathcal{L}}(f)| + \{\hat{\mathcal{L}}(\hat{f}_{n_s}) - \hat{\mathcal{L}}(f^*)\} \\ &\leq \mathcal{L}(f^*) + 2 \sup_{f \in \mathcal{F}} |\mathcal{L}(f) - \hat{\mathcal{L}}(f)| + \{\hat{\mathcal{L}}(\bar{f}) - \hat{\mathcal{L}}(f^*)\}. \end{aligned}$$

Taking the expectation regarding to  $\tilde{D}_s$  on both sides yields  $\mathbb{E}_{\tilde{D}_s}\{\mathcal{L}(\hat{f}_{n_s})\} \leq \mathcal{L}(f^*) + 2\mathbb{E}_{\tilde{D}_s}\{\sup_{f \in \mathcal{F}} |\mathcal{L}(f) - \hat{\mathcal{L}}(f)|\} + \inf_{f \in \mathcal{F}}\{\mathcal{L}(f) - \mathcal{L}(f^*)\}$ . As observed, the first term is a typical problem in the area of empirical process, where the sample bound also require unbiasedness; not to mention that such typical risk decomposition itself necessitates a guarantee of unbiasedness. In contrast, based on the modification of ACT, we develop an novel error decomposition as follows:

$$\begin{aligned} \mathbb{E}_{\tilde{D}_s}[\mathcal{L}(\hat{f}_{n_s})] &\lesssim \mathcal{L}(f^*) + \inf_{f \in \mathcal{F}}\{\mathcal{L}(f) - \mathcal{L}(f^*)\} \\ &\quad + \mathbb{E}_{\tilde{D}_s}\left\{\sup_{f \in \mathcal{F}, G \in \hat{\mathcal{G}}(f)} |\mathcal{L}(f, G) - \hat{\mathcal{L}}(f, G)|\right\} \\ &\quad + \mathbb{E}_{\tilde{D}_s}\left[\sup_{f \in \mathcal{F}}\{G^*(f) - \hat{G}(f)\}\right] \end{aligned} \quad (8)$$

where  $G^*(f) = \mathbb{E}_{\mathbf{x}}\mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})}\{f(\mathbf{x}_1)f(\mathbf{x}_2)^\top - I_{d^*}\}$  and  $\hat{G}(f) = \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{x}_1^{(i)})f(\mathbf{x}_2^{(i)})^\top - I_{d^*}$ . The proof details can be found in Section A.3.3. Through this decomposition in (8), we can systematically analysis  $\mathbb{E}_{\tilde{D}_s}[\mathcal{L}(\hat{f}_{n_s})]$ . Particularly, the first term in (8) can be bounded under Assumption 3.3, as will be demonstrated subsequently. The

second term, known as approximation error, represents the error introduced by using  $\mathcal{F}$  to approximate  $f^*$ . Utilizing the unbiasedness of  $\hat{\mathcal{L}}(f, G)$ , the third term can be bounded using standard techniques from empirical process theory, while the last term can be reformulated as a common problem regarding the rate of convergence of the empirical mean to the population mean.

### 3.3. Assumptions

We begin with introducing the Hölder class, which plays a curial role in bounding the approximation error, i.e., the second term in (8).

**Definition 3.1** (Hölder class). Let  $d \in \mathbb{N}$  and  $\alpha = r + \beta > 0$ , where  $r \in \mathbb{N}_0$  and  $\beta \in (0, 1]$ . We assert  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  belongs to the Hölder class  $\mathcal{H}^\alpha(\mathbb{R}^d)$  if and only if

$$|\partial^s f(\mathbf{x})| \leq 1 \text{ and } \max_{\|\mathbf{s}\|_1=r} \sup_{\mathbf{x} \neq \mathbf{y}} \frac{\partial^s f(\mathbf{x}) - \partial^s f(\mathbf{y})}{\|\mathbf{x} - \mathbf{y}\|_\infty^\beta} \leq 1,$$

where for a multi-index  $\mathbf{s} = (s_1, \dots, s_d) \in \mathbb{N}_0^d$  and  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , the symbol  $\partial^s f$  denotes the partial differential operator  $\partial^s = \frac{\partial^{s_1}}{\partial x_1^{s_1}} \frac{\partial^{s_2}}{\partial x_2^{s_2}} \dots \frac{\partial^{s_d}}{\partial x_d^{s_d}}$ . Furthermore, we define  $\mathcal{H}^\alpha := \{f : [0, 1]^d \rightarrow \mathbb{R}, f \in \mathcal{H}^\alpha(\mathbb{R}^d)\}$  as the restriction of  $\mathcal{H}^\alpha(\mathbb{R}^d)$  to  $[0, 1]^d$ .

The Hölder class is known to be a highly comprehensive functional class, providing a precise characterization of the low-order regularity of functions. In this regard, we make following assumption:

**Assumption 3.2.** There exists  $\alpha = r + \beta$  with  $r \in \mathbb{N}_0$  and  $\beta \in (0, 1]$  s.t.  $f_i^* \in \mathcal{H}^\alpha$  for each  $i \in [d^*]$ .

Assumption 3.2 is standard and mild in the context of nonparametric statistics (Tsybakov, 2008; Schmidt-Hieber, 2020) due to the universality of the Hölder class.

As for the term  $\mathcal{L}(f^*)$  in eq (8), we make following Assumption 3.3 to ensure  $\mathcal{L}(f^*) = 0$ .

**Assumption 3.3.** Assume there exists a measurable partition  $\{\mathcal{P}_1, \dots, \mathcal{P}_{d^*}\}$  of  $\mathcal{X}_s$ , such that  $1/B_2^2 \leq \mathbb{P}_s(\mathcal{P}_i) \leq 1/B_1^2$  for each  $i \in [d^*]$ .

Assumption 3.3 suggests that the data distribution in the source domain should not be overly singular. All common continuous distributions defined on Borel algebra satisfy these requirements, as the measure of any single point is zero. More details are deferred to Section A.3.2.

*Remark 3.4.* (HaoChen & Ma, 2023, Assumption 4.2) assumes that the term  $\mathcal{L}(f)$  can be sufficiently minimized by a specific network. Constructing a network  $f_\theta \in \mathcal{F}$  such that population statistic  $\mathcal{L}(f_\theta)$  is sufficiently small is too complex. In contrast, we consider a more general setting where  $f^*$  may not belong to  $\mathcal{F}$ . Based on the mild Assumption 3.3, we can theoretically illustrate that  $\mathcal{L}(f^*)$  vanishes. This is crucial for subsequent theoretical analysis.

Additionally, we need to introduce two assumptions regarding to the data augmentation.

**Assumption 3.5.** Assume any data augmentation  $A_i \in \mathcal{A}$  is  $M$ -Lipschitz map, i.e.,  $\|A_i(\mathbf{v}_1) - A_i(\mathbf{v}_2)\|_2 \leq M\|\mathbf{v}_1 - \mathbf{v}_2\|_2$  for any  $\mathbf{v}_1, \mathbf{v}_2 \in [0, 1]^d$ .

A typical example to illustrate Assumption 3.5 is that the augmented view yielded by cropping should not undergo drastic changes when minor perturbations are applied to the original image.

In addition to the Lipschitz property of data augmentation, we adopt Definition 3.6 to mathematically quantify the quality of data augmentations.

**Definition 3.6** ( $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -Augmentation). The augmentations in  $\mathcal{A}$  is  $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -augmentations, that is, for each  $k \in [K]$ , there exists a subset  $\tilde{C}_s(k) \subseteq C_s(k)$  and  $\tilde{C}_t(k) \subseteq C_t(k)$ , such that (i)  $\mathbb{P}_s\{\mathbf{x} \in \tilde{C}_s(k)\} \geq \sigma_s \mathbb{P}_s\{\mathbf{x} \in C_s(k)\}$ , (ii)  $\sup_{\mathbf{x}_1, \mathbf{x}_2 \in \tilde{C}_s(k)} \min_{\mathbf{x}'_1 \in \mathcal{A}(\mathbf{x}_1), \mathbf{x}'_2 \in \mathcal{A}(\mathbf{x}_2)} \|\mathbf{x}'_1 - \mathbf{x}'_2\|_2 \leq \delta_s$ ; (iii)  $\mathbb{P}_t\{\mathbf{z} \in \tilde{C}_t(k)\} \geq \sigma_t \mathbb{P}_t\{\mathbf{z} \in C_t(k)\}$ , (iv)  $\sup_{\mathbf{z}_1, \mathbf{z}_2 \in \tilde{C}_t(k)} \min_{\mathbf{z}'_1 \in \mathcal{A}(\mathbf{z}_1), \mathbf{z}'_2 \in \mathcal{A}(\mathbf{z}_2)} \|\mathbf{z}'_1 - \mathbf{z}'_2\|_2 \leq \delta_t$  and (v)  $\mathbb{P}_t\{\cup_{k=1}^K \tilde{C}_t(k)\} \geq \sigma_t$ , where  $\sigma_s, \sigma_t \in (0, 1]$  and  $\delta_s, \delta_t \geq 0$ .

The  $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -augmentation is an extensive version of the  $(\sigma, \delta)$ -augmentation proposed by Huang et al. (2023). This definition emphasizes that a robust data augmentation should consistently produce distance-closed augmented views for semantically similar original images. Therein, condition (v) replaces the assumption  $\mathcal{A}(C_t(i)) \cap \mathcal{A}(C_t(j)) = \emptyset$  proposed by Huang et al. (2023). This implies that the augmentation methods used should be intelligent enough to recognize objects that align with the image labels in multi-objective images. A straightforward alternative to this requirement is to assume that different classes  $C_t(k)$  are pairwise disjoint, meaning that for all  $i \neq j$ ,  $C_t(i) \cap C_t(j) = \emptyset$ . This implies that  $\mathbb{P}_t\{\cup_{k=1}^K \tilde{C}_t(k)\} = \sum_{k=1}^K \mathbb{P}_t\{\tilde{C}_t(k)\} \geq \sigma_t \sum_{k=1}^K \mathbb{P}_t\{C_t(k)\} = \sigma_t$ .

In the context of Definition 3.6, we introduce the following assumption to delineate the data augmentation necessary for the end-to-end theory of ACT.

**Assumption 3.7** (Existence of augmentation sequence). Assume there exists a sequence of  $(\sigma_s^{(n)}, \sigma_t^{(n)}, \delta_s^{(n)}, \delta_t^{(n)})$ -data augmentations  $\mathcal{A}_n = \{A_i^{(n)}\}_{i \in [m]}$  and  $\tau > 0$  such that (i)  $\max\{\delta_s^{(n)}, \delta_t^{(n)}\} \leq n^{-\frac{\tau+d+1}{2(\alpha+d+1)}}$ , (ii)  $\min\{\sigma_s^{(n)}, \sigma_t^{(n)}\} \rightarrow 1$  as  $n \rightarrow \infty$ .

It is noteworthy that this assumption closely aligns with Assumption 3.5 in HaoChen et al. (2021) and Assumption 3.6 in HaoChen & Ma (2023), both of which stipulate that

the augmentations must be sufficiently robust to ensure that the internal connections within latent classes remain strong enough to prevent the separation of instance clusters. Recently, various methods for developing more effective data augmentations, as discussed by Jahanian et al. (2022) and Trabucco et al. (2024), have been proposed, making it increasingly feasible to satisfy the theoretical requirements for data augmentation. Next, we will introduce the assumption related to distribution shift.

Prior to characterizing the transferability from the source domain to the target domain, we must first quantify the similarity between these domains.

**Assumption 3.8** (Domain shift). Assume there exists  $\nu > 0$  and  $\varsigma > 0$  such that (i)  $\epsilon_1 \lesssim n_s^{-\frac{\nu+d+1}{2(\alpha+d+1)}}$  and (ii)  $\epsilon_2 \lesssim n_s^{-\frac{\varsigma}{2(\alpha+d+1)}}$ , where  $\epsilon_1$  and  $\epsilon_2$  measure the divergence between the source and the target domains defined as (7).

As shown in (7), smaller values of  $\epsilon_1$  and  $\epsilon_2$  indicate less discrepancy between the source and target domains. Similar assumptions using alternative divergence measures have been proposed in Ben-David et al. (2010); Germain et al. (2013); Cortes et al. (2019).

### 3.4. End-to-end Theoretical Guarantee

Let  $\mu_t(k) = \mathbb{E}_{\mathbf{x} \in C_t(k)} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} \{\hat{f}_{n_s}(\mathbf{x}')\}$ , which is the representation center of  $k$ -th class  $C_t(k)$ . We present the end-to-end theoretical guarantee of ACT as follows:

**Theorem 3.9.** Suppose Assumptions 3.2, 3.3, 3.5, 3.7 and 3.8 all hold. Set the width, depth and the Lipschitz constraint of the deep neural network as

$$W \geq \mathcal{O}\left(n_s^{\frac{2d+\alpha}{4(\alpha+d+1)}}\right), \quad L \geq \mathcal{O}(1), \quad \mathcal{K} = \mathcal{O}\left(n_s^{\frac{d+1}{2(\alpha+d+1)}}\right),$$

then the following inequality holds

$$\mathbb{E}_{\tilde{D}_s} \left\{ \max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \right\} \lesssim (1 - \sigma_s^{(n_s)}) + n_s^{-\frac{\min\{\alpha, 2\tau\}}{4(\alpha+d+1)}}. \quad (9)$$

Furthermore, regarding to the misclassification rate of  $Q_{\hat{f}_{n_s}}$ , we have

$$\mathbb{E}_{\tilde{D}_s} \left\{ \text{Err}(Q_{\hat{f}_{n_s}}) \right\} \leq (1 - \sigma_t^{(n_s)}) + \mathcal{O}\left(n_s^{-\frac{\min\{\alpha, \nu, \varsigma\}}{8(\alpha+d+1)}}\right),$$

with probability at least  $\sigma_s^{(n_s)} - \mathcal{O}\left(n_s^{-\frac{\min\{\alpha, \nu, \varsigma, \tau\}}{16(\alpha+d+1)}}\right) - \mathcal{O}\left(\frac{1}{\sqrt{\min_k n_t(k)}}\right)$  for  $n_s$  sufficiently large.

**Provable Advantages of ACT** Theorem 3.9 demonstrates how the abundance of unlabeled data in the source domain leveraged by ACT benefits downstream tasks in the target domain. Specifically, the quantity  $\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)|$  in eq (9) reflects the angle between different representation

centers when  $R_1 \approx R_2$ . A smaller value indicates that the centers approach orthogonality, enhancing discriminability among categories and thus improving classification accuracy. eq (9) essentially indicates that minimaxing the loss function of ACT can lead to the downstream data distribution being clustered in the representation space, provided that the upstream unlabeled sample size is sufficient. On the other hand, the theorem shows that only the failure probability depends on the downstream sample size  $n_t$  with fast convergence rate, while the misclassification rate converges with respect to the number of unlabeled samples in ACT. This finding indicates that the build classifier based on ACT can achieve excellent performance with a few labeled samples. In summary, this theorem not only demonstrates the provable advantages of ACT but also provides rigorous theoretical understandings for few-shot learning (Liu et al., 2021; Rizve et al., 2021; Yang et al., 2022; Lim et al., 2023).

**Over-parametrization** Theorem 3.9 does not impose an upper bound constraint on either the width  $W$  or depth  $L$  of the deep neural network, implying that the number of network parameters can grow arbitrarily large when only the weight norm is constrained. This aligns with the over-parameterization regime commonly concerned in deep learning.

#### 4. Comparison with Existing Methods

As the experiments conducted in existing self-supervised learning methods, we pretrain the representation on CIFAR-10, CIFAR-100 and Tiny ImageNet, and subsequently conduct fine-tuning on each dataset with annotations. Table 2 shows the classification accuracy of representations learned by ACT, compared with the results reported in Ermolov et al. (2021). We can see that ACT consistently outperforms previous mainstream self-supervised methods across various datasets and evaluation metrics.

Table 2. Classification accuracy (top 1) of a linear classifier and a 5-nearest neighbors classifier for different loss functions and datasets.

Method	CIFAR-10		CIFAR-100		Tiny ImageNet	
	Linear	$k$ -nn	Linear	$k$ -nn	Linear	$k$ -nn
SimCLR	91.80	88.42	66.83	56.56	48.84	32.86
BYOL	91.73	89.45	66.60	56.82	<b>51.00</b>	36.24
WMSE2	91.55	89.69	66.10	56.69	48.20	34.16
WMSE4	91.99	89.87	67.64	56.45	49.20	35.44
<b>ACT</b>	<b>92.11</b>	<b>90.01</b>	<b>68.24</b>	<b>58.35</b>	49.72	<b>36.40</b>

**Implementation details.** Except for tuning  $\lambda$  for different datasets, all other hyperparameters used in our experiments align with Ermolov et al. (2021). Before each iteration, we first standardize the representations and then calculate the

loss of ACT. We train for 1,000 epochs with a learning rate of  $3 \times 10^{-3}$  for CIFAR-10 and CIFAR-100, and  $2 \times 10^{-3}$  for Tiny ImageNet. A learning rate warm-up is applied for the first 500 iterations of the optimizer, in addition to a 0.2 learning rate drop at 50 and 25 epochs before the training end. We use a mini-batch size of 256, and the dimension of the hidden layer in the projection head is set to 1024. The weight decay is set to  $10^{-6}$ . We adopt an embedding size ( $d^*$ ) of 64 for CIFAR10, CIFAR100 and 128 for Tiny ImageNet during the pretraining process. The backbone network used in our implementation is ResNet-18.

**Image transformation details.** We randomly extract crops with sizes ranging from 0.08 to 1.0 of the original area and aspect ratios ranging from 3/4 to 4/3 of the original aspect ratio. Furthermore, we apply horizontal mirroring with a probability of 0.5. Additionally, color jittering is applied with a configuration of (0.4; 0.4; 0.4; 0.1) and a probability of 0.8, while grayscaling is applied with a probability of 0.2. For CIFAR-10 and CIFAR-100, random Gaussian blurring is adopted with a probability of 0.5 and a kernel size of 0.1. During testing, only one crop is used for evaluation.

**Evaluation protocol.** During evaluation, we freeze the network encoder and remove the projection head after pre-training, then train a supervised linear classifier on top of it, which is a fully-connected layer followed by softmax. we train the linear classifier for 500 epochs using the Adam optimizer with corresponding labeled training set without data augmentation. The learning rate is exponentially decayed from  $10^{-2}$  to  $10^{-6}$ . The weight decay is set as  $10^{-6}$ . we also include the accuracy of a  $k$ -nearest neighbors classifier with  $k = 5$ , which does not require fine tuning.

All experiments were conducted using a single Tesla V100 GPU unit. The PyTorch implementations can be found in <https://anonymous.4open.science/r/ACT-1F45>.

#### 5. Conclusion

In this paper, we propose a novel adversarial contrastive learning method for unsupervised transfer learning. Our experimental results achieved state-of-the-art classification accuracy under both fine-tuned linear probe and  $k$ -nn protocol on various real datasets, comparing with the self-supervised learning methods. Meanwhile, we present end to end theoretical guarantee for the downstream classification task under misspecified and over-parameterized setting. Our theoretical results not only indicate that the misclassification rate of downstream task solely depends on the strength of data augmentation on the large amount of unlabeled data, but also bridge the gap in the theoretical understanding of the effectiveness of few-shot learning for downstream tasks with small sample size.

## References

- Ash, J. T., Goel, S., Krishnamurthy, A., and Misra, D. Investigating the role of negatives in contrastive representation learning. In Camps-Valls, G., Ruiz, F. J. R., and Valera, I. (eds.), *International Conference on Artificial Intelligence and Statistics, AISTATS 2022, 28-30 March 2022, Virtual Event*, volume 151 of *Proceedings of Machine Learning Research*, pp. 7187–7209. PMLR, 2022. URL <https://proceedings.mlr.press/v151/ash22a.html>.
- Awasthi, P., Dikkala, N., and Kamath, P. Do more negative samples necessarily hurt in contrastive learning? In *International conference on machine learning*, pp. 1101–1116. PMLR, 2022.
- Bardes, A., Ponce, J., and LeCun, Y. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=xm6YD62D1Ub>.
- Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., and Vaughan, J. W. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.
- Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., and Joulin, A. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020.
- Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., and Joulin, A. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020a.
- Chen, X. and He, K. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.
- Chen, X., Fan, H., Girshick, R., and He, K. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.
- Chuang, C.-Y., Robinson, J., Lin, Y.-C., Torralba, A., and Jegelka, S. Debaised contrastive learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 8765–8775. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/63c3ddcc7b23daa1e42dc41f9a44a873-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/63c3ddcc7b23daa1e42dc41f9a44a873-Paper.pdf).
- Chuang, C.-Y., Hjelm, R. D., Wang, X., Vineet, V., Joshi, N., Torralba, A., Jegelka, S., and Song, Y. Robust contrastive learning against noisy views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16670–16681, 2022.
- Cortes, C., Mohri, M., and Medina, A. M. Adaptation based on generalized discrepancy. *Journal of Machine Learning Research*, 20(1):1–30, 2019.
- Ermolov, A., Siarohin, A., Sangineto, E., and Sebe, N. Whitening for self-supervised representation learning. In *International conference on machine learning*, pp. 3015–3024. PMLR, 2021.
- Garrido, Q., Chen, Y., Bardes, A., Najman, L., and LeCun, Y. On the duality between contrastive and non-contrastive self-supervised learning. *arXiv preprint arXiv:2206.02574*, 2022.
- Germain, P., Habrard, A., Laviolette, F., and Morvant, E. A pac-bayesian approach for domain adaptation with specialization to linear classifiers. In Dasgupta, S. and McAllester, D. (eds.), *Proceedings of the 30th International Conference on Machine Learning*, number 3 in *Proceedings of Machine Learning Research*, pp. 738–746. Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- Giné, E. and Nickl, R. *Mathematical Foundations of Infinite-Dimensional Statistical Models*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2016. ISBN 9781107043169. URL <https://books.google.com.hk/books?id=ywFGrgEACAAJ>.
- Golowich, N., Rakhlin, A., and Shamir, O. Size-independent sample complexity of neural networks. In Bubeck, S., Perchet, V., and Rigollet, P. (eds.), *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pp. 297–299. PMLR, 06–09 Jul 2018. URL <https://proceedings.mlr.press/v75/golowich18a.html>.
- Grill, J.-B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.

- HaoChen, J. Z. and Ma, T. A theoretical study of inductive biases in contrastive learning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/forum?id=AuEgNlEAmed>.
- HaoChen, J. Z., Wei, C., Gaidon, A., and Ma, T. Provable guarantees for self-supervised deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*, 34:5000–5011, 2021.
- HaoChen, J. Z., Wei, C., Kumar, A., and Ma, T. Beyond separability: Analyzing the linear transferability of contrastive representations to related subpopulations. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/ac112e8fffc4e5b9ece32070440a8ca43-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/ac112e8fffc4e5b9ece32070440a8ca43-Abstract-Conference.html).
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Huang, W., Yi, M., Zhao, X., and Jiang, Z. Towards the generalization of contrastive self-supervised learning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/forum?id=XDJwuEYHhme>.
- Jahaniyan, A., Puig, X., Tian, Y., and Isola, P. Generative models as a data source for multiview representation learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=qhAeZjs7dCL>.
- Jiao, Y., Wang, Y., and Yang, Y. Approximation bounds for norm constrained neural networks with applications to regression and gans. *Applied and Computational Harmonic Analysis*, 65:249–278, 2023.
- Lei, Y., Yang, T., Ying, Y., and Zhou, D.-X. Generalization analysis for contrastive representation learning. In *International Conference on Machine Learning*, pp. 19200–19227. PMLR, 2023.
- Lim, J. Y., Lim, K. M., Lee, C. P., and Tan, Y. X. Self-supervised contrastive learning for few-shot image classification. *Neural Networks*, 165:19–30, 2023.
- Liu, C., Fu, Y., Xu, C., Yang, S., Li, J., Wang, C., and Zhang, L. Learning a few-shot embedding model with contrastive learning. In *AAAI Conference on Artificial Intelligence*, 2021. URL <https://api.semanticscholar.org/CorpusID:235349153>.
- Maurer, A. A vector-contraction inequality for rademacher complexities. In *Algorithmic Learning Theory: 27th International Conference, ALT 2016, Bari, Italy, October 19-21, 2016, Proceedings 27*, pp. 3–17. Springer, 2016.
- Rizve, M. N., Khan, S., Khan, F. S., and Shah, M. Exploring complementary strengths of invariant and equivariant representations for few-shot learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10836–10846, 2021.
- Saunshi, N., Plevrakis, O., Arora, S., Khodak, M., and Khandeparkar, H. A theoretical analysis of contrastive unsupervised representation learning. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5628–5637. PMLR, 2019. URL <http://proceedings.mlr.press/v97/saunshi19a.html>.
- Schmidt-Hieber, J. Nonparametric regression using deep neural networks with ReLU activation function. *The Annals of Statistics*, 48(4):1875 – 1897, 2020. doi: 10.1214/19-AOS1875. URL <https://doi.org/10.1214/19-AOS1875>.
- Trabucco, B., Doherty, K., Gurinas, M., and Salakhutdinov, R. Effective data augmentation with diffusion models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=ZWzUA9zeAg>.
- Tsybakov, A. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer New York, 2008. ISBN 9780387790527. URL <https://books.google.com/books?id=mwB8rUBsbqoC>.
- Vershynin, R. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.
- Yang, Z., Wang, J., and Zhu, Y. Few-shot classification with contrastive learning. In *European conference on computer vision*, pp. 293–309. Springer, 2022.
- Ye, M., Zhang, X., Yuen, P. C., and Chang, S.-F. Unsupervised embedding learning via invariant and spreading

550 instance feature. In *Proceedings of the IEEE/CVF con-*  
551 *ference on computer vision and pattern recognition*, pp.  
552 6210–6219, 2019.

553 Zbontar, J., Jing, L., Misra, I., LeCun, Y., and Deny, S.  
554 Barlow twins: Self-supervised learning via redundancy  
555 reduction. In *International conference on machine learn-*  
556 *ing*, pp. 12310–12320. PMLR, 2021.

558 Zhang, Q., Wang, Y., and Wang, Y. Identifiable contrastive  
559 learning with automatic feature importance discovery. In  
560 *NeurIPS*, 2023.

561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604

## A. Deferred Proof

### A.1. $\mathcal{K}$ -Lipschitz property of $\mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}, B_1, B_2)$

*Proof.* To demonstrate that any function  $f_\theta \in \mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}, B_1, B_2)$  is a  $\mathcal{K}$ -Lipschitz function, we first define two special classes. The first class is given by

$$\mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}) := \{f_\theta(\mathbf{x}) = A_L \sigma(A_{L-1} \sigma(\cdots \sigma(A_0 \mathbf{x}))) : \kappa(\theta) \leq \mathcal{K}\}, \quad (10)$$

which is equivalent to  $\mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}, B_1, B_2)$  when ignoring the condition  $\|f_\theta\|_2 \in [B_1, B_2]$ . The second class is defined as

$$\mathcal{SNN}_{d_1, d_2}(W, L, \mathcal{K}) := \{\check{f}(\mathbf{x}) = \check{A}_L \sigma(\check{A}_{L-1} \sigma(\cdots \sigma(\check{A}_0 \check{\mathbf{x}}))) : \prod_{l=1}^L \|\check{A}_l\|_\infty \leq \mathcal{K}\}, \quad \check{\mathbf{x}} := \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix},$$

where  $\check{A}_l \in \mathbb{R}^{N_{l+1} \times N_l}$  with  $N_0 = d_1 + 1$ .

It is clear that  $\mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}, B_1, B_2) \subseteq \mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K})$ , and every element in  $\mathcal{SNN}_{d_1, d_2}(W, L, \mathcal{K})$  is a  $\mathcal{K}$ -Lipschitz function due to the 1-Lipschitz property of the ReLU activation function. Thus, it suffices to show that

$$\mathcal{SNN}_{d_1, d_2}(W, L, \mathcal{K}) \subseteq \mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}) \subseteq \mathcal{SNN}_{d_1, d_2}(W + 1, L, \mathcal{K})$$

to establish our claim.

To begin, any function  $f_\theta(\mathbf{x}) = A_L \sigma(A_{L-1} \sigma(\cdots \sigma(A_0 \mathbf{x} + \mathbf{b}_0))) + \mathbf{b}_{L-1} \in \mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K})$  can be restructured as  $\check{f}(\mathbf{x}) = \check{A}_L \sigma(\check{A}_{L-1} \sigma(\cdots \sigma(\check{A}_0 \check{\mathbf{x}})))$ , where

$$\check{\mathbf{x}} := \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}, \quad \check{A} = (A_L, \mathbf{0}), \quad \check{A}_l = \begin{pmatrix} A_l & \mathbf{b}_l \\ \mathbf{0} & 1 \end{pmatrix}, \quad l = 0, \dots, L-1.$$

Notably, we have  $\prod_{l=0}^L \|\check{A}_l\|_\infty = \|A_L\|_\infty \prod_{l=0}^{L-1} \max\{\|(A_l, \mathbf{b}_l)\|_\infty, 1\} = \kappa(\theta) \leq \mathcal{K}$ , which implies that  $f_\theta \in \mathcal{SNN}_{d_1, d_2}(W + 1, L, \mathcal{K})$ .

Conversely, since any  $\check{f} \in \mathcal{SNN}(W, L, \mathcal{K})$  can also be parameterized as  $A_L \sigma(A_{L-1} \sigma(\cdots \sigma(A_0 \mathbf{x} + \mathbf{b}_0))) + \mathbf{b}_{L-1}$  with  $\theta = (\check{A}_0, (\check{A}_1, \mathbf{0}), \dots, (\check{A}_{L-1}, \mathbf{0}), \check{A}_L)$ , we can use the absolute homogeneity of the ReLU function to rescale  $\check{A}_l$  such that  $\|\check{A}_L\|_\infty \leq \mathcal{K}$  and  $\|\check{A}_l\|_\infty = 1$  for  $l \neq L$ . Consequently, we have  $\kappa(\theta) = \prod_{l=0}^L \|\check{A}_l\|_\infty \leq \mathcal{K}$ , which yields  $\check{f} \in \mathcal{NN}(W, L, \mathcal{K})$ . This completes the proof.  $\square$

### A.2. Proof of population theorem

In this section, we aim to present the population theorem of ACT and its proof. we begin by exploring the sufficient condition for achieving small  $\text{Err}(Q_f)$  in A.2.1. Following that, we build the connection between the required condition and optimizing our adversarial self-supervised learning loss in Lemma A.4 of A.2.3, it reveals that small value of  $\mathcal{L}(f)$  may induce significant class divergence and highly augmented concentration. Lastly, by combining Lemma A.1 and Lemma A.4, we present the population theorem as Theorem A.5.

#### A.2.1. SUFFICIENT CONDITION OF SMALL MISCLASSIFICATION RATE

**Lemma A.1.** *Given a  $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -augmentation, if the encoder  $f$  such that  $B_1 \leq \|f\|_2 \leq B_2$  is  $\mathcal{K}$ -Lipschitz and*

$$\mu_t(i)^\top \mu_t(j) < B_2^2 \psi(\sigma_t, \delta_t, \varepsilon, f),$$

*holds for any pair of  $(i, j)$  with  $i \neq j$ , then the downstream error rate of  $Q_f$*

$$\text{Err}(Q_f) \leq (1 - \sigma_t) + R_t(\varepsilon, f),$$

*where  $\Delta_{\hat{\mu}_t} = 1 - \frac{\min_{k \in [K]} \|\hat{\mu}_t(k)\|_2^2}{B_2^2}$ . For any  $\varepsilon > 0$ ,  $R_t(\varepsilon, f) = \mathbb{P}_t(\mathbf{z} \in \cup_{k=1}^K C_t(k) : \sup_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2 >$*

*$\varepsilon)$  and  $\psi(\sigma_t, \delta_t, \varepsilon, f) = \Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, f) - \sqrt{2 - 2\Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, f)} - \frac{\Delta_{\hat{\mu}_t}}{2} - \frac{2 \max_{k \in [K]} \|\hat{\mu}_t(k) - \mu_t(k)\|_2}{B_2}$ , wherein*

$$\Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, f) = \left( \sigma_t - \frac{R_t(\varepsilon, f)}{\min_i p_t(i)} \right) \left( 1 + \left( \frac{B_1}{B_2} \right)^2 - \frac{\mathcal{K} \delta_t}{B_2} - \frac{2\varepsilon}{B_2} \right) - 1.$$

660 *Proof.* For any encoder  $f$ , let  $S_t(\varepsilon, f) := \{z \in \cup_{k=1}^K C_t(k) : \sup_{z_1, z_2 \in \mathcal{A}(z)} \|f(z_1) - f(z_2)\|_2 \leq \varepsilon\}$ , if any  $z \in$   
 661  $\{\tilde{C}_t(1) \cup \dots \cup \tilde{C}_t(K)\} \cap S_t(\varepsilon, f)$  can be correctly classified by  $Q_f$ , it turns out that  $\text{Err}(Q_f)$  can be bounded by  
 662  $(1 - \sigma_t) + R_t(\varepsilon, f)$ . In fact,

$$\begin{aligned} 663 \text{Err}(Q_f) &= \sum_{k=1}^K \mathbb{P}_t \{Q_f(z) \neq k, \forall z \in C_t(k)\} \leq \mathbb{P}_t \left[ \{\tilde{C}_t(1) \cup \dots \cup \tilde{C}_t(K) \cap S_t(\varepsilon, f)\}^c \right] \\ 664 &= \mathbb{P}_t \left[ \{\tilde{C}_t(1) \cup \dots \cup \tilde{C}_t(K)\}^c \cup \{S_t(\varepsilon, f)\}^c \right] \leq (1 - \sigma_t) + \mathbb{P}_t \left[ \{S_t(\varepsilon, f)\}^c \right] \\ 665 &= (1 - \sigma_t) + R_t(\varepsilon, f). \end{aligned}$$

666 The first row is derived from the definition of  $\text{Err}(Q_f)$ . Since any  $z \in \{\tilde{C}_t(1) \cup \dots \cup \tilde{C}_t(K)\} \cap S_t(\varepsilon, f)$  can be correctly  
 667 classified by  $Q_f$ , we obtain the second row. De Morgan's laws imply the third row. The fourth row follows from  
 668 Definition 3.6. Finally, noting that  $R_t(\varepsilon, f) = \mathbb{P}_t[\{S_t(\varepsilon, f)\}^c]$  yields the last line.

669 Hence it suffices to show for given  $i \in [K]$ ,  $z \in \tilde{C}_t(i) \cap S_t(\varepsilon, f)$  can be correctly classified by  $Q_f$  if for any  $j \neq i$ ,

$$670 \mu_t(i)^\top \mu_t(j) < B_2^2 \left( \Gamma_i(\sigma_t, \delta_t, \varepsilon, f) - \sqrt{2 - 2\Gamma_i(\sigma_t, \delta_t, \varepsilon, f)} - \frac{\Delta \hat{\mu}_t}{2} - \frac{\|\hat{\mu}_t(i) - \mu_t(i)\|_2}{B_2} - \frac{\|\hat{\mu}_t(j) - \mu_t(j)\|_2}{B_2} \right),$$

671 where  $\Gamma_i(\sigma_t, \delta_t, \varepsilon, f) = \left( \sigma_t - \frac{R_t(\varepsilon, f)}{p_t(i)} \right) \left( 1 + \left( \frac{B_1}{B_2} \right)^2 - \frac{\kappa \delta_t}{B_2} - \frac{2\varepsilon}{B_2} \right) - 1$ .

672 To this end, without losing generality, consider the case  $i = 1$ . To turn out  $z_0 \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)$  can be correctly classified  
 673 by  $Q_f$ , by the definition of  $\tilde{C}_t(1)$  and  $S_t(\varepsilon, f)$ , It just need to show  $\forall k \neq 1, \|f(z_0) - \hat{\mu}_t(1)\|_2 < \|f(z_0) - \hat{\mu}_t(k)\|_2$ , which  
 674 is equivalent to

$$675 f(z_0)^\top \hat{\mu}_t(1) - f(z_0)^\top \hat{\mu}_t(k) - \left( \frac{1}{2} \|\hat{\mu}_t(1)\|_2^2 - \frac{1}{2} \|\hat{\mu}_t(k)\|_2^2 \right) > 0.$$

676 We first deal with the term  $f(z_0)^\top \hat{\mu}_t(1)$ ,

$$\begin{aligned} 677 f(z_0)^\top \hat{\mu}_t(1) &= f(z_0)^\top \mu_t(1) + f(z_0)^\top (\hat{\mu}_t(1) - \mu_t(1)) \\ 678 &\geq f(z_0)^\top \mathbb{E}_{z \in C_t(1)} \mathbb{E}_{z' \in \mathcal{A}(z)} \{f(z')\} - \|f(z_0)\|_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 \\ 679 &\geq \frac{1}{p_t(1)} f(z_0)^\top \mathbb{E}_z \mathbb{E}_{z' \in \mathcal{A}(z)} [f(z') \mathbb{1}\{z \in C_t(1)\}] - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 \\ 680 &= \frac{1}{p_t(1)} f(z_0)^\top \mathbb{E}_z \mathbb{E}_{z' \in \mathcal{A}(z)} [f(z') \mathbb{1}\{z \in C_t(1) \cap \tilde{C}_t(1) \cap S_t(\varepsilon, f)\}] \\ 681 &\quad + \frac{1}{p_t(1)} f(z_0)^\top \mathbb{E}_z \mathbb{E}_{z' \in \mathcal{A}(z)} [f(z') \mathbb{1}\{z \in C_t(1) \cap \{\tilde{C}_t(1) \cap S_t(\varepsilon, f)\}^c\}] - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 \\ 682 &= \frac{\mathbb{P}_t\{\tilde{C}_t(1) \cap S_t(\varepsilon, f)\}}{p_t(1)} f(z_0)^\top \mathbb{E}_{z \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{z' \in \mathcal{A}(z)} \{f(z')\} \\ 683 &\quad + \frac{1}{p_t(1)} \mathbb{E}_z [\mathbb{E}_{z' \in \mathcal{A}(z)} \{f(z_0)^\top f(z')\} \mathbb{1}\{z \in C_t(1) \setminus \{\tilde{C}_t(1) \cap S_t(\varepsilon, f)\}\}] - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 \\ 684 &\geq \frac{\mathbb{P}_t\{\tilde{C}_t(1) \cap S_t(\varepsilon, f)\}}{p_t(1)} f(z_0)^\top \mathbb{E}_{z \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{z' \in \mathcal{A}(z)} [f(z')] - \frac{B_2^2}{p_t(1)} \mathbb{P}_t[C_t(1) \setminus \{\tilde{C}_t(1) \cap S_t(\varepsilon, f)\}] \\ 685 &\quad - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2. \end{aligned} \tag{11}$$

686 The second row follows from the Cauchy–Schwarz inequality. The third and last rows are derived from the condition  
 687  $\|f\|_2 \leq B_2$ . Note that

$$688 \mathbb{P}_t[C_t(1) \setminus \{\tilde{C}_t(1) \cap S_t(\varepsilon, f)\}] = \mathbb{P}_t[\{C_t(1) \setminus \tilde{C}_t(1)\} \cup \{\tilde{C}_t(1) \cap \{S_t(\varepsilon, f)\}^c\}] \leq (1 - \sigma_t) p_t(1) + R_t(\varepsilon, f), \tag{12}$$

689 and

$$690 \mathbb{P}_t(\tilde{C}_t(1) \cap S_t(\varepsilon, f)) = \mathbb{P}_t(C_t(1)) - \mathbb{P}_t(C_t(1) \setminus (\tilde{C}_t(1) \cap S_t(\varepsilon, f))) \geq p_t(1) - \{(1 - \sigma_t) p_t(1) + R_t(\varepsilon, f)\}$$

$$= \sigma_t p_t(1) - R_t(\varepsilon, f). \quad (13)$$

Plugging (12) and (13) into (11) yields

$$f(\mathbf{z}_0)^\top \hat{\mu}_t(1) \geq \left( \sigma_t - \frac{R_t(\varepsilon, f)}{p_t(1)} \right) f(\mathbf{z}_0)^\top \mathbb{E}_{\mathbf{z} \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} \{f(\mathbf{z}')\} - B_2^2 \left( 1 - \sigma_t + \frac{R_t(\varepsilon, f)}{p_t(1)} \right) - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2. \quad (14)$$

Notice that  $\mathbf{z}_0 \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)$ . Thus, for any  $\mathbf{z} \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)$ , by the definition of  $\tilde{C}_t(1)$ , we have  $\min_{\mathbf{z}'_0 \in \mathcal{A}(\mathbf{z}_0), \mathbf{z}' \in \mathcal{A}(\mathbf{z})} \|\mathbf{z}'_0 - \mathbf{z}'\|_2 \leq \delta_t$ . Further, denote  $(\mathbf{z}'_0, \mathbf{z}^*) = \arg \min_{\mathbf{z}'_0 \in \mathcal{A}(\mathbf{z}_0), \mathbf{z}' \in \mathcal{A}(\mathbf{z})} \|\mathbf{z}'_0 - \mathbf{z}'\|_2$ . Then, we have  $\|\mathbf{z}'_0 - \mathbf{z}^*\|_2 \leq \delta_t$ . Combining this with the  $\mathcal{K}$ -Lipschitz property of  $f$ , we obtain  $\|f(\mathbf{z}'_0) - f(\mathbf{z}^*)\|_2 \leq \mathcal{K} \|\mathbf{z}'_0 - \mathbf{z}^*\|_2 \leq \mathcal{K} \delta_t$ . Moreover, since  $\mathbf{z} \in S_t(\varepsilon, f)$ , it follows that for all  $\mathbf{z}' \in \mathcal{A}(\mathbf{z})$ ,  $\|f(\mathbf{z}') - f(\mathbf{z}^*)\|_2 \leq \varepsilon$ . Similarly, as  $\mathbf{z}_0 \in S_t(\varepsilon, f)$  and both  $\mathbf{z}_0$  and  $\mathbf{z}'_0$  belong to  $\mathcal{A}(\mathbf{z}_0)$ , we know  $\|f(\mathbf{z}_0) - f(\mathbf{z}'_0)\|_2 \leq \varepsilon$ .

Therefore,

$$\begin{aligned} f(\mathbf{z}_0)^\top \mathbb{E}_{\mathbf{z} \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} \{f(\mathbf{z}')\} &= \mathbb{E}_{\mathbf{z} \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} \{f(\mathbf{z}_0)^\top f(\mathbf{z}')\} \\ &= \mathbb{E}_{\mathbf{z} \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} [f(\mathbf{z}_0)^\top \{f(\mathbf{z}') - f(\mathbf{z}_0) + f(\mathbf{z}_0)\}] \\ &\geq B_1^2 + \mathbb{E}_{\mathbf{z} \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} [f(\mathbf{z}_0)^\top \{f(\mathbf{z}') - f(\mathbf{z}_0)\}] \\ &= B_1^2 + \mathbb{E}_{\mathbf{z} \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} [f(\mathbf{z}_0)^\top \underbrace{\{f(\mathbf{z}') - f(\mathbf{z}^*)\}}_{\|\cdot\|_2 \leq \varepsilon} + \underbrace{f(\mathbf{z}^*) - f(\mathbf{z}'_0)}_{\|\cdot\|_2 \leq \mathcal{K} \delta_t} + \underbrace{f(\mathbf{z}'_0) - f(\mathbf{z}_0)}_{\|\cdot\|_2 \leq \varepsilon}] \\ &\geq B_1^2 - (B_2 \varepsilon + B_2 \mathcal{K} \delta_t + B_2 \varepsilon) \\ &= B_1^2 - B_2 (\mathcal{K} \delta_t + 2\varepsilon), \end{aligned} \quad (15)$$

where the fourth row is derived from  $\|f\|_2 \geq B_1$ .

Plugging (15) into the inequality (14) yields

$$\begin{aligned} f(\mathbf{z}_0)^\top \hat{\mu}_t(1) &\geq \left( \sigma_t - \frac{R_t(\varepsilon, f)}{p_t(1)} \right) f(\mathbf{z}_0)^\top \mathbb{E}_{\mathbf{z} \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} \{f(\mathbf{z}')\} - B_2^2 \left( 1 - \sigma_t + \frac{R_t(\varepsilon, f)}{p_t(1)} \right) - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 \\ &\geq \left( \sigma_t - \frac{R_t(\varepsilon, f)}{p_t(1)} \right) (B_1^2 - B_2 (\mathcal{K} \delta_t + 2\varepsilon)) - B_2^2 \left\{ 1 - \sigma_t + \frac{R_t(\varepsilon, f)}{p_t(1)} \right\} - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 \\ &= B_2^2 \left\{ \left( 1 + \left( \frac{B_1}{B_2} \right)^2 \right) \left( \sigma_t - \frac{R_t(\varepsilon, f)}{p_t(1)} \right) - \left( \sigma_t - \frac{R_t(\varepsilon, f)}{p_t(1)} \right) \left( \frac{\mathcal{K} \delta_t}{B_2} + \frac{2\varepsilon}{B_2} \right) - 1 \right\} - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 \\ &= B_2^2 \left\{ \left( \sigma_t - \frac{R_t(\varepsilon, f)}{p_t(1)} \right) \left( 1 + \left( \frac{B_1}{B_2} \right)^2 - \frac{\mathcal{K} \delta_t}{B_2} - \frac{2\varepsilon}{B_2} \right) - 1 \right\} - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 \\ &= B_2^2 \Gamma_1(\sigma_t, \delta_t, \varepsilon, f) - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2. \end{aligned}$$

Similar process can also turn out

$$f(\mathbf{z}_0)^\top \mu_t(1) \geq B_2^2 \Gamma_1(\sigma_t, \delta_t, \varepsilon, f). \quad (16)$$

Combining the fact that  $\|\mu_t(k)\|_2 = \|\mathbb{E}_{\mathbf{z} \in \tilde{C}_t(k)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} [f(\mathbf{z}')] \|_2 \leq \mathbb{E}_{\mathbf{x} \in \tilde{C}_t(k)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}')\|_2 \leq B_2$  yields

$$\begin{aligned} f(\mathbf{z}_0)^\top \hat{\mu}_t(k) &\leq f(\mathbf{z}_0)^\top \mu_t(k) + f(\mathbf{z}_0)^\top (\hat{\mu}_t(k) - \mu_t(k)) \\ &\leq f(\mathbf{z}_0)^\top \mu_t(k) + \|f(\mathbf{z}_0)\|_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 \\ &\leq f(\mathbf{z}_0)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 \\ &= \{f(\mathbf{z}_0) - \mu_t(1)\}^\top \mu_t(k) + \mu_t(1)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 \\ &\leq \|f(\mathbf{z}_0) - \mu_t(1)\|_2 \cdot \|\mu_t(k)\|_2 + \mu_t(1)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 \\ &\leq B_2 \sqrt{\|f(\mathbf{z}_0)\|_2^2 - 2f(\mathbf{z}_0)^\top \mu_t(1) + \|\mu_t(1)\|_2^2} + \mu_t(1)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 \end{aligned}$$

$$\begin{aligned}
 &\leq B_2 \sqrt{2B_2^2 - 2f(\mathbf{z}_0)^\top \mu_t(1) + \mu_t(1)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2} \\
 &\leq B_2 \sqrt{2B_2^2 - 2B_2^2 \Gamma_1(\sigma_t, \delta_t, \varepsilon, f) + \mu_t(1)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2} \\
 &= \sqrt{2B_2^2} \sqrt{1 - \Gamma_1(\sigma_t, \delta_t, \varepsilon, f) + \mu_t(1)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2},
 \end{aligned}$$

where the inequality in eighth row stems from (16). Hence, by  $\Delta_{\hat{\mu}_t} = 1 - \min_{k \in [K]} \|\hat{\mu}_t(k)\|_2^2 / B_2^2$ , we can conclude

$$\begin{aligned}
 &f(\mathbf{z}_0)^\top \hat{\mu}_t(1) - f(\mathbf{z}_0)^\top \hat{\mu}_t(k) - \left( \frac{1}{2} \|\hat{\mu}_t(1)\|_2^2 - \frac{1}{2} \|\hat{\mu}_t(k)\|_2^2 \right) = f(\mathbf{z}_0)^\top \hat{\mu}_t(1) - f(\mathbf{z}_0)^\top \hat{\mu}_t(k) - \frac{1}{2} \|\hat{\mu}_t(1)\|_2^2 + \frac{1}{2} \|\hat{\mu}_t(k)\|_2^2 \\
 &\geq f(\mathbf{z}_0)^\top \hat{\mu}_t(1) - f(\mathbf{z}_0)^\top \hat{\mu}_t(k) - \frac{1}{2} B_2^2 + \frac{1}{2} \min_{k \in [K]} \|\hat{\mu}_t(k)\|_2^2 \\
 &= f(\mathbf{z}_0)^\top \hat{\mu}_t(1) - f(\mathbf{z}_0)^\top \hat{\mu}_t(k) - \frac{1}{2} B_2^2 \Delta_{\hat{\mu}_t} \\
 &\geq B_2^2 \Gamma_1(\sigma_t, \delta_t, \varepsilon, f) - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 - \sqrt{2} B_2 \sqrt{1 - \Gamma_1(\sigma_t, \delta_t, \varepsilon, f)} - \mu_t(1)^\top \mu_t(k) - B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 - \\
 &\quad - \frac{1}{2} B_2^2 \Delta_{\hat{\mu}_t} > 0,
 \end{aligned}$$

which is what we desire. Here the last inequality is derived from the given condition.  $\square$

#### A.2.2. PRELIMINARIES FOR LEMMA A.4

To obtain Theorem A.5, we need to bridge the gap between the condition in Lemma A.1 and the insights provided by ACT in Lemma A.4. To this end, we first introduce Lemma A.2 and Lemma A.3.

Following the notations in the target domain, we denote the center of the  $k$ -th latent class in the representation space as  $\mu_s(k) := \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}')\} = \frac{1}{p_s(k)} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}') \mathbb{1}\{\mathbf{x} \in C_s(k)\}]$ . Then Lemma A.2 can be presented as follows:

**Lemma A.2.** *If the encoder  $f$  is  $\mathcal{K}$ -Lipschitz continuous, then for any  $k \in [K]$ ,*

$$\|\mu_s(k) - \mu_t(k)\|_2 \leq \sqrt{d^*} MK \epsilon_1.$$

*Proof.* For all  $k \in [K]$ ,

$$\begin{aligned}
 \|\mu_s(k) - \mu_t(k)\|_2^2 &= \sum_{l=1}^{d^*} [\{\mu_s(k)\}_l - \{\mu_t(k)\}_l]^2 = \sum_{l=1}^{d^*} [\mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} \{f_l(\mathbf{x}')\} - \mathbb{E}_{\mathbf{z} \in C_t(k)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} \{f_l(\mathbf{z}')\}]^2 \\
 &= \sum_{l=1}^{d^*} \left[ \frac{1}{m} \sum_{\gamma=1}^m \left( \mathbb{E}_{\mathbf{x} \in C_s(k)} \{f_l(A_i(\mathbf{x}))\} - \mathbb{E}_{\mathbf{z} \in C_t(k)} \{f_l(A_i(\mathbf{z}))\} \right) \right]^2 \leq d^* M^2 \mathcal{K}^2 \epsilon_1^2.
 \end{aligned}$$

The final inequality is obtained from  $\epsilon_1 = \max_{k \in [K]} \mathcal{W}(\mathbb{P}_s(k), \mathbb{P}_t(k))$  and the definition of Wasserstein distance, along with the fact that  $f(A_i(\cdot))$  is  $M\mathcal{K}$ -Lipschitz continuous. In fact, since  $f \in \text{Lip}(\mathcal{K})$ , it follows that for every  $l \in [d^*]$ ,  $f_l \in \text{Lip}(\mathcal{K})$ . Combining this with the property that  $A_i(\cdot) \in \text{Lip}(M)$  stated in Assumption 3.5, we conclude that  $f(A_i(\cdot))$  is  $M\mathcal{K}$ -Lipschitz continuous. So that

$$\|\mu_s(k) - \mu_t(k)\|_2 \leq \sqrt{d^*} MK \epsilon_1.$$

Next we present Lemma A.3.

**Lemma A.3.** *Given a  $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -augmentation, if the encoder  $f$  with  $\|f\|_2 \leq B_2$  is  $\mathcal{K}$ -Lipschitz continuous, then*

$$\mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 \leq 4B_2^2 \left\{ \left( 1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2} + \frac{R_s(\varepsilon, f)}{p_s(k)} \right)^2 + \left( 1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)} \right) \right\},$$

where  $R_s(\varepsilon, f) = \mathbb{P}_s \{ \mathbf{x} \in \cup_{k=1}^K C_s(k) : \sup_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 > \varepsilon \}$ .

*Proof.* Let  $S_s(\varepsilon, f) := \{\mathbf{x} \in \cup_{k=1}^K C_s(k) : \sup_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 \leq \varepsilon\}$ , for each  $k \in [K]$ ,

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 = \frac{1}{p_s(k)} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\mathbb{1}\{\mathbf{x} \in C_s(k)\} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2] \\
 & = \frac{1}{p_s(k)} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\mathbb{1}\{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)\} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2] \\
 & \quad + \frac{1}{p_s(k)} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\mathbb{1}\{\mathbf{x} \in C_s(k) \setminus (\tilde{C}_s(k) \cap S_s(\varepsilon, f))\} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2] \\
 & \leq \frac{1}{p_s(k)} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\mathbb{1}\{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)\} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2] + \frac{4B_2^2 \mathbb{P}_s[C_s(k) \setminus \{\tilde{C}_s(k) \cap S_s(\varepsilon, f)\}]}{p_s(k)} \\
 & \leq \frac{1}{p_s(k)} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\mathbb{1}\{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)\} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2] + 4B_2^2 \left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right) \\
 & \leq \frac{\mathbb{P}_s(\tilde{C}_s(k) \cap S_s(\varepsilon, f))}{p_s(k)} \mathbb{E}_{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 + 4B_2^2 \left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right) \\
 & \leq \mathbb{E}_{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 + 4B_2^2 \left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right), \tag{17}
 \end{aligned}$$

where the second inequality is due to

$$\mathbb{P}_s[C_s(k) \setminus \{\tilde{C}_s(k) \cap S_s(\varepsilon, f)\}] = \mathbb{P}_s[\{C_s(k) \setminus \tilde{C}_s(k)\} \cup \{C_s(k) \setminus S_s(\varepsilon, f)\}] \leq (1 - \sigma_s)p_s(k) + R_s(\varepsilon, f).$$

Furthermore,

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 = \mathbb{E}_{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - \mathbb{E}_{\mathbf{x}' \in C_s(k)} \mathbb{E}_{\mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')} f(\mathbf{x}_2)\|_2^2 \\
 & = \mathbb{E}_{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \left\| f(\mathbf{x}_1) - \frac{\mathbb{P}_s\{\tilde{C}_s(k) \cap S_s(\varepsilon, f)\}}{p_s(k)} \mathbb{E}_{\mathbf{x}' \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')} f(\mathbf{x}_2) \right. \\
 & \quad \left. - \frac{\mathbb{P}_s[C_s(k) \setminus \{\tilde{C}_s(k) \cap S_s(\varepsilon, f)\}]}{p_s(k)} \mathbb{E}_{\mathbf{x}' \in C_s(k) \setminus \{\tilde{C}_s(k) \cap S_s(\varepsilon, f)\}} \mathbb{E}_{\mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')} f(\mathbf{x}_2) \right\|_2^2 \\
 & = \mathbb{E}_{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \left\| \frac{\mathbb{P}_s\{\tilde{C}_s(k) \cap S_s(\varepsilon, f)\}}{p_s(k)} \left( f(\mathbf{x}_1) - \mathbb{E}_{\mathbf{x}' \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')} f(\mathbf{x}_2) \right) \right. \\
 & \quad \left. - \frac{\mathbb{P}_s[C_s(k) \setminus \{\tilde{C}_s(k) \cap S_s(\varepsilon, f)\}]}{p_s(k)} \left( f(\mathbf{x}_1) - \mathbb{E}_{\mathbf{x}' \in C_s(k) \setminus \{\tilde{C}_s(k) \cap S_s(\varepsilon, f)\}} \mathbb{E}_{\mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')} f(\mathbf{x}_2) \right) \right\|_2^2 \\
 & \leq \mathbb{E}_{\mathbf{x} \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \left[ \left\| f(\mathbf{x}_1) - \mathbb{E}_{\mathbf{x}' \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)} \mathbb{E}_{\mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')} f(\mathbf{x}_2) \right\|_2 + 2B_2 \left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right) \right]^2 \tag{18}
 \end{aligned}$$

For any  $\mathbf{x}, \mathbf{x}' \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)$ , by the definition of  $\tilde{C}_s(k)$ , we can yield that

$$\min_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x}), \mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')} \|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq \delta_s,$$

Thus, let  $(\mathbf{x}_1^*, \mathbf{x}_2^*) = \arg \min_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x}), \mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')} \|\mathbf{x}_1 - \mathbf{x}_2\|_2$ , we have  $\|\mathbf{x}_1^* - \mathbf{x}_2^*\|_2 \leq \delta_s$ . Furthermore, by the  $\mathcal{K}$ -Lipschitz continuity of  $f$ , we yield  $\|f(\mathbf{x}_1^*) - f(\mathbf{x}_2^*)\|_2 \leq \mathcal{K}\|\mathbf{x}_1^* - \mathbf{x}_2^*\|_2 \leq \mathcal{K}\delta_s$ . In addition, since  $\mathbf{x} \in S_s(\varepsilon, f)$ , we know for any  $\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})$ ,  $\|f(\mathbf{x}_1) - f(\mathbf{x}_1^*)\|_2 \leq \varepsilon$ . Similarly,  $\mathbf{x}' \in S_s(\varepsilon, f)$  implies  $\|f(\mathbf{x}_2) - f(\mathbf{x}_2^*)\|_2 \leq \varepsilon$  for any  $\mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')$ . Therefore, for any  $\mathbf{x}, \mathbf{x}' \in \tilde{C}_s(k) \cap S_s(\varepsilon, f)$  and  $\mathbf{x}_1 \in \mathcal{A}(\mathbf{x}), \mathbf{x}_2 \in \mathcal{A}(\mathbf{x}')$ ,

$$\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 \leq \|f(\mathbf{x}_1) - f(\mathbf{x}_1^*)\|_2 + \|f(\mathbf{x}_1^*) - f(\mathbf{x}_2^*)\|_2 + \|f(\mathbf{x}_2^*) - f(\mathbf{x}_2)\|_2 \leq 2\varepsilon + \mathcal{K}\delta_s. \tag{19}$$

Combining inequalities (17), (18) and (19) concludes

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 & \leq \left[ 2\varepsilon + \mathcal{K}\delta_s + 2B_2 \left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right) \right]^2 + 4B_2^2 \left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right) \\
 & = 4B_2^2 \left[ \left(1 - \sigma_s + \frac{\mathcal{K}\delta_s}{2B_2} + \frac{\varepsilon}{B_2} + \frac{R_s(\varepsilon, f)}{p_s(k)}\right)^2 + \left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right) \right]
 \end{aligned}$$

□

Subsequently, we state Lemma A.4 to establish the connection between ACT and the requirements outlined in Lemma A.1.

### A.2.3. THE EFFECT OF MINIMAXING OUR LOSS

**Lemma A.4.** *Given a  $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -augmentation, if  $d^* > K$  and the encoder  $f$  with  $B_1 \leq \|f\|_2 \leq B_2$  is  $\mathcal{K}$ -Lipschitz continuous, then for any  $\varepsilon > 0$ ,*

$$\begin{aligned} R_s^2(\varepsilon, f) &\leq \frac{m^4}{\varepsilon^2} \mathcal{L}_{\text{align}}(f), \\ R_t^2(\varepsilon, f) &\leq \frac{m^4}{\varepsilon^2} \mathcal{L}_{\text{align}}(f) + \frac{8m^4}{\varepsilon^2} B_2 d^* M \mathcal{K} \varepsilon_1 + \frac{4m^4}{\varepsilon^2} B_2^2 d^* K \varepsilon_2, \end{aligned}$$

and

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \leq \sqrt{\frac{2}{\min_{i \neq j} p_s(i) p_s(j)}} \left\{ \mathcal{R}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f) \right\} + 2\sqrt{d^*} B_2 M \mathcal{K} \varepsilon_1.$$

where  $\varphi(\sigma_s, \delta_s, \varepsilon, f) := 4B_2^2 \left[ \left(1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2}\right)^2 + (1 - \sigma_s) + K R_s(\varepsilon, f) \left(3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2}\right) + R_s^2(\varepsilon, f) \left(\sum_{k=1}^K \frac{1}{p_s(k)}\right) \right] + B_2(\varepsilon^2 + 4B_2^2 R_s(\varepsilon, f))^{\frac{1}{2}}$ .

*Proof.* Since the measure on  $\mathcal{A}$  is uniform, we have

$$\mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2 = \frac{1}{m^2} \sum_{\gamma=1}^m \sum_{\beta=1}^m \|f(A_\gamma(\mathbf{z})) - f(A_\beta(\mathbf{z}))\|_2,$$

hence,

$$\begin{aligned} \sup_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2 &= \sup_{\gamma, \beta \in [m]} \|f(A_\gamma(\mathbf{z})) - f(A_\beta(\mathbf{z}))\|_2 \leq \sum_{\gamma=1}^m \sum_{\beta=1}^m \|f(A_\gamma(\mathbf{z})) - f(A_\beta(\mathbf{z}))\|_2 \\ &= m^2 \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2. \end{aligned}$$

Denote  $S := \{\mathbf{z} : \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2 > \frac{\varepsilon}{m^2}\}$ , by the definition of  $R_t(\varepsilon, f)$  along with Markov inequality, we have

$$\begin{aligned} R_t^2(\varepsilon, f) &\leq \mathbb{P}_t^2(S) \leq \left( \frac{\mathbb{E}_{\mathbf{z}} \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2}{\frac{\varepsilon}{m^2}} \right)^2 \leq \frac{\mathbb{E}_{\mathbf{z}} \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2^2}{\frac{\varepsilon^2}{m^4}} \\ &= \frac{m^4}{\varepsilon^2} \mathbb{E}_{\mathbf{z}} \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2^2. \end{aligned} \quad (20)$$

Apart from that, similar process yields the first inequity to be justified in Lemma A.4:

$$R_s^2(\varepsilon, f) \leq \frac{m^4}{\varepsilon^2} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2 = \frac{m^4}{\varepsilon^2} \mathcal{L}_{\text{align}}(f).$$

Furthermore, we can turn out

$$\begin{aligned} &\mathbb{E}_{\mathbf{z}} \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2^2 \\ &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2 + \mathbb{E}_{\mathbf{z}} \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2^2 - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2 \\ &= \frac{1}{m^2} \sum_{\gamma=1}^m \sum_{\beta=1}^m \left\{ \mathbb{E}_{\mathbf{z}} \|f(A_\gamma(\mathbf{z})) - f(A_\beta(\mathbf{z}))\|_2^2 - \mathbb{E}_{\mathbf{x}} \|f(A_\gamma(\mathbf{x})) - f(A_\beta(\mathbf{x}))\|_2^2 \right\} + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2 \\ &= \frac{1}{m^2} \sum_{\gamma=1}^m \sum_{\beta=1}^m \sum_{l=1}^{d^*} \left[ \mathbb{E}_{\mathbf{z}} \{f_l(A_\gamma(\mathbf{z})) - f_l(A_\beta(\mathbf{z}))\}^2 - \mathbb{E}_{\mathbf{x}} \{f_l(A_\gamma(\mathbf{x})) - f_l(A_\beta(\mathbf{x}))\}^2 \right] + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2, \end{aligned}$$

we subsequently focus on dealing with the first term. Since for all  $\gamma \in [m], \beta \in [m]$  and  $l \in [d^*]$ ,

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{z}} [f_l(A_i(\mathbf{z})) - f_l(A_j(\mathbf{z}))]^2 - \mathbb{E}_{\mathbf{x}} [f_l(A_i(\mathbf{x})) - f_l(A_j(\mathbf{x}))]^2 \\
 &= \sum_{k=1}^K \left[ p_t(k) \mathbb{E}_{\mathbf{z} \in C_t(k)} \{f_l(A_i(\mathbf{z})) - f_l(A_j(\mathbf{z}))\}^2 - p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \{f_l(A_i(\mathbf{x})) - f_l(A_j(\mathbf{x}))\}^2 \right] \\
 &= \sum_{k=1}^K \left[ p_t(k) \left\{ \mathbb{E}_{\mathbf{z} \in C_t(k)} \{f_l(A_i(\mathbf{z})) - f_l(A_j(\mathbf{z}))\}^2 - \underbrace{\mathbb{E}_{\mathbf{x} \in C_s(k)} \{f_l(A_i(\mathbf{x})) - f_l(A_j(\mathbf{x}))\}^2}_{g(\mathbf{x})} \right\} \right. \\
 &\quad \left. + \{p_t(k) - p_s(k)\} \mathbb{E}_{\mathbf{x} \in C_s(k)} \{f_l(A_i(\mathbf{x})) - f_l(A_j(\mathbf{x}))\}^2 \right] \\
 &\leq 8B_2MK\epsilon_1 + 4B_2^2K\epsilon_2.
 \end{aligned}$$

To obtain the last inequality, it suffices to show  $g(\mathbf{x}) \in \text{Lip}(8B_2MK)$ . In fact, we know  $\forall l \in [d^*], f_l \in \text{Lip}(K)$  as  $f \in \text{Lip}(K)$ , along with the fact that  $A_i(\cdot)$  and  $A_j(\cdot)$  are both  $M$ -Lipschitz continuous according to Assumption 3.5, we can conclude  $f_l(A_i(\cdot)) - f_l(A_j(\cdot)) \in \text{Lip}(2MK)$ . Additionally, note that  $|f_l(A_i(\cdot)) - f_l(A_j(\cdot))| \leq 2B_2$  as  $\|f\|_2 \leq B_2$ , we can turn out outermost quadratic function remains locally  $4B_2$ -Lipschitz continuity in  $[-2B_2, 2B_2]$ , which implies that  $g \in \text{Lip}(8B_2MK)$ . Furthermore, by the definition of Wasserstein distance, we yield

$$\sum_{k=1}^K \left[ p_t(k) \left( \mathbb{E}_{\mathbf{z} \in C_t(k)} \{f_l(A_i(\mathbf{z})) - f_l(A_j(\mathbf{z}))\}^2 - \mathbb{E}_{\mathbf{x} \in C_s(k)} \{f_l(A_i(\mathbf{x})) - f_l(A_j(\mathbf{x}))\}^2 \right) \right] \leq 8B_2MK\epsilon_1 \sum_{k=1}^K p_t(k) = 8B_2MK\epsilon_1,$$

As for the second term in the last inequality, note that  $f_l(A_i(\mathbf{x})) - f_l(A_j(\mathbf{x})) \leq 2B_2$  to yield

$$\sum_{k=1}^K \left[ \{p_t(k) - p_s(k)\} \mathbb{E}_{\mathbf{x} \in C_s(k)} \{f_l(A_i(\mathbf{x})) - f_l(A_j(\mathbf{x}))\}^2 \right] \leq 4B_2^2K\epsilon_2.$$

Therefore,

$$\mathbb{E}_{\mathbf{z}} \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}_1) - f(\mathbf{z}_2)\|_2^2 \leq \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2 + 8B_2d^*MK\epsilon_1 + 4B_2^2d^*K\epsilon_2. \quad (21)$$

Combining (20) and (21) turns out the second inequality of Lemma A.4.

$$R_t^2(\varepsilon, f) \leq \frac{m^4}{\varepsilon^2} \mathcal{L}_{\text{align}}(f) + \frac{8m^4}{\varepsilon^2} B_2d^*MK\epsilon_1 + \frac{4m^4}{\varepsilon^2} B_2^2d^*K\epsilon_2.$$

To justify the third part of this Lemma, first recall Lemma A.2 that  $\forall k \in [K], \|\mu_s(k) - \mu_t(k)\|_2 \leq \sqrt{d^*}MK\epsilon_1$ . Hence, for any  $i \neq j$ , we have

$$\begin{aligned}
 & |\mu_t(i)^\top \mu_t(j) - \mu_s(i)^\top \mu_s(j)| = |\mu_t(i)^\top \mu_t(j) - \mu_t(i)^\top \mu_s(j) + \mu_t(i)^\top \mu_s(j) - \mu_s(i)^\top \mu_s(j)| \\
 & \leq \|\mu_t(i)\|_2 \|\mu_t(j) - \mu_s(j)\|_2 + \|\mu_s(j)\|_2 \|\mu_t(i) - \mu_s(i)\|_2 \leq 2\sqrt{d^*}B_2MK\epsilon_1,
 \end{aligned}$$

so that we can further yield the relationship of class center divergence between the source domain and the target domain as follows:

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \leq \max_{i \neq j} |\mu_s(i)^\top \mu_s(j)| + 2\sqrt{d^*}B_2MK\epsilon_1. \quad (22)$$

Next, we will attempt to derive an upper bound for  $\max_{i \neq j} |\mu_s(i)^\top \mu_s(j)|$ . Let  $U = (\sqrt{p_s(1)}\mu_s(1), \dots, \sqrt{p_s(K)}\mu_s(K)) \in \mathbb{R}^{d^* \times K}$ , then

$$\begin{aligned}
 & \left\| \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top - I_{d^*} \right\|_F^2 = \|UU^\top - I_{d^*}\|_F^2 \\
 & = \text{Tr}(UU^\top UU^\top - 2UU^\top + I_{d^*}) \quad (\|A\|_F^2 = \text{Tr}(A^\top A))
 \end{aligned}$$

$$\begin{aligned}
 &= \text{Tr}(U^\top U U^\top U - 2U^\top U) + \text{Tr}(I_K) + d^* - K \quad (\text{Tr}(AB) = \text{Tr}(BA)) \\
 &\geq \|U^\top U - I_K\|_F^2 \quad (d^* > K) \\
 &= \sum_{k=1}^K \sum_{l=1}^K (\sqrt{p_s(k)p_s(l)} \mu_s(k)^\top \mu_s(l) - \delta_{kl})^2 \\
 &\geq p_s(i)p_s(j) (\mu_s(i)^\top \mu_s(j))^2.
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 (\mu_s(i)^\top \mu_s(j))^2 &\leq \frac{\left\| \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top - I_{d^*} \right\|_F^2}{p_s(i)p_s(j)} \\
 &= \frac{\left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_2)^\top\} - I_{d^*} + \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_2)^\top\} \right\|_F^2}{p_s(i)p_s(j)} \\
 &\leq \frac{2 \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_2)^\top\} - I_{d^*} \right\|_F^2 + 2 \left\| \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_2)^\top\} \right\|_F^2}{p_s(i)p_s(j)} \quad (23)
 \end{aligned}$$

For the term  $\left\| \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1) f(\mathbf{x}_2)^\top] \right\|_F^2$ , note that

$$\begin{aligned}
 &= \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_1)^\top\} - \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top \\
 &\quad + \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1) \{f(\mathbf{x}_2) - f(\mathbf{x}_1)\}^\top] \\
 &= \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\{f(\mathbf{x}_1) - \mu_s(k)\} \{f(\mathbf{x}_1) - \mu_s(k)\}^\top] + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1) \{f(\mathbf{x}_2) - f(\mathbf{x}_1)\}^\top], \quad (24)
 \end{aligned}$$

where the last equation is derived from

$$\begin{aligned}
 &\mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_1)^\top\} - \mu_s(k) \mu_s(k)^\top = \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_1)^\top\} + \mu_s(k) \mu_s(k)^\top \\
 &\quad - (\mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1)\}) \mu_s(k)^\top - \mu_s(k) (\mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1)\})^\top \\
 &= \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\{f(\mathbf{x}_1) - \mu_s(k)\} \{f(\mathbf{x}_1) - \mu_s(k)\}^\top].
 \end{aligned}$$

So its norm is

$$\begin{aligned}
 &\left\| \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1) f(\mathbf{x}_2)^\top] \right\|_F \\
 &\leq \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\| \{f(\mathbf{x}_1) - \mu_s(k)\} \{f(\mathbf{x}_1) - \mu_s(k)\}^\top \|_F] + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [\| f(\mathbf{x}_1) \{f(\mathbf{x}_2) - f(\mathbf{x}_1)\}^\top \|_F] \\
 &\leq \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \{ \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 \} + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{ \|f(\mathbf{x}_1)\|_2 \|f(\mathbf{x}_2) - f(\mathbf{x}_1)\|_2 \} \\
 &\leq \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \{ \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 \} + \left\{ \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1)\|_2^2 \right\}^{\frac{1}{2}} \left\{ \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_2) - f(\mathbf{x}_1)\|_2^2 \right\}^{\frac{1}{2}} \\
 &\leq \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} [\|f(\mathbf{x}_1) - \mu_s(k)\|_2^2] + B_2 \left[ \varepsilon^2 + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [\|f(\mathbf{x}_2) - f(\mathbf{x}_1)\|_2^2 \mathbb{1}\{\mathbf{x} \notin S_s(\varepsilon, f)\}] \right]^{\frac{1}{2}}
 \end{aligned}$$

$$\begin{aligned}
 & \left( \text{Review that } S_s(\varepsilon, f) := \{\mathbf{x} \in \cup_{k=1}^K C_s(k) : \sup_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 \leq \varepsilon\} \right) \\
 & \leq \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \left\{ \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 \right\} + B_2 \left[ \varepsilon^2 + 4B_2^2 \mathbb{E}_{\mathbf{x}} [\mathbb{1}\{\mathbf{x} \notin S_s(\varepsilon, f)\}] \right]^{\frac{1}{2}} \\
 & = \sum_{k=1}^K p_s(k) \mathbb{E}_{\mathbf{x} \in C_s(k)} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \left[ \|f(\mathbf{x}_1) - \mu_s(k)\|_2^2 \right] + B_2 (\varepsilon^2 + 4B_2^2 R_s(\varepsilon, f))^{\frac{1}{2}} \\
 & \leq 4B_2^2 \sum_{k=1}^K p_s(k) \left\{ \left( 1 - \sigma_s + \frac{\mathcal{K}\delta_s}{2B_2} + \frac{\varepsilon}{B_2} + \frac{R_s(\varepsilon, f)}{p_s(k)} \right)^2 + \left( 1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)} \right) \right\} + B_2 \{ \varepsilon^2 + 4B_2^2 R_s(\varepsilon, f) \}^{\frac{1}{2}} \\
 & = 4B_2^2 \left\{ \left( 1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2} \right)^2 + (1 - \sigma_s) + KR_s(\varepsilon, f) \left( 3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2} \right) + R_s^2(\varepsilon, f) \left( \sum_{k=1}^K \frac{1}{p_s(k)} \right) \right\} \\
 & \quad + B_2 \{ \varepsilon^2 + 4B_2^2 R_s(\varepsilon, f) \}^{\frac{1}{2}} \tag{Lemma A.3}
 \end{aligned}$$

If we define  $\varphi(\sigma_s, \delta_s, \varepsilon, f) := 4B_2^2 \left\{ \left( 1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2} \right)^2 + (1 - \sigma_s) + KR_s(\varepsilon, f) \left( 3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2} \right) + R_s^2(\varepsilon, f) \left( \sum_{k=1}^K \frac{1}{p_s(k)} \right) \right\} + B_2 (\varepsilon^2 + 4B_2^2 R_s(\varepsilon, f))^{\frac{1}{2}}$ , above derivation implies

$$\left\| \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{ f(\mathbf{x}_1) f(\mathbf{x}_2)^\top \} \right\|_F \leq \varphi(\sigma_s, \delta_s, \varepsilon, f). \tag{25}$$

Besides that, Note that

$$\mathcal{R} = \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1) f(\mathbf{x}_2)^\top] - I_{d^*} \right\|_F^2, \tag{26}$$

Combining (23), (24), (25) and (26) yields for any  $i \neq j$

$$(\mu_s(i)^\top \mu_s(j))^2 \leq \frac{2}{p_s(i)p_s(j)} \left\{ \mathcal{R}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f) \right\},$$

which implies that

$$\max_{i \neq j} |\mu_s(i)^\top \mu_s(j)| \leq \sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} \left\{ \mathcal{R}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f) \right\}}.$$

So we can get what we desired according to (22)

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \leq \sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} \left\{ \mathcal{R}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f) \right\}} + 2\sqrt{d^*} B_2 M \mathcal{K} \epsilon_1.$$

□

Next we present the population theorem as follows, which is a direct corollary of Lemma A.4 because of the facts that  $\mathcal{R}(f) \lesssim \mathcal{L}(f)$  and  $\mathcal{L}_{\text{align}}(f) \lesssim \mathcal{L}(f)$ .

**Theorem A.5.** *Given a  $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -augmentation, if  $d^* > K$ , Assumption 3.5 holds and the encoder  $f$  with  $B_1 \leq \|f\|_2 \leq B_2$  is  $\mathcal{K}$ -Lipschitz continuous, then for any  $\varepsilon > 0$ ,*

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \lesssim \sqrt{\mathcal{L}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f)} + \mathcal{K}\epsilon_1.$$

Furthermore, if  $\max_{i \neq j} \mu_t(i)^\top \mu_t(j) < B_2^2 \psi(\sigma_t, \delta_t, \varepsilon, f)$ , then the misclassification rate of  $Q_f$

$$\text{Err}(Q_f) \leq (1 - \sigma_t) + \mathcal{O}(\{\mathcal{L}_{\text{align}}(f) + \mathcal{K}\epsilon_1 + \epsilon_2\}/\varepsilon^2),$$

where the specific formulations of  $\varphi(\sigma_s, \delta_s, \varepsilon, f)$  and  $\psi(\sigma_t, \delta_t, \varepsilon, f)$  can be found in Lemma A.4 and Lemma A.1, respectively.

### A.3. Proof of Theorem 3.9

In this section, we focus on providing the proof of Theorem 3.9. Although Lemma A.4 elucidates some essential factors behind the success of our method, its analysis remains at the population level, leaving the impact of sample size on  $\text{Err}(Q_{\hat{f}_{n_s}})$  unresolved. To further explore this, applying Theorem A.5 yields Lemma A.6, indicating that investigating the sample complexity of  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}$  is a correct direction towards our goal.

However, there are two main challenges that hinder this exploration. The first is  $\mathcal{L}(f^*)$ , which represents the gap between  $\mathcal{L}(\hat{f}_{n_s})$  and the excess risk  $\mathcal{E}(\hat{f}_{n_s})$  defined in the Definition A.13. Since the excess risk can be addressed through typical error decomposition techniques and tools from nonparametric statistics, we aim to construct a measurable function under Assumption 3.3 that makes  $\mathcal{L}(f^*)$  vanish, rather than directly assuming this term can be well-implemented by a specific neural network in HaoChen & Ma (2023).

The second issue stems from bias. To tackle this problem, we develop a novel risk decomposition in Section A.3.3. Utilizing this technique,  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{E}(\hat{f}_{n_s})\}$  can be decomposed into three parts: statistical error:  $\mathcal{E}_{\text{sta}}$ , approximation error brought by  $\mathcal{F}$ :  $\mathcal{E}_{\mathcal{F}}$  and the error induced by using  $\hat{G}(f)$  to approximate  $\mathcal{G}(f)$ :  $\mathcal{E}_G$ . We subsequently address each term in succession. For  $\mathcal{E}_{\text{sta}}$ , we apply standard empirical process techniques and leverage results from Golowich et al. (2018) in A.3.4 to bound it by  $\frac{\mathcal{K}\sqrt{L}}{\sqrt{n_s}}$ . Regarding  $\mathcal{E}_{\mathcal{F}}$ , we first reformulate the problem as a function approximation issue and adopt existing conclusions from Jiao et al. (2023), yielding a bound of  $\mathcal{E}_{\mathcal{F}}$  in Section A.3.5. By leveraging the property  $\mathbb{E}_{\tilde{D}_s} [\hat{\mathcal{L}}(f, G)] = \mathcal{L}(f, G)$ , we transform the problem of bounding  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{E}_G\}$  into a common problem of mean convergence rate, further controlling it by  $\frac{1}{n_s^{1/4}}$  in Section A.3.6.

After completing these preliminaries, we balance these errors to determine appropriate values for the width  $W$ , depth  $D$  and the Lipschitz constant  $\mathcal{K}$  of the neural network while establishing an end-to-end upper bound for  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{E}(\hat{f}_{n_s})\}$ . More details are deferred to Section A.3.7. Finally, Lemma A.20 presents the formal version of Theorem 3.9, with the connection between Lemma A.20 and Theorem 3.9 detailed in A.3.8.

As stated above, we apply Theorem A.5 to the sample optimizer  $\hat{f}_{n_s}$  to yield Lemma A.6.

**Lemma A.6.** *Given a  $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -augmentation, for any  $\varepsilon > 0$ , if  $\psi(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s}) > 0$ , then with probability at least*

$$1 - \frac{\sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} [\frac{1}{\lambda} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + \phi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})] + 2\sqrt{d^*} B_2 M K \epsilon_1}}{B_2^2 \psi(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})}, \text{ we have}$$

$$\mathbb{E}_{\tilde{D}_s} \{\text{Err}(Q_{\hat{f}_{n_s}})\} \leq (1 - \sigma_t) + \frac{m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + 8B_2 d^* M K \epsilon_1 + 4B_2^2 d^* K \epsilon_2},$$

where

$$\begin{aligned} \phi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}) := & B_2 \left( \varepsilon^2 + 4B_2^2 \frac{m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}} \right)^{\frac{1}{2}} + 4B_2^2 \left[ \left( 1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2} \right)^2 \right. \\ & + (1 - \sigma_s) + \frac{K m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}} \left( 3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2} \right) \\ & \left. + \frac{m^4}{\varepsilon^2} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} \left( \sum_{k=1}^K \frac{1}{p_s(k)} \right) \right]. \end{aligned}$$

In addition, the following inequalities always hold

$$\mathbb{E}_{\tilde{D}_s} \{R_t^2(\varepsilon, \hat{f}_{n_s})\} \leq \frac{m^4}{\varepsilon^2} \left( \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + 8B_2 d^* M K \epsilon_1 + 4B_2^2 d^* K \epsilon_2 \right).$$

*Proof.* Applying Lemma A.4 to  $\hat{f}_{n_s}$  yields

$$R_s^2(\varepsilon, \hat{f}_{n_s}) \leq \frac{m^4}{\varepsilon^2} \mathcal{L}(\hat{f}_{n_s}) \tag{27}$$

$$R_t^2(\varepsilon, \hat{f}_{n_s}) \leq \frac{m^4}{\varepsilon^2} \mathcal{L}(\hat{f}_{n_s}) + \frac{8m^4}{\varepsilon^2} B_2 d^* M K \epsilon_1 + \frac{4m^4}{\varepsilon^2} B_2^2 d^* K \epsilon_2 \tag{28}$$

1155 and

$$1156 \max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \leq \sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} \left( \frac{1}{\lambda} \mathcal{L}(\hat{f}_{n_s}) + \varphi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}) \right) + 2\sqrt{d^*} B_2 M K \epsilon_1} \quad (29)$$

1159 Take expectation regarding to  $D_s$  on the both sides of (27), (28) and (29), along with the Jensen's inequality to obtain

$$1160 \mathbb{E}_{\tilde{D}_s} \{R_s^2(\varepsilon, \hat{f}_{n_s})\} \leq \frac{m^4}{\varepsilon^2} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}$$

$$1161 \mathbb{E}_{\tilde{D}_s} \{R_t^2(\varepsilon, \hat{f}_{n_s})\} \leq \frac{m^4}{\varepsilon^2} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + \frac{8m^4}{\varepsilon^2} B_2 d^* M K \epsilon_1 + \frac{4m^4}{\varepsilon^2} B_2^2 d^* K \epsilon_2$$

$$1162 \mathbb{E}_{\tilde{D}_s} [\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)|] \leq \sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} \left( \frac{1}{\lambda} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + \mathbb{E}_{\tilde{D}_s} [\varphi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})] \right) + 2\sqrt{d^*} B_2 M K \epsilon_1}$$

1169 where  $\mathbb{E}_{\tilde{D}_s} \{\varphi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})\} = 4B_2^2 \left[ \left(1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2}\right)^2 + (1 - \sigma_s) + K \mathbb{E}_{\tilde{D}_s} \{R_s(\varepsilon, \hat{f}_{n_s})\} \left(3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2}\right) \right]$   
 1170  $\mathbb{E}_{\tilde{D}_s} \{R_s^2(\varepsilon, \hat{f}_{n_s})\} \left( \sum_{k=1}^K \frac{1}{p_s(k)} \right) + B_2 \mathbb{E}_{\tilde{D}_s} [\{\varepsilon^2 + 4B_2^2 R_s(\varepsilon, \hat{f}_{n_s})\}^{\frac{1}{2}}]$ .

1172 Therefore, by Jensen inequality, we have

$$1173 \mathbb{E}_{\tilde{D}_s} \{\varphi(\sigma_s, \delta_s, \varepsilon, R_s(\varepsilon, \hat{f}_{n_s}))\} \leq 4B_2^2 \left[ \left(1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2}\right)^2 + (1 - \sigma_s) + K \mathbb{E}_{\tilde{D}_s} [R_s(\varepsilon, \hat{f}_{n_s})] \left(3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2}\right) \right]$$

$$1174 + \mathbb{E}_{\tilde{D}_s} [R_s^2(\varepsilon, \hat{f}_{n_s})] \left( \sum_{k=1}^K \frac{1}{p_s(k)} \right) + B_2 [\varepsilon^2 + 4B_2^2 \mathbb{E}_{\tilde{D}_s} \{R_s(\varepsilon, \hat{f}_{n_s})\}]^{\frac{1}{2}}$$

$$1175 \leq 4B_2^2 \left[ \left(1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2}\right)^2 + \frac{K m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}} \left(3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2}\right) + \frac{m^4}{\varepsilon^2} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} \left( \sum_{k=1}^K \frac{1}{p_s(k)} \right) \right]$$

$$1176 + (1 - \sigma_s) + B_2 \left( \varepsilon^2 + \frac{4B_2^2 m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}} \right)^{\frac{1}{2}} := \phi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}).$$

1185 Since Lemma A.1 reveals that if  $\max_{i \neq j} |(\mu_t(i)^\top \mu_t(j))| < B_2^2 \psi(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})$ , then  $\text{Err}(Q_{\hat{f}_{n_s}}) \leq (1 - \sigma_t) +$   
 1186  $R_t(\varepsilon, \hat{f}_{n_s})$ . Thus, if  $\psi(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s}) > 0$ , by Markov inequality, we know that with probability at least  $1 -$   
 1187  $\frac{\sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} \left( \frac{1}{\lambda} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + \phi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}) \right) + 2\sqrt{d^*} B_2 M K \epsilon_1}}{B_2^2 \psi(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})}$ ,  $\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| < B_2^2 \psi(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})$ , which im-  
 1188 plies that

$$1189 \mathbb{E}_{\tilde{D}_s} \{\text{Err}(Q_{\hat{f}_{n_s}})\} \leq (1 - \sigma_t) + R_t(\varepsilon, \hat{f}_{n_s}) \leq (1 - \sigma_t) + \frac{m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + 8B_2 d^* M K \epsilon_1 + 4B_2^2 d^* K \epsilon_2},$$

1190 where the last inequality stems from (28).  $\square$

1196 Therefore, to justify Theorem 3.9, we need to explore the sample complexity of  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}$ . To this end, it is necessary  
 1197 to introduce some basic facts about ACT and learning theory.

### 1199 A.3.1. PRELIMINARIES FOR PROVING THEOREM 3.9

1200 Recall that for any  $\mathbf{x} \in \mathcal{X}_s$ ,  $\mathbf{x}_1, \mathbf{x}_2 \stackrel{\text{i.i.d.}}{\sim} A(\mathbf{x})$ ,  $\tilde{\mathbf{x}} = (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{R}^{2d^*}$ . If we define  $\ell(\tilde{\mathbf{x}}, G) := \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2 +$   
 1201  $\lambda \langle f(\mathbf{x}_1) f(\mathbf{x}_2)^\top - I_{d^*}, G \rangle_F$ , then our loss function at sample level can be rewritten as

$$1202 \hat{\mathcal{L}}(f, G) := \frac{1}{n_s} \sum_{i=1}^{n_s} \left\{ \|f(\mathbf{x}_1^{(i)}) - f(\mathbf{x}_2^{(i)})\|_2^2 + \lambda \langle f(\mathbf{x}_1^{(i)}) f(\mathbf{x}_2^{(i)})^\top - I_{d^*}, G \rangle_F \right\} = \frac{1}{n_s} \sum_{i=1}^{n_s} \ell(\tilde{\mathbf{x}}^{(i)}, G),$$

1206 moreover, let  $\mathcal{G}_1 := \{G \in \mathbb{R}^{d^* \times d^*} : \|G\|_F \leq B_2^2 + \sqrt{d^*}\}$ . It is obvious that both  $\mathcal{G}(f)$  for any  $f : \|f\|_2 \leq B_2$  and  $\hat{\mathcal{G}}(f)$   
 1207 for any  $f \in \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}, B_1, B_2)$  are the subset of  $\mathcal{G}_1$ . In this regard, following Proposition A.7 reveals that  $\ell(\mathbf{u}, G)$   
 1208 is a Lipschitz function on the domain  $\{\mathbf{u} \in \mathbb{R}^{2d^*} : \|\mathbf{u}\|_2 \leq \sqrt{2}B_2\} \times \mathcal{G}_1 \subseteq \mathbb{R}^{2d^* + (d^*)^2}$ .

1210 **Proposition A.7.**  $\ell$  is a Lipschitz function on the domain  $\{\mathbf{u} \in \mathbb{R}^{2d^*} : \|\mathbf{u}\|_2 \leq \sqrt{2}B_2\} \times \mathcal{G}_1$ .

1211  
1212 *Proof.* We begin by proving that  $\|\ell(\cdot, G)\|_{\text{Lip}} < \infty$  for any fixed  $G \in \mathcal{G}_1$ . To this end, let  $\mathbf{u} = (\mathbf{u}_1, \mathbf{u}_2)$ , where  
1213  $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^{d^*}$ . We first demonstrate that  $J(\mathbf{u}) = \|\mathbf{u}_1 - \mathbf{u}_2\|_2^2$  is a Lipschitz function. Define  $g(\mathbf{u}) := \mathbf{u}_1 - \mathbf{u}_2$ . We have:

$$\begin{aligned} 1215 \quad \|g(\mathbf{u}) - g(\mathbf{v})\|_2^2 &= \|\mathbf{u}_1 - \mathbf{u}_2 - \mathbf{v}_1 + \mathbf{v}_2\|_2^2 \leq (\|\mathbf{u}_1 - \mathbf{v}_1\|_2 + \|\mathbf{u}_2 - \mathbf{v}_2\|_2)^2 \\ 1216 &= \|\mathbf{u}_1 - \mathbf{v}_1\|_2^2 + \|\mathbf{u}_2 - \mathbf{v}_2\|_2^2 + 2\|\mathbf{u}_1 - \mathbf{v}_1\|_2\|\mathbf{u}_2 - \mathbf{v}_2\|_2 \\ 1217 &\leq 2(\|\mathbf{u}_1 - \mathbf{v}_1\|_2^2 + \|\mathbf{u}_2 - \mathbf{v}_2\|_2^2) = 2\|\mathbf{u} - \mathbf{v}\|_2^2, \end{aligned}$$

1220 which implies that  $g(\mathbf{u}) \in \text{Lip}(\sqrt{2})$ . Furthermore,  $g$  possesses the property that  $\|g(\mathbf{u})\|_2 = \|\mathbf{u}_1 - \mathbf{u}_2\|_2 \leq \|\mathbf{u}_1\|_2 + \|\mathbf{u}_2\|_2 \leq$   
1221  $2\|\mathbf{u}\|_2 \leq 2\sqrt{2}B_2$ . Next, let  $h(\mathbf{v}) := \|\mathbf{v}\|_2^2$ . We have:

$$1223 \quad \left\| \frac{\partial h}{\partial \mathbf{v}}(g(\mathbf{u})) \right\|_2 = 2\|g(\mathbf{u})\|_2 \leq 4\sqrt{2}B_2.$$

1226 Thus,  $J(\mathbf{u}) = h(g(\mathbf{u})) = \|\mathbf{u}_1 - \mathbf{u}_2\|_2^2 \in \text{Lip}(8B_2)$ . Now, we show that  $Q(\mathbf{u}) = \langle \mathbf{u}_1 \mathbf{u}_2^\top - I_{d^*}, G \rangle_F$  is also a Lipschitz  
1227 function. Define  $\tilde{g}(\mathbf{u}) := \mathbf{u}_1 \mathbf{u}_2^\top$ . We have:

$$\begin{aligned} 1229 \quad \|\tilde{g}(\mathbf{u}) - \tilde{g}(\mathbf{v})\|_F &= \|\mathbf{u}_1 \mathbf{u}_2^\top - \mathbf{v}_1 \mathbf{v}_2^\top\|_F = \|\mathbf{u}_1 \mathbf{u}_2^\top - \mathbf{u}_1 \mathbf{v}_2^\top + \mathbf{u}_1 \mathbf{v}_2^\top - \mathbf{v}_1 \mathbf{v}_2^\top\|_F \\ 1230 &= \|\mathbf{u}_1(\mathbf{u}_2 - \mathbf{v}_2)^\top + (\mathbf{u}_1 - \mathbf{v}_1)\mathbf{v}_2^\top\|_F \leq \|\mathbf{u}_1\|_F \|\mathbf{u}_2 - \mathbf{v}_2\|_F + \|\mathbf{u}_1 - \mathbf{v}_1\|_F \|\mathbf{v}_2\|_F \\ 1231 &\leq (\|\mathbf{u}_1\|_2 + \|\mathbf{v}_2\|_2) \|\mathbf{u} - \mathbf{v}\|_2 \leq 2\sqrt{2}B_2 \|\mathbf{u} - \mathbf{v}\|_2. \end{aligned}$$

1234 Subsequently, denote  $\tilde{h}(A) := \langle A - I_{d^*}, G \rangle_F$ . Then, we find that  $\|\nabla \tilde{h}(A)\|_F = \|G\|_F \leq B_2^2 + \sqrt{d^*}$ . Therefore, we  
1235 conclude that  $Q(\mathbf{u}) = \tilde{h}(\tilde{g}(\mathbf{u})) \in \text{Lip}(2\sqrt{2}B_2(B_2^2 + \sqrt{d^*}))$ . Combining the above results, we establish that for any  $G \in \mathcal{G}_1$ ,  
1236 we have  $\|\ell(\cdot, G)\|_{\text{Lip}} < \infty$  on the domain  $\{\mathbf{u} : \|\mathbf{u}\|_2 \leq \sqrt{2}B_2\}$ . Next, for a fixed  $\mathbf{u} \in \mathbb{R}^{2d^*}$  such that  $\|\mathbf{u}\|_2 \leq \sqrt{2}B_2$ , we  
1237 obtain:

$$1239 \quad |\ell(\mathbf{u}, G_1) - \ell(\mathbf{u}, G_2)| = |\langle \mathbf{u}, G_1 - G_2 \rangle_F| \leq \|\mathbf{u}\|_2 \|G_1 - G_2\|_F = \sqrt{2}B_2 \|G_1 - G_2\|_F,$$

1241 which implies that  $\ell(\mathbf{u}, \cdot) \in \text{Lip}(\sqrt{2}B_2)$ . Finally, we note that:

$$\begin{aligned} 1243 \quad |\ell(\mathbf{u}_1, G_1) - \ell(\mathbf{u}_2, G_2)|^2 &\leq \{|\ell(\mathbf{u}_1, G_1) - \ell(\mathbf{u}_2, G_1)| + |\ell(\mathbf{u}_2, G_1) - \ell(\mathbf{u}_2, G_2)|\}^2 \\ 1244 &\leq \left[ \{\sqrt{2} + 2\sqrt{2}B_2(B_2^2 + \sqrt{d^*})\} \|\mathbf{u}_1 - \mathbf{u}_2\|_2 + \sqrt{2}B_2 \|G_1 - G_2\|_F \right]^2 \\ 1245 &\leq 2\{\sqrt{2} + 2\sqrt{2}B_2(B_2^2 + \sqrt{d^*})\}^2 \|\mathbf{u}_1 - \mathbf{u}_2\|_2^2 + 4B_2^2 \|G_1 - G_2\|_F^2 \\ 1246 &\leq C \|\text{vec}(\mathbf{u}_1, G_1) - \text{vec}(\mathbf{u}_2, G_2)\|_2^2, \end{aligned}$$

1250 where  $C$  is a constant such that  $C \geq \max\{2\{\sqrt{2} + 2\sqrt{2}B_2(B_2^2 + \sqrt{d^*})\}^2, 4B_2^2\}$ , thus yielding the desired result.  $\square$

1252 We summary the Lipschitz constants of  $\ell(\mathbf{u}, G)$  with respect to  $\mathbf{u} \in \{\mathbf{u} \in \mathbb{R}^{2d^*} : \|\mathbf{u}\|_2 \leq \sqrt{2}B_2\}$  and  $G \in \mathcal{G}_1$  in Table 3.

Table 3. Lipschitz constant of  $\ell$  with respect to each component

Function	Lipschitz Constant
$\ell(\mathbf{u}, \cdot)$	$\sqrt{2}B_2$
$\ell(\cdot, G)$	$2\sqrt{2}B_2(B_2^2 + \sqrt{d^*})$
$\ell(\cdot)$	$\max\{\sqrt{2}B_2, 2\sqrt{2}B_2(B_2^2 + \sqrt{d^*})\}$

1263 Following Definition A.8, A.10 and Lemma A.9, A.11 are all typical elements in the area of learning theory.

1265 **Definition A.8** (Rademacher complexity). Given a set  $S \subseteq \mathbb{R}^n$ , the Rademacher complexity of  $S$  is denoted by

$$1266 \mathcal{R}_n(S) := \mathbb{E}_\xi \left[ \sup_{(s_1, \dots, s_n) \in S} \frac{1}{n} \sum_{i=1}^n \xi_i s_i \right],$$

1270 where  $\{\xi_i\}_{i \in [n]}$  is a sequence of i.i.d Radmacher random variables which take the values 1 and  $-1$  with equal probability  
1271  $1/2$ .

1273 Following vector-contraction principle of Rademacher complexity will be used in later contents.

1274 **Lemma A.9** (Vector-contraction principle). Let  $\mathcal{X}$  be any set,  $(x_1, \dots, x_n) \in \mathcal{X}^n$ , let  $F$  be a class of functions  $f : \mathcal{X} \rightarrow \ell_2$   
1275 and let  $h_i : \ell_2 \rightarrow \mathbb{R}$  have Lipschitz norm  $L$ . Then

$$1277 \mathbb{E} \sup_{f \in F} \left| \sum_i \epsilon_i h_i(f(x_i)) \right| \leq 2\sqrt{2}L \mathbb{E} \sup_{f \in F} \left| \sum_{i,k} \epsilon_{ik} f_k(x_i) \right|,$$

1280 where  $\epsilon_{ik}$  is an independent doubly indexed Rademacher sequence and  $f_k(x_i)$  is the  $k$ -th component of  $f(x_i)$ .

1283 *Proof.* Combining Maurer (2016) and Theorem 3.2.1 of Giné & Nickl (2016) obtains the desired result.  $\square$

1285 **Definition A.10** (Covering number).  $\forall n \in \mathbb{N}^+$ , Fix  $S \subseteq \mathbb{R}^n$  and  $\rho > 0$ , the set  $\mathcal{N}$  is called an  $\rho$ -net of  $S$  with respect to a  
1286 norm  $\|\cdot\|$  on  $\mathbb{R}^n$ , if  $\mathcal{N} \subseteq S$  and for any  $\mathbf{u} \in S$ , there exists  $\mathbf{v} \in \mathcal{N}$  such that  $\|\mathbf{u} - \mathbf{v}\| \leq \rho$ . The covering number of  $S$  is  
1287 defined as

$$1289 \mathcal{N}(S, \|\cdot\|, \rho) := \min\{|\mathcal{Q}| : \mathcal{Q} \text{ is an } \rho\text{-cover of } S\}$$

1291 where  $|\mathcal{Q}|$  is the cardinality of the set  $\mathcal{Q}$ .

1293 According to the Corollary 4.2.13 of Vershynin (2018),  $|\mathcal{N}(\mathcal{B}_2, \|\cdot\|_2, \rho)|$ , which is the the covering number of 2-norm unit  
1294 ball in  $\mathbb{R}^{(d^*)^2}$ , can be bounded by  $(\frac{3}{\rho})^{(d^*)^2}$ , so that if we denote  $\mathcal{N}_{\mathcal{G}_1}(\rho)$  is a cover of  $\mathcal{G}_1$  with radius  $\rho$  whose cardinality  
1295  $|\mathcal{N}_{\mathcal{G}_1}(\rho)|$  is equal to the covering number of  $\mathcal{G}_1$ , then  $|\mathcal{N}_{\mathcal{G}_1}(\rho)| \leq \left\{ \frac{3}{(B_2^2 + \sqrt{d^*})\rho} \right\}^{(d^*)^2}$ .

1297 **Lemma A.11** (Finite maximum inequality). For any  $N \geq 1$ , if  $X_i, i \leq N$ , are sub-Gaussian random variables admitting  
1298 constants  $\sigma_i$ , then

$$1300 \mathbb{E} \max_{i \leq N} |X_i| \leq \sqrt{2 \log 2N} \max_{i \leq N} \sigma_i$$

1303 The proof of this lemma can be found in Giné & Nickl (2016), Lemma 2.3.4.

1304 Recall  $\mathcal{NN}_{d_1, d_2}(W, L, \mathcal{K}) := \{f_\theta(\mathbf{x}) = A_L \sigma(A_{L-1} \sigma(\dots \sigma(A_0 \mathbf{x}))) : \kappa(\theta) \leq \mathcal{K}\}$ , as defined in eq 10. The second  
1305 lemma we will employ is related to the upper bound for the Rademacher complexity of the hypothesis space consisting of  
1306 norm-constrained neural networks, which was provided by Golowich et al. (2018).

1308 **Lemma A.12** (Theorem 3.2 of Golowich et al. (2018)).  $\forall n \in \mathbb{N}^+, \forall \mathbf{x}_1, \dots, \mathbf{x}_n \in [-B, B]^d$  with  $B \geq 1, S :=$   
1309  $\{(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) : f \in \mathcal{NN}_{d,1}(W, L, \mathcal{K})\} \subseteq \mathbb{R}^n$ , then

$$1311 \mathcal{R}_n(S) \leq \frac{1}{n} \mathcal{K} \sqrt{2(L+2+\log(d+1))} \max_{1 \leq j \leq d+1} \sqrt{\sum_{i=1}^n x_{i,j}^2} \leq \frac{BK \sqrt{2(L+2+\log(d+1))}}{\sqrt{n}},$$

1314 where  $x_{i,j}$  is the  $j$ -th coordinate of the vector  $(\mathbf{x}_i^\top, 1)^\top \in \mathbb{R}^{d+1}$ .

1316 **Definition A.13** (Excess risk). The difference between  $\mathcal{L}(\hat{f}_{n_s})$  and  $\mathcal{L}(f^*)$  is called excess risk, i.e.,

$$1318 \mathcal{E}(\hat{f}_{n_s}) = \mathcal{L}(\hat{f}_{n_s}) - \mathcal{L}(f^*).$$

1320 A.3.2. DEAL WITH  $\mathcal{L}(f^*)$ 

1321 Since our objective is to explore the sample complexity of  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}$ , it is essential to assert that  $\mathcal{L}(f^*) = 0$ . This ensures  
 1322 that the tools used to analyze  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{E}(\hat{f}_{n_s})\}$  are also applicable for handling  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}$ . The justification comprises a  
 1323 total of two steps. First, we assert that if there exists a measurable map  $f$  such that  $\Sigma = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_s} [f(\mathbf{x})f(\mathbf{x})^\top]$  be positive  
 1324 definite, then we can make minor modifications to obtain  $\tilde{f}$  such that  $\mathcal{L}(\tilde{f}) = 0$ . In the second step, we will demonstrate that  
 1325 the required  $f$  exists under Assumption 3.3, and that the modified  $\tilde{f}$  also satisfies the condition  $B_1 \leq \|\tilde{f}\|_2 \leq B_2$ , which  
 1326 implies that  $\mathcal{L}(f^*) = 0$ , since the definition of  $f^*$  indicates that  $\mathcal{L}(f^*) \leq \mathcal{L}(\tilde{f})$ .

1327 Our final target is to result in a measurable map  $f$ , s.t  $B_1 \leq \|f\|_2 \leq B_2$  and  $\sup_{f \in \mathcal{G}(f)} \mathcal{L}(f) = 0$ , it suffices to find a  
 1328  $f : B_1 \leq \|f\|_2 \leq B_2$  satisfying both  $\mathcal{L}_{\text{align}}(f) = 0$  and  $\left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1)f(\mathbf{x}_2)^\top] - I_{d^*} \right\|_F = 0$ . Note that

$$\begin{aligned}
 & \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1)f(\mathbf{x}_2)^\top\} - I_{d^*} \right\|_F \\
 &= \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1)f(\mathbf{x}_1)^\top\} + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1)\{f(\mathbf{x}_2) - f(\mathbf{x}_1)\}^\top] - I_{d^*} \right\|_F \\
 &\leq \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1)f(\mathbf{x}_1)^\top\} - I_{d^*} \right\|_F + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} \{\|f(\mathbf{x}_1)\|_2 \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2\} \\
 &\leq \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}')f(\mathbf{x}')^\top\} - I_{d^*} \right\|_F + B_2 \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2. \quad (\|f\|_2 \leq B_2)
 \end{aligned}$$

1340 The above deduction indicates that finding a measurable map  $f$  such that  $B_1 \leq \|f\|_2 \leq B_2$  and ensuring both  $\mathcal{L}_{\text{align}}(f)$   
 1341 and  $\left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}')f(\mathbf{x}')^\top\} - I_{d^*} \right\|_F$  vanish is sufficient to achieve our goal.

1342 **Lemma A.14.** *If there exists a measurable map  $f$  making  $\Sigma = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_s} \{f(\mathbf{x})f(\mathbf{x})^\top\}$  positive definite, then there exists a  
 1343 measurable map  $\tilde{f}$  such that*

$$\mathcal{L}_{\text{align}}(\tilde{f}) = 0, \quad \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} \{\tilde{f}(\mathbf{x}')\tilde{f}(\mathbf{x}')^\top\} - I_{d^*} \right\|_F = 0.$$

1344 *Proof.* We conduct modifications for given  $f$  as follows: For any  $\mathbf{x} \in \mathcal{X}$ , define

$$\tilde{f}_{\mathbf{x}}(\mathbf{x}') = \begin{cases} V^{-1}f(\mathbf{x}) & \text{if } \mathbf{x}' \in \mathcal{A}(\mathbf{x}) \\ f(\mathbf{x}) & \text{if } \mathbf{x}' \notin \mathcal{A}(\mathbf{x}) \end{cases}$$

1345 where  $\Sigma = VV^\top$ , which is the Cholesky decomposition of  $\Sigma$ , which is evident well-defined as  $\Sigma$  is positive definite.  
 1346 Iteratively repeat this argument for all  $\mathbf{x} \in \mathcal{X}$  to yield  $\tilde{f}$ , then we have

$$\mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} \{\tilde{f}(\mathbf{x}')\tilde{f}(\mathbf{x}')^\top\} = V^{-1} \mathbb{E}_{\mathbf{x}} \{f(\mathbf{x})f(\mathbf{x})^\top\} V^{-T} = I_{d^*}$$

1347 and

$$\forall \mathbf{x} \in \mathcal{X}, \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x}), \|\tilde{f}(\mathbf{x}_1) - \tilde{f}(\mathbf{x}_2)\|_2 = \|f(\mathbf{x}) - f(\mathbf{x})\|_2 = 0.$$

1348 That is precisely what we desire. □

1349 **Remark A.15.** If we have a measurable partition  $\mathcal{X} = \cup_{i=1}^{d^*} \mathcal{P}_i$  stated in Assumption 3.3 such that  $\mathcal{P}_i \cap \mathcal{P}_j = \emptyset$  and  
 1350  $\forall i \in [d^*], \frac{1}{B_2^2} \leq \mathbb{P}_s(\mathcal{P}_i) \leq \frac{1}{B_1^2}$ , just set the  $f(\mathbf{x}) = \mathbf{e}_i$  if  $\mathbf{x} \in \mathcal{P}_i$ , where  $\mathbf{e}_i$  is the standard basis of  $\mathbb{R}^{d^*}$ , then  $\Sigma =$   
 1351  $\text{diag}\{\mathbb{P}_s(\mathcal{P}_1), \dots, \mathbb{P}_s(\mathcal{P}_i), \dots, \mathbb{P}_s(\mathcal{P}_{d^*})\}$ ,  $V^{-1} = \text{diag}\{\sqrt{\frac{1}{\mathbb{P}_s(\mathcal{P}_1)}}, \dots, \sqrt{\frac{1}{\mathbb{P}_s(\mathcal{P}_i)}}, \dots, \sqrt{\frac{1}{\mathbb{P}_s(\mathcal{P}_{d^*})}}\}$ ,  $\tilde{f}(\mathbf{x}) = \sqrt{\frac{1}{\mathbb{P}_s(\mathcal{P}_i)}} \mathbf{e}_i$  if  
 1352  $\mathbf{x} \in \mathcal{P}_i$ , it is obviously that  $B_1 \leq \|\tilde{f}\|_2 \leq B_2$ .

1353 In this context, exploring the sample complexity of  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{E}(\hat{f}_{n_s})\}$  is equivalent to investigating  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}$ . However,  
 1354 the unbiasedness between  $\hat{\mathcal{L}}(f)$  and  $\mathcal{L}(f)$  hinders our ability to analyze  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{E}(\hat{f}_{n_s})\}$ . To address this issue, we develop  
 1355 the following novel risk decomposition.

## 1375 A.3.3. RISK DECOMPOSITION

 1376 If denote  $\widehat{G}(f) = \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{x}_1^{(i)})f(\mathbf{x}_2^{(i)})^\top - I_{d^*}$  and  $G^*(f) = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1)f(\mathbf{x}_2)^\top] - I_{d^*}$ , we can decompose  
 1377  $\mathcal{E}(\widehat{f}_{n_s})$  into three terms shown as follow and then deal each term successively. To achieve conciseness in subsequent  
 1378 conclusions, we employ  $X \lesssim Y$  or  $Y \gtrsim X$  to indicate the statement that  $X \leq CY$  form some  $C > 0$  if  $X$  and  $Y$  are two  
 1379 quantities.  
 1380

 1381 **Lemma A.16.** *The excess risk  $\mathcal{E}(\widehat{f}_{n_s})$  satisfies*

1382 
$$\mathbb{E}_{\widetilde{D}_s} \{\mathcal{E}(\widehat{f}_{n_s})\} \lesssim \underbrace{2 \mathbb{E}_{\widetilde{D}_s} \left\{ \sup_{f \in \mathcal{F}, G \in \widehat{\mathcal{G}}(f)} |\mathcal{L}(f, G) - \widehat{\mathcal{L}}(f, G)| \right\}}_{\text{statistical error : } \mathcal{E}_{\text{sta}}} + \underbrace{\inf_{f \in \mathcal{F}} \{\mathcal{L}(f) - \mathcal{L}(f^*)\}}_{\text{approximation error of } \mathcal{F} : \mathcal{E}_{\mathcal{F}}} + \underbrace{\mathbb{E}_{\widetilde{D}_s} \left[ \sup_{f \in \mathcal{F}} \{G^*(f) - \widehat{G}(f)\} \right]}_{\text{approximation error of } \mathcal{G} : \mathcal{E}_{\mathcal{G}}},$$

1387 That is,

1388 
$$\mathcal{E}(\widehat{f}_{n_s}) \leq \mathcal{E}_{\text{sta}} + \mathcal{E}_{\mathcal{F}} + \mathcal{E}_{\mathcal{G}}.$$

 1391 *Proof.* Recall  $\mathcal{F} = \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}, B_1, B_2)$ , for any  $f \in \mathcal{F}$ ,

1392 
$$\begin{aligned} \mathcal{L}(\widehat{f}_{n_s}) - \mathcal{L}(f^*) &= \sup_{G \in \mathcal{G}(\widehat{f}_{n_s})} \mathcal{L}(\widehat{f}_{n_s}, G) - \sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G) \\ &= \left[ \sup_{G \in \mathcal{G}(\widehat{f}_{n_s})} \mathcal{L}(\widehat{f}_{n_s}, G) - \sup_{G \in \widehat{\mathcal{G}}(\widehat{f}_{n_s})} \mathcal{L}(\widehat{f}_{n_s}, G) \right] + \left[ \sup_{G \in \widehat{\mathcal{G}}(\widehat{f}_{n_s})} \mathcal{L}(\widehat{f}_{n_s}, G) - \sup_{G \in \widehat{\mathcal{G}}(\widehat{f}_{n_s})} \widehat{\mathcal{L}}(\widehat{f}_{n_s}, G) \right] \\ &\quad + \left[ \sup_{G \in \widehat{\mathcal{G}}(\widehat{f}_{n_s})} \widehat{\mathcal{L}}(\widehat{f}_{n_s}, G) - \sup_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{L}}(f, G) \right] + \left[ \sup_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{L}}(f, G) - \sup_{G \in \widehat{\mathcal{G}}(f)} \mathcal{L}(f, G) \right] \\ &\quad + \left[ \sup_{G \in \widehat{\mathcal{G}}(f)} \mathcal{L}(f, G) - \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) \right] + \left[ \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) - \sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G) \right], \end{aligned}$$

 1403 where the second and fourth terms can be bounded by  $\mathcal{E}_{\text{sta}}$ . In fact, regarding to the fourth term, we have

1404 
$$\begin{aligned} \sup_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{L}}(f, G) - \sup_{G \in \widehat{\mathcal{G}}(f)} \mathcal{L}(f, G) &\leq \sup_{G \in \widehat{\mathcal{G}}(f)} \{\widehat{\mathcal{L}}(f, G) - \mathcal{L}(f, G)\} \leq \sup_{G \in \widehat{\mathcal{G}}(f)} |\widehat{\mathcal{L}}(f, G) - \mathcal{L}(f, G)| \\ &\leq \sup_{f \in \mathcal{F}, G \in \widehat{\mathcal{G}}(f)} |\widehat{\mathcal{L}}(f, G) - \mathcal{L}(f, G)|, \end{aligned}$$

1409 and the same conclusion holds for the second term.

 1411 The summation of first term and fifth term can be bounded by  $\mathcal{E}_{\mathcal{G}}$ . Actually, for the first term

1412 
$$\begin{aligned} \sup_{G \in \mathcal{G}(\widehat{f}_{n_s})} \mathcal{L}(\widehat{f}_{n_s}, G) - \sup_{G \in \widehat{\mathcal{G}}(\widehat{f}_{n_s})} \mathcal{L}(\widehat{f}_{n_s}, G) &\leq \sup_{f \in \mathcal{F}} \left\{ \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) - \sup_{G \in \widehat{\mathcal{G}}(f)} \mathcal{L}(f, G) \right\} \\ &\leq \sup_{f \in \mathcal{F}} \left\{ \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) - \mathcal{L}(f, \widehat{G}(f)) \right\} = \sup_{f \in \mathcal{F}} \left\{ \mathcal{L}(f, G^*(f)) - \mathcal{L}(f, \widehat{G}(f)) \right\} \\ &\leq \sqrt{2} B_2 \sup_{f \in \mathcal{F}} \|G^*(f) - \widehat{G}(f)\|_F \leq \sqrt{2} B_2 \sup_{f \in \mathcal{F}} \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} [f(\mathbf{x}_1)f(\mathbf{x}_2)^\top] - \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{x}_1^{(i)})f(\mathbf{x}_2^{(i)})^\top \right\|_F. \quad (30) \end{aligned}$$

 1420 where the second inequity stems from  $\widehat{G}(f) \in \widehat{\mathcal{G}}(f)$  and the third inequality is due to  $\ell(\mathbf{u}, \cdot) \in \text{Lip}(\sqrt{2}B_2)$ , as outlined in  
 1421 Table 3. and for the fifth term, we turn out

1422 
$$\begin{aligned} \sup_{G \in \widehat{\mathcal{G}}(f)} \mathcal{L}(f, G) - \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) &= \sup_{G \in \widehat{\mathcal{G}}(f)} \mathbb{E}_{\widetilde{D}_s} \{ \langle \widehat{G}(f), G \rangle_F \} - \sup_{G \in \mathcal{G}(f)} \langle G^*(f), G \rangle_F \\ &\leq \mathbb{E}_{\widetilde{D}_s} \left\{ \sup_{G \in \widehat{\mathcal{G}}(f)} \langle \widehat{G}(f), G \rangle_F \right\} - \sup_{G \in \mathcal{G}(f)} \langle G^*(f), G \rangle_F = \mathbb{E}_{\widetilde{D}_s} \{ \|\widehat{G}(f)\|_F^2 \} - \|G^*(f)\|_F^2 \\ &\leq 2(B_2^2 + \sqrt{d^*}) (\mathbb{E}_{\widetilde{D}_s} \{ \|\widehat{G}(f)\|_F \} - \|G^*(f)\|_F) \leq 2(B_2^2 + \sqrt{d^*}) \left( \sup_{f \in \mathcal{F}} [\mathbb{E}_{\widetilde{D}_s} \{ \|\widehat{G}(f)\|_F \}] - \|G^*(f)\|_F \right) \end{aligned}$$

$$\begin{aligned}
 & \lesssim \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{\tilde{D}_s} \left[ \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{x}_1^{(i)}) f(\mathbf{x}_2^{(i)})^\top - I_{d^*} \right\|_F - \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_2)^\top\} - I_{d^*} \right\|_F \right] \right\} \\
 & \leq \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{\tilde{D}_s} \left[ \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{x}_1^{(i)}) f(\mathbf{x}_2^{(i)})^\top - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_2)^\top\} \right\|_F \right] \right\} \\
 & \leq \mathbb{E}_{\tilde{D}_s} \left[ \sup_{f \in \mathcal{F}} \left\{ \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\mathbf{x}_1^{(i)}) f(\mathbf{x}_2^{(i)})^\top - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1) f(\mathbf{x}_2)^\top\} \right\|_F \right\} \right] \tag{31}
 \end{aligned}$$

where the first equality is due to  $\langle G^*(f), G \rangle_F = \mathbb{E}_{\tilde{D}_s} \{ \langle \hat{G}(f), G \rangle_F \}$  and the second inequality is derived from the facts that  $\|\hat{G}(f)\|_F \leq B_2^2 + \sqrt{d^*}$  and  $\|G^*(f)\|_F \leq B_2^2 + \sqrt{d^*}$ . Combining (30) and (31) yields  $\mathbb{E}_{\tilde{D}_s} \{ \mathcal{E}_{\mathcal{G}} \}$ .

Furthermore, the third term  $\sup_{G \in \hat{\mathcal{G}}(\hat{f}_{n_s})} \hat{\mathcal{L}}(\hat{f}_{n_s}, G) - \sup_{G \in \hat{\mathcal{G}}(f)} \hat{\mathcal{L}}(f, G) \leq 0$  because of the definition of  $\hat{f}_{n_s}$ . Taking infimum over all  $f \in \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}, B_1, B_2)$  yields

$$\mathcal{E}(\hat{f}_{n_s}) \lesssim \mathcal{E}_{\text{sta}} + \mathcal{E}_{\mathcal{F}} + \mathcal{E}_{\mathcal{G}},$$

which completes the proof.  $\square$

#### A.3.4. BOUND $\mathcal{E}_{\text{sta}}$

**Lemma A.17.** *Regarding to  $\mathcal{E}_{\text{sta}}$ , we have*

$$\mathbb{E}_{\tilde{D}_s} [\mathcal{E}_{\text{sta}}] \lesssim \frac{\mathcal{K} \sqrt{L}}{\sqrt{n_s}}.$$

*Proof.* We are going to be introducing the relevant notations at first.

For any  $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d^*}$ , let  $\tilde{f} : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d^*}$  such that  $\tilde{f}(\tilde{\mathbf{x}}) = (f(\mathbf{x}_1), f(\mathbf{x}_2))$ , where  $\tilde{\mathbf{x}} = (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{R}^{2d}$ . Furthermore, define  $\tilde{\mathcal{F}} := \{ \tilde{f} : f \in \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}) \}$  and denote  $D'_s = \{ \tilde{\mathbf{x}}^{(i)} \}_{i=1}^{n_s}$  as an independent identically distributed samples to  $D_s$ , which is called as ghost samples of  $D_s$ .

Next, we attempt to establish the relationship between  $\mathbb{E}_{\tilde{D}_s} [\mathcal{E}_{\text{sta}}]$  and the Rademacher complexity of  $\mathcal{NN}_{d, d^*}(W, L, \mathcal{K})$ . By the definition of  $\mathcal{E}_{\text{sta}}$ , we have

$$\begin{aligned}
 \mathbb{E}_{\tilde{D}_s} [\mathcal{E}_{\text{sta}}] &= \mathbb{E}_{\tilde{D}_s} \left[ \sup_{f \in \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}, B_1, B_2), G \in \hat{\mathcal{G}}(f)} |\mathcal{L}(f, G) - \hat{\mathcal{L}}(f, G)| \right] \\
 &\leq \mathbb{E}_{\tilde{D}_s} \left[ \sup_{(f, G) \in \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}, B_1, B_2) \times \mathcal{G}_1} |\mathcal{L}(f, G) - \hat{\mathcal{L}}(f, G)| \right] \\
 &\hspace{15em} (\text{As } \hat{\mathcal{G}}(f) \subseteq \mathcal{G}_1 \text{ for any } f \in \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}, B_1, B_2)) \\
 &\leq \mathbb{E}_{\tilde{D}_s} \left[ \sup_{(f, G) \in \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}) \times \mathcal{G}_1} |\mathcal{L}(f, G) - \hat{\mathcal{L}}(f, G)| \right] \\
 &\hspace{15em} (\text{As } \mathcal{NN}_{d, d^*}(W, L, \mathcal{K}, B_1, B_2) \subseteq \mathcal{NN}_{d, d^*}(W, L, \mathcal{K})) \\
 &= \mathbb{E}_{\tilde{D}_s} \left[ \sup_{(\tilde{f}, G) \in \tilde{\mathcal{F}} \times \mathcal{G}_1} \left| \frac{1}{n_s} \sum_{i=1}^{n_s} \mathbb{E}_{D'_s} [\ell(\tilde{f}(\tilde{\mathbf{x}}'^{(i)}), G)] - \frac{1}{n_s} \sum_{i=1}^{n_s} \ell(\tilde{f}(\tilde{\mathbf{x}}^{(i)}), G) \right| \right] \\
 &\leq \mathbb{E}_{D_s, D'_s} \left[ \sup_{(\tilde{f}, G) \in \tilde{\mathcal{F}} \times \mathcal{G}_1} \left| \frac{1}{n_s} \sum_{i=1}^{n_s} \ell(\tilde{f}(\tilde{\mathbf{x}}'^{(i)}), G) - \frac{1}{n_s} \sum_{i=1}^{n_s} \ell(\tilde{f}(\tilde{\mathbf{x}}^{(i)}), G) \right| \right] \\
 &= \mathbb{E}_{D_s, D'_s, \xi} \left[ \sup_{(\tilde{f}, G) \in \tilde{\mathcal{F}} \times \mathcal{G}_1} \left| \frac{1}{n_s} \sum_{i=1}^{n_s} \xi_i (\ell(\tilde{f}(\tilde{\mathbf{x}}'^{(i)}), G) - \ell(\tilde{f}(\tilde{\mathbf{x}}^{(i)}), G)) \right| \right] \tag{32} \\
 &\leq 2 \mathbb{E}_{D_s, \xi} \left[ \sup_{(\tilde{f}, G) \in \tilde{\mathcal{F}} \times \mathcal{G}_1} \left| \frac{1}{n_s} \sum_{i=1}^{n_s} \xi_i \ell(\tilde{f}(\tilde{\mathbf{x}}^{(i)}), G) \right| \right]
 \end{aligned}$$



1540 A.3.5. BOUND  $\mathcal{E}_{\mathcal{F}}$ 

1541 If we denote

1542 
$$\mathcal{E}(\mathcal{H}^\alpha, \mathcal{NN}_{d,1}(W, L, \mathcal{K})) := \sup_{g \in \mathcal{H}^\alpha} \inf_{f \in \mathcal{NN}_{d,1}(W, L, \mathcal{K})} \|f - g\|_{C([0,1]^d)},$$

 1543 where  $C([0,1]^d)$  is the space of continuous functions on  $[0,1]^d$  equipped with the sup-norm. Theorem 3.2 of [Jiao et al. \(2023\)](#) has already proven  $\mathcal{E}(\mathcal{H}^\alpha, \mathcal{NN}_{d,1}(W, L, \mathcal{K}))$  can be bound by a quantity related to  $\mathcal{K}$  when setting appropriate architecture of network, that is

 1544 **Lemma A.18** (Theorem 3.2 of [Jiao et al. \(2023\)](#)). *Let  $d \in \mathbb{N}$  and  $\alpha = r + \beta > 0$ , where  $r \in \mathbb{N}_0$  and  $\beta \in (0, 1]$ . There exists  $c > 0$  such that for any  $\mathcal{K} \geq 1$ , any  $W \geq c\mathcal{K}^{(2d+\alpha)/(2d+2)}$  and  $L \geq 2\lceil \log_2(d+r) \rceil + 2$ ,*

1545 
$$\mathcal{E}(\mathcal{H}^\alpha, \mathcal{NN}_{d,1}(W, L, \mathcal{K})) \lesssim \mathcal{K}^{-\alpha/(d+1)}.$$

1546 For utilizing this conclusion, first notice that

1547 
$$\begin{aligned} & \inf_{f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})} \|f(\mathbf{u}) - f^*(\mathbf{u})\|_2 = \inf_{f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})} \sqrt{\sum_{i=1}^{d^*} \{f_i(\mathbf{u}) - f_i^*(\mathbf{u})\}^2} \\ & \leq \inf_{f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})} \sqrt{\sum_{i=1}^{d^*} \|f_i - f_i^*\|_{C([0,1]^d)}^2} \leq \sup_{g \in \mathcal{H}^\alpha} \inf_{f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})} \sqrt{\sum_{i=1}^{d^*} \|f_i - g\|_{C([0,1]^d)}^2} \\ & \leq \sup_{g \in \mathcal{H}^\alpha} \sqrt{\sum_{i=1}^{d^*} \inf_{f \in \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K})} \|f - g\|_{C([0,1]^d)}^2} \leq \sqrt{d^*} \mathcal{E}(\mathcal{H}^\alpha, \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K})) \lesssim \mathcal{K}^{-\alpha/(d+1)}, \end{aligned}$$

 1548 where the third to last line inequality is from following reason: if  $f_i \in \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K})$ , where  $i \in [d^*]$ , whose parameter are independent with each other, then their concatenation  $f = (f_1, f_2, \dots, f_{d^*})^\top$  can be regarded as an elements of  $\mathcal{NN}_{d,d^*}(W, D, \mathcal{K})$  with specific parameters, by following Proposition A.19, we have  $f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})$ .

 1549 **Proposition A.19** ((iii) of Proposition 2.5 in [Jiao et al. \(2023\)](#)). *Let  $f_1 \in \mathcal{NN}_{d,d_1^*}(w_1, L_1, \mathcal{K}_1)$  and  $f_2 \in \mathcal{NN}_{d,d_2^*}(w_2, L_2, \mathcal{K}_2)$ , define  $f(\mathbf{x}) := (f_1(\mathbf{x}), f_2(\mathbf{x}))$ , then  $f \in \mathcal{NN}_{d,d_1^*+d_2^*}(W_1 + W_2, \max\{L_1, L_2\}, \max\{\mathcal{K}_1, \mathcal{K}_2\})$ .*

 1550 Above conclusion implies optimal approximation element of  $f^*$  in  $\mathcal{NN}_{d,d^*}(W, L, \mathcal{K})$  can be arbitrarily close to  $f^*$  under the setting that  $\mathcal{K}$  is large enough. Hence we can conclude optimal approximation element of  $f^*$  is also contained in  $\mathcal{F} = \mathcal{NN}_{d,d^*}(W, L, \mathcal{K}, B_1, B_2)$  as the setting that  $B_1 \leq \|f^*\|_2 \leq B_2$ .

1551 Therefore, if we denote

1552 
$$\mathcal{T}(f) := \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2^2\} + \lambda \|\mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{f(\mathbf{x}_1)f(\mathbf{x}_2)^\top\} - I_{d^*}\|_F^2,$$

 1553 we can yield the upper bound of  $\mathcal{E}_{\mathcal{F}}$  by following deduction

1554 
$$\begin{aligned} \mathcal{E}_{\mathcal{F}} &= \inf_{f \in \mathcal{F}} \{ \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) - \sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G) \} = \inf_{f \in \mathcal{F}} \{ \mathcal{T}(f) - \mathcal{T}(f^*) \} = \inf_{f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})} \{ \mathcal{T}(f) - \mathcal{T}(f^*) \} \\ &\leq \|\ell\|_{\text{Lip}} \inf_{f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\tilde{\mathbf{x}}} \|\tilde{f}(\tilde{\mathbf{x}}) - \tilde{f}^*(\tilde{\mathbf{x}})\|_2 \leq \|\ell\|_{\text{Lip}} \inf_{f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}' \in \mathcal{A}(\mathbf{x})} \sqrt{2 \sum_{i=1}^{d^*} \{f_i(\mathbf{x}') - f_i^*(\mathbf{x}')\}^2} \\ &\leq \sqrt{2d^*} \|\ell\|_{\text{Lip}} \sup_{g \in \mathcal{H}^\alpha} \inf_{f \in \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K}/\sqrt{d^*})} \|f - g\|_{C([0,1]^d)} \leq \sqrt{2d^*} \|\ell\|_{\text{Lip}} \mathcal{E}(\mathcal{H}^\alpha, \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K}/\sqrt{d^*})) \\ &\lesssim \mathcal{K}^{-\alpha/(d+1)}. \end{aligned}$$

1555 where the first inequality is because of Proposition A.7.

A.3.6. BOUND  $\mathcal{E}_{\mathcal{G}}$ 

Let  $\mathcal{M}(\mathbf{u}) = \mathbf{u}_1 \mathbf{u}_2^\top$ ,  $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^{d^*}$ , which is a Lipchitz map on  $\{\mathbf{u} \in \mathbb{R}^{2d^*} : \mathbf{u} \leq \sqrt{2}B_2\}$ , as presented in Proposition A.7. Then

$$\begin{aligned} \mathbb{E}_{\tilde{D}_s} \{\mathcal{E}_{\mathcal{G}}\} &\lesssim \mathbb{E}_{\tilde{D}_s} \left[ \sup_{f \in \mathcal{F}} \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \left[ \frac{1}{n_s} \sum_{i=1}^{n_s} \{\mathcal{M}(\tilde{f}(\tilde{\mathbf{x}})) - \mathcal{M}(\tilde{f}(\tilde{\mathbf{x}}^{(i)}))\} \right] \right\|_F \right] \\ &\leq \|\mathcal{M}\|_{\text{Lip}} \mathbb{E}_{\tilde{D}_s} \left[ \left\| \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{\tilde{f}(\tilde{\mathbf{x}})\} - \frac{1}{n_s} \sum_{i=1}^{n_s} \tilde{f}(\tilde{\mathbf{x}}^{(i)}) \right\|_2 \right] \end{aligned}$$

Furthermore, according to the multidimensional Chebyshev's inequality, we turn out that  $\mathbb{P}_s \left( \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} \tilde{f}(\tilde{\mathbf{x}}^{(i)}) - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{\tilde{f}(\tilde{\mathbf{x}})\} \right\|_2 \geq \frac{1}{n_s^{1/4}} \right) \leq \frac{\mathbb{E} \|\tilde{f}(\tilde{\mathbf{x}}) - \mathbb{E}\{\tilde{f}(\tilde{\mathbf{x}})\}\|_2^2}{\sqrt{n_s}} \leq \frac{8B_2^2}{\sqrt{n_s}}$  as  $\|\tilde{f}(\tilde{\mathbf{x}})\|_2 \leq \sqrt{2}B_2$ . Thus we have

$$\begin{aligned} \mathbb{E}_{\tilde{D}_s} \{\mathcal{E}_{\mathcal{G}}\} &\lesssim \frac{1}{n_s^{1/4}} \cdot \mathbb{P}_s \left( \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} \tilde{f}(\tilde{\mathbf{x}}^{(i)}) - \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \{\tilde{f}(\tilde{\mathbf{x}})\} \right\|_2 \geq \frac{1}{n_s^{1/4}} \right) + 2\sqrt{2}B_2 \cdot \frac{8B_2^2}{\sqrt{n_s}} \\ &\leq \frac{1}{n_s^{1/4}} + 16\sqrt{2}B_2^3 \frac{1}{\sqrt{n_s}} \lesssim \frac{1}{n_s^{1/4}}. \end{aligned}$$

where the first inequity is due to  $\|\tilde{f}(\tilde{\mathbf{x}})\|_2 \leq \sqrt{2}B_2$ .

## A.3.7. SUBSECTION: TRADE OFF BETWEEN STATISTICAL ERROR AND APPROXIMATION ERROR

Let  $W \geq c\mathcal{K}^{(2d+\alpha)/(2d+2)}$  and  $L \geq 2\lceil \log_2(d+r) \rceil + 2$ , combine the bound results of statistical error and approximation error to yield

$$\mathbb{E}_{\tilde{D}_s} \{\mathcal{E}(\hat{f}_{n_s})\} \lesssim \mathbb{E}_{\tilde{D}_s} [\mathcal{E}_{\text{sta}}] + \mathcal{E}_{\mathcal{F}} + \mathbb{E}_{\tilde{D}_s} \{\mathcal{E}_{\mathcal{G}}\} \lesssim \frac{\mathcal{K}}{\sqrt{n_s}} + \mathcal{K}^{-\alpha/(d+1)}.$$

Taking  $\mathcal{K} = n_s^{\frac{d+1}{2(\alpha+d+1)}}$  to yield  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{E}(\hat{f}_{n_s})\} \lesssim n_s^{-\frac{\alpha}{2(\alpha+d+1)}}$ . As we have shown that  $\mathcal{L}(f^*) = 0$ , above inequality implies  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} \lesssim n_s^{-\frac{\alpha}{2(\alpha+d+1)}}$ . To ensure above deduction holds, We need to set  $W \geq cn_s^{\frac{2d+\alpha}{4(\alpha+d+1)}}$  and  $L \geq 2\lceil \log_2(d+r) \rceil + 2$ .

## A.3.8. THE PROOF OF MAIN THEOREM

Next, we are going to prove our main theorem 3.9. We will state its formal version at first and then conclude Theorem 3.9 as a corollary.

To notation conciseness, let  $p = \frac{\sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} \left( \frac{C}{\lambda} n_s^{-\frac{\alpha}{2(\alpha+d+1)}} + \phi(n_s) \right) + 2\sqrt{d^*} B_2 M n_s^{-\frac{\nu}{2(\alpha+d+1)}}}}{B_2^2 \psi(\sigma_s^{(n_s)}, \delta_s^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s})}$ , where  $C$  is a constant,  $0 \leq \phi(n_s) \lesssim (1 - \sigma_s^{(n_s)} + n_s^{-\frac{\min\{\alpha, \nu, \varsigma, \tau\}}{4(\alpha+d+1)}})^2 + (1 - \sigma_s^{(n_s)} + n_s^{-\frac{\min\{\alpha, \nu, \varsigma, \tau\}}{8(\alpha+d+1)}})$ , then the formal version of our main theoretical result can be stated as follow.

**Lemma A.20.** *When Assumptions 3.5, 3.3, 3.2, 3.7 and 3.8 all hold, set  $\varepsilon_{n_s} = m^2 n_s^{-\frac{\min\{\alpha, \nu, \varsigma, \tau\}}{8(\alpha+d+1)}}$ ,  $W \geq cn_s^{\frac{2d+\alpha}{4(\alpha+d+1)}}$ ,  $L \geq 2\lceil \log_2(d+r) \rceil + 2$ ,  $\mathcal{K} = n_s^{\frac{d+1}{2(\alpha+d+1)}}$  and  $\mathcal{A} = \mathcal{A}_{n_s}$  in Assumption 3.7, then we have*

$$\mathbb{E}_{\tilde{D}_s} [R_t^2(\varepsilon_{n_s}, \hat{f}_{n_s})] \lesssim n_s^{-\frac{\min\{\alpha, \nu, \varsigma\}}{4(\alpha+d+1)}} \quad (36)$$

and

$$\mathbb{E}_{\tilde{D}_s} \left\{ \max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \right\} \lesssim (1 - \sigma_s^{(n_s)}) + n_s^{-\frac{\min\{\alpha, 2\tau\}}{4(\alpha+d+1)}}. \quad (37)$$

Furthermore, If  $\psi(\sigma_s^{(n_s)}, \delta_s^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) > 0$ , then with probability at least  $1 - p$ , we have

$$\mathbb{E}_{\tilde{D}_s} \{\text{Err}(Q_{\hat{f}_{n_s}})\} \leq (1 - \sigma_t^{(n_s)}) + \mathcal{O}(n_s^{-\frac{\min\{\alpha, \nu, \varsigma\}}{8(\alpha+d+1)}}).$$

1650 *Proof.* First recall the conclusion we've got in Lemma A.6

$$1651 \mathbb{E}_{\tilde{D}_s} \{R_t^2(\varepsilon, \hat{f}_{n_s})\} \leq \frac{m^4}{\varepsilon^2} (\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + 8B_2 d^* MK \varepsilon_1 + 4B_2^2 d^* K \varepsilon_2),$$

$$1652 \mathbb{E}_{\tilde{D}_s} [\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)|] \leq \sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} \left( \frac{1}{\lambda} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + \mathbb{E}_{\tilde{D}_s} \{\phi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})\} \right) + 2\sqrt{d^*} B_2 MK \varepsilon_1},$$

1653 and with probability at least  $1 - \frac{\sqrt{\frac{2}{\min_{i \neq j} p_s(i)p_s(j)} \left( \frac{1}{\lambda} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + \phi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}) \right) + 2\sqrt{d^*} B_2 MK \varepsilon_1}}{B_2^2 \psi(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})}$ , we have

$$1654 \mathbb{E}_{\tilde{D}_s} \{\text{Err}(Q_{\hat{f}_{n_s}})\} \leq (1 - \sigma_t) + \frac{m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} + 8B_2 d^* MK \varepsilon_1 + 4B_2^2 d^* K \varepsilon_2},$$

1655 where  $\phi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}) = 4B_2^2 \left[ \left(1 - \sigma_s + \frac{\mathcal{K} \delta_s + 2\varepsilon}{2B_2}\right)^2 + (1 - \sigma_s) + \frac{\mathcal{K} m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}} \left(3 - 2\sigma_s + \frac{\mathcal{K} \delta_s + 2\varepsilon}{B_2}\right) + \frac{m^4}{\varepsilon^2} \mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} \left( \sum_{k=1}^K \frac{1}{p_s(k)} \right) \right] + B_2 \left( \varepsilon^2 + \frac{4B_2^2 m^2}{\varepsilon} \sqrt{\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\}} \right)^{\frac{1}{2}}$ .

1656 To obtain the conclusion shown in this theorem from above formulations, first we plug  $\varepsilon_1 \leq n_s^{-\frac{\nu+d+1}{2(\alpha+d+1)}}$  and  $\varepsilon_2 \leq n_s^{-\frac{\tau+d+1}{2(\alpha+d+1)}}$  into it. apart from that, we have shown  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} \lesssim n_s^{-\frac{\alpha}{2(\alpha+d+1)}}$  in A.3.7 and known  $\delta_s^{(n_s)} \leq n_s^{-\frac{\tau+d+1}{2(\alpha+d+1)}}$ , combining with the setting  $\varepsilon_{n_s} = m^2 n_s^{-\frac{\min\{\alpha, \nu, \varsigma, \tau\}}{8(\alpha+d+1)}}$ ,  $\mathcal{K} = n_s^{\frac{d+1}{2(\alpha+d+1)}}$  implies that  $\mathcal{K} \varepsilon_1 / \varepsilon_{n_s}^2 \leq n_s^{-\frac{\tau}{2(\alpha+d+1)}}$ ,  $\varepsilon_2 / \varepsilon_{n_s}^2 \leq n_s^{-\frac{\tau}{2(\alpha+d+1)}}$ ,  $\mathcal{K} \delta_s^{(n_s)} \leq n_s^{-\frac{\tau}{2(\alpha+d+1)}}$  and  $\mathbb{E}_{\tilde{D}_s} \{\mathcal{L}(\hat{f}_{n_s})\} / \varepsilon_{n_s}^2 \leq n_s^{-\frac{\alpha}{4(\alpha+d+1)}}$ .

1657 Plugin these facts into the corresponding term of above formulations to get what we desired.  $\square$

1658 Let us first state the formal version of Theorem 3.9 and then prove it.

1659 **Theorem A.21** (Formal version of Theorem 3.9). *If Assumptions 3.5, 3.3, 3.2, 3.7 and 3.8 all hold, set  $W \geq cn_s^{\frac{2d+\alpha}{4(\alpha+d+1)}}$ ,  $L \geq 2[\log_2(d+r)] + 2$ ,  $\mathcal{K} = n_s^{\frac{d+1}{2(\alpha+d+1)}}$  and  $\mathcal{A} = \mathcal{A}_{n_s}$  in Assumption 3.7, then, provided that  $n_s$  is sufficiently large, with probability at least  $\sigma_s^{(n_s)} - \mathcal{O}(n_s^{-\frac{\min\{\alpha, \nu, \varsigma, \tau\}}{16(\alpha+d+1)}}) - \mathcal{O}(\frac{1}{\sqrt{\min_k n_t(k)}})$ , we have*

$$1660 \mathbb{E}_{\tilde{D}_s} \{\text{Err}(Q_{\hat{f}_{n_s}})\} \leq (1 - \sigma_t^{(n_s)}) + \mathcal{O}(n_s^{-\frac{\min\{\alpha, \nu, \varsigma\}}{8(\alpha+d+1)}}).$$

1661 *Proof of Theorem 3.9.* Note that the main difference between Theorem A.20 and Theorem 3.9 is the condition  $\psi(\sigma_s^{(n_s)}, \delta_s^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) > 0$ , so we are going to focus on whether this condition holds under the condition of Theorem 3.9.

1662 To show this, first recall  $\psi(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) = \Gamma_{\min}(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) - \sqrt{2 - 2\Gamma_{\min}(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s})} - \frac{\Delta_{\hat{\mu}_t}}{2} - \frac{2 \max_{k \in [K]} \|\hat{\mu}_t(k) - \mu_t(k)\|_2}{B_2}$ . Note (28) and dominated convergence theorem imply  $R_t(\varepsilon_{n_s}, \hat{f}_{n_s}) \rightarrow 0$  a.s., thus

$$1663 \Gamma_{\min}(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) = \left( \sigma_t^{(n_s)} - \frac{R_t(\varepsilon_{n_s}, \hat{f}_{n_s})}{\min_i p_t(i)} \right) \left( 1 + \left( \frac{B_1}{B_2} \right)^2 - \frac{\mathcal{K} \delta_t^{(n_s)}}{B_2} - \frac{2\varepsilon_{n_s}}{B_2} \right) - 1 \rightarrow \left( \frac{B_1}{B_2} \right)^2$$

1664 Combining with the fact that  $\frac{\Delta_{\hat{\mu}_t}}{2} = \frac{1 - \min_{k \in [K]} \|\hat{\mu}_t(k)\|_2^2 / B_2^2}{2} < \frac{1}{2}$  can yield

$$1665 \Gamma_{\min}(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) - \sqrt{2 - 2\Gamma_{\min}(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s})} - \frac{\Delta_{\hat{\mu}_t}}{2} > 1/2$$

1666 if we select proper  $B_1$  and  $B_2$ .

1705 Besides that, by Multidimensional Chebyshev's inequality, we know that

$$1706 \mathbb{P}_t(\|\hat{\mu}_t(k) - \mu_t(k)\|_2 \geq \frac{B_2}{8}) \leq \frac{64 \sqrt{\mathbb{E}_{\mathbf{z} \in \tilde{C}_t(k)} \mathbb{E}_{\mathbf{z}' \in \mathcal{A}(\mathbf{z})} \|f(\mathbf{z}') - \mu_t(k)\|_2^2}}{B_2^2 \sqrt{2n_t(k)}} \leq \frac{128}{B_2 \sqrt{n_t(k)}},$$

1710 so that  $\psi(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) \geq \frac{1}{4}$  with probability at least  $1 - \frac{128K}{B_2 \sqrt{\min_k n_t(k)}}$  if  $n_s$  is large enough, of course the  
 1712 condition  $\psi(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) > 0$  in Theorem A.20 can be satisfied.

1714 Therefore, with probability at least

$$1715 1 - p - \frac{128K}{B_2 \sqrt{\min_k n_t(k)}} \gtrsim 1 - (1 - \sigma_s^{(n_s)}) - \mathcal{O}\left(n_s^{-\frac{\min\{\alpha, \nu, \varsigma, \tau\}}{16(\alpha+d+1)}}\right) - \mathcal{O}\left(\frac{1}{\sqrt{\min_k n_t(k)}}\right)$$

$$1716 = \sigma_s^{(n_s)} - \mathcal{O}\left(n_s^{-\frac{\min\{\alpha, \nu, \varsigma, \tau\}}{16(\alpha+d+1)}}\right) - \mathcal{O}\left(\frac{1}{\sqrt{\min_k n_t(k)}}\right).$$

1721 we have the conclusions shown in Theorem 3.9, which completes the proof. □