SELF-SUPERVISED TRANSFER LEARNING VIA ADVER SARIAL CONTRASTIVE TRAINING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023 024

025

Paper under double-blind review

ABSTRACT

Learning a data representation with strong transferability from an unlabeled scenario is both crucial and challenging. In this paper, we propose a novel self-supervised transfer learning approach via Adversarial Contrastive Training (ACT). Additionally, we establish an end-to-end theoretical understanding for self-supervised contrastive pretraining and its implications for downstream classification tasks in a misspecified, over-parameterized setting. Our theoretical findings highlight the provable advantages of adversarial contrastive training in the source domain towards improving the accuracy of downstream tasks in the target domain. Furthermore, we illustrate that downstream tasks necessitate only a minimal sample size when working with a well-trained representation, offering valuable insights on few-shot learning. Last but not least, extensive experiments across various datasets demonstrate a significant enhancement in classification accuracy when compared to existing state-of-the-art self-supervised learning methods.

1 INTRODUCTION

Collecting unlabeled data is far more convenient and cost-effective than gathering labeled data in real-world applications. As a result, learning representations from abundant unlabeled data has become a critical and foundational challenge. Pretraining on unlabeled data enables the capture of more general, abstract features without the need for specific labels. Consequently, the learned task-invariant representations demonstrate superior transferability to unseen data, making them highly effective in transfer learning scenarios.

- 032 One of the most popular approaches to learning representations from unlabeled data is self-033 supervised contrastive learning, which has garnered significant attention due to its impressive perfor-034 mance. The rationale behind contrastive learning involves acquiring a representation that maintains augmentation invariance while preventing model collapse. The latter aspect is crucial, as solely bringing positive pairs closer could result in trivial solutions. The initial body of work heavily re-037 lies on the utilization of negative samples, such as Ye et al. (2019); He et al. (2020); Chen et al. 038 (2020a;b); HaoChen et al. (2021); Zhang et al. (2023). These studies prevent representation collapse by pushing negative pairs apart in the feature space. However, the construction of negative pairs poses significant challenges. Firstly, augmented views from distinct data points sharing the 040 same semantic meaning may inadvertently be treated as negative pairs, impeding semantic extrac-041 tion. Secondly, the quality of the representation is highly dependent on the number of negative pairs, 042 necessitating substantial computational and memory resources. 043
- In recent years, there has been a surge of interests in developing self-supervised learning methods
 that eschew the use of negative samples (Grill et al., 2020; Caron et al., 2020; 2021; Ermolov et al.,
 2021; Zbontar et al., 2021; Chen & He, 2021; Bardes et al., 2022; Ozsoy et al., 2022; HaoChen
 et al., 2022; Wang et al., 2024). Among above mentioned studies, the most prominent works include Zbontar et al. (2021); Bardes et al. (2022); Ozsoy et al. (2022); Zhang
 et al. (2023), which prevent the model collapse by incorporating a regularization term into the loss
 function.
- In this study, we introduce a novel self-supervised learning approach called <u>A</u>dversarial <u>C</u>ontrastive <u>Training (ACT)</u>, designed to learn representations without the need for constructing negative samples, while avoiding the bias between population loss and sample-level loss. The core idea of ACT is the introduction of a novel regularization term that encourages the separation of category centers

within the latent space, thereby improving classification accuracy in downstream tasks. Moreover, this regularization term incorporates an unbiased sample version, enabling rigorous theoretical analysis. Let

$$\mathcal{R}(f,G) = \left\langle \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} \left[f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top \right] - I_{d^*}, G \right\rangle_F$$

where $f : \mathbb{R}^d \to \mathbb{R}^{d^*}$ is a representation function, G is a matrix in $\mathbb{R}^{d^* \times d^*}$, and the Frobenius inner product is defined as $\langle A, B \rangle_F := \operatorname{tr}(A^\top B)$ for any $A, B \in \mathbb{R}^{d_1 \times d_2}$. Then we learn the contrastive representation through a minimax optimization problem

$$\min_{f} \max_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) = \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} \left[\|f(\boldsymbol{x}_1) - f(\boldsymbol{x}_2)\|_2^2 \right] + \lambda \mathcal{R}(f, G),$$
(1)

where the first term in (1) facilitates achieving augmentation invariance in the representation, similar with the previous works (Zbontar et al., 2021; Bardes et al., 2022; HaoChen et al., 2022). Here $\mathcal{A}(x)$ denotes the set of augmentations of a sample x, and $\lambda > 0$ is the regularization parameter. Thanks to the minimax formulation in (1), we propose the following loss of our ACT at the sample level

$$\widehat{\mathcal{L}}(f,G) := \frac{1}{n_s} \sum_{i=1}^{n_s} \left[\|f(\boldsymbol{x}_1^{(i)}) - f(\boldsymbol{x}_2^{(i)})\|_2^2 + \lambda \langle f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - I_{d^*}, G \rangle_F \right],$$
(2)

where $\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(n_s)}$ are unlabeled data, $\boldsymbol{x}_1^{(i)}$ and $\boldsymbol{x}_2^{(i)}$ are independent augmentations of $\boldsymbol{x}^{(i)}$. It is crucial that (2) is unbiased since $\mathbb{E}_{D_s}[\widehat{\mathcal{L}}(f,G)] = \mathcal{L}(f,G)$ for each fixed $G \in \mathbb{R}^{d^* \times d^*}$.

In fact, the inner maximization problem has a explicit solution that $G = \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*}$, therefore (1) is equivalent to minimizing following loss

$$\mathcal{L}(f) := \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} \left[\|f(\boldsymbol{x}_1) - f(\boldsymbol{x}_2)\|_2^2 \right] + \lambda \left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*} \right\|_F^2.$$
(3)

Directly discretizing the expectation in (3) yields a biased sample-level loss as

$$\widehat{\mathcal{L}}(f) := \frac{1}{n_s} \sum_{i=1}^{n_s} \|f(\boldsymbol{x}_1^{(i)}) - f(\boldsymbol{x}_2^{(i)})\|_2^2 + \lambda \left\|\frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top\right] - I_{d^*} \right\|_F^2.$$

Specifically, we have $\mathbb{E}_{D_s}[\widehat{\mathcal{L}}(f)] \neq \mathcal{L}(f)$ due to the non-commutativity between the expectation 084 085 and the Frobenius norm, where D_s represents the dataset used for pretraining. While this biased discretization method has been employed in previous studies (HaoChen et al., 2022; HaoChen & Ma, 2023), its application presents a significant challenge in terms of theoretical analysis. For 087 instance, despite that Huang et al. (2023) establish a theoretical analysis for Zbontar et al. (2021) at 880 the population-level, the extensions of these findings to the sample-level is not straightforward due to 089 the bias of the estimation. HaoChen & Ma (2023) establish a theoretical understanding for HaoChen 090 et al. (2022) at the sample-level, nonetheless, the results are subject to strong assumptions given the 091 biased nature of the estimation. 092

From a theoretical perspective, we establish a rigorous end-to-end theoretical analysis for both the contrastive pre-training and the downstream classification under mild conditions. Further, our findings demonstrate the provable advantages of self-supervised contrastive pre-training and provides theoretical insights into determining the sample size and selecting the appropriate scale for deep neural networks. Our experiment yields remarkable classification accuracy when employing both fine-tuned linear probes and the K-nearest neighbor (K-NN) protocol across a range of benchmark datasets. These results showcase a high level of competitiveness with current state-of-the-art selfsupervised learning methodologies, as illustrated in Table 1.

101

102

058

063

069 070 071

072

073 074

075

076

077 078 079

1.1 RELATED WORK

Self-supervised transfer learning Thanks to the robust transferability inherent in representations learned by self-supervised learning, the field of few-shot learning, which aims to train models with only a limited number of labeled samples, has significantly advanced through self-supervised methodologies. This progression is evidenced by the contributions of Liu et al. (2021); Rizve et al. (2021); Yang et al. (2022); Lim et al. (2023). However, current work only demonstrates the effective-ness of self-supervised learning for few-shot learning mainly empirically. Theoretical explanations

remain scarce. Understanding how the learned representations from unlabeled data enhance prediction performance with only a few labeled samples in downstream tasks is a critical question that requires further investigation. Especially investigating the impact of upstream samples on downstream samples. Therefore, a thorough theoretical analysis at sample level is urgently needed.

Although Saunshi et al. (2019); HaoChen et al. (2021); Garrido et al. (2022); Awasthi et al. (2022);
Ash et al. (2022); HaoChen et al. (2022); Lei et al. (2023); Huang et al. (2023) have offered some theoretical progresses in understanding self-supervised learning, all these studies either remain at the population level, or focus solely on the generalization property of hypothesis space with a finite complexity measure. The effects of both upstream and downstream sample sizes are still unknown.

117 HaoChen & Ma (2023) use augmented graphs to provide a more thorough theoretical analysis at 118 sample level for the self-supervised learning loss proposed in HaoChen et al. (2022). They establish 119 a theoretical guarantees at the sample level, under certain strong assumptions, including Assump-120 tions 4.2 and 4.4. Assumption 4.2 assumes the existence of a neural network capable of sufficiently 121 minimize the loss. In contrast, we demonstrate the existence of a measurable function that can 122 vanish our loss by accounting for additional approximation error. This necessitates an extension 123 of the well-specified setting to a misspecified setting. Moreover, the most important problem in 124 self-supervised transfer learning theory pertains to elucidating the mechanism through which the 125 representation acquired from the upstream task facilitates the learning process of the downstream task. While HaoChen & Ma (2023) assume this relationship as Assumption 4.4 in their research, 126 our study surpasses the current body of literature by conducting a comprehensive investigation into 127 the impact of approximation error and generalization error during the pre-training phase on down-128 stream test error. This analysis sheds light on how the size of the upstream sample influences the 129 downstream task, particularly in scenarios where the availability of downstream samples is con-130 strained.

131

138

140

142

143

144

145

146

147

148

149

Comparison with existing contrastive learning algorithms HaoChen et al. (2022) can be regarded as a special version of our model with the constraint $x_1 = x_2$ at the population level. However, its loss at the sample level adopts a biased discretization method, which leads to a different optimization direction compared to ACT, especially in the mini-batch scenario. More discussion can be found in Remark 2.1. Besides that, the loss at the sample level provided by Zbontar et al. (2021) is also similar to our loss, but its unbiased counterpart at the population level is still unknown.

139 1.2 CONTRIBUTIONS

141 Our main contributions can be summarized as follows.

- We introduce a novel self-supervised transfer learning method called <u>A</u>dversarial <u>C</u>ontrastive <u>T</u>raining (ACT). This approach learns representations from unlabeled data by tackling a minimax optimization problem, which aims to de-bias the initially proposed risk, thereby providing a foundation for establishing a thorough theoretical understanding.
- Our experimental results demonstrate outstanding classification accuracy using both finetuned linear probe and *K*-nearest neighbor (*K*-NN) protocol on various benchmark datasets, showing competitiveness with existing state-of-the-art self-supervised learning methods.
- In the context of transfer learning, we present a thorough theoretical understanding for both ACT and its downstream classification tasks within a misspecified and overparameterized scenario. Our theoretical results offer insights into determining the samples size for pre-training and appropriate depth, width, and norm restrictions of neural networks. These findings illuminate the advantages of ACT in enhancing the accuracy of downstream tasks. Furthermore, we demonstrate that leveraging the representations learned by ACT in the source domain enables high accuracy in the downstream tasks of the target domain, even when only a small amount of data is available.
- 158
- 159 1.3 ORGANIZATIONS
- 161 The remainder of this paper is organized as follows. In Section 2, we introduce basic notations and presents the adversarial self-supervised learning loss, along with an alternating optimization algo-

rithm to address the minimax problem. Section 3 showcases experimental results for representations
 learned by ACT across various real datasets and evaluation protocols. Section 4 provides an end-to end theoretical guarantee for ACT. Conclusions are discussed in Section 5, respectively. Detailed
 proofs and experimental details are differed to Section B and C respectively.

166 167

168 169

170

171

2 ADVERSARIAL CONTRASTIVE TRAINING

In this section, we provide a novel method for unsupervised transfer learning via adversarial contrastive training (ACT). We begin with some notations in Section 2.1. Then, we introduce ACT method and alternating optimization algorithm in Section 2.2. Finally, we outline the setup of the downstream task in Section 2.3.

2.1 PRELIMINARIES AND NOTATIONS

176 Denote by $\|\cdot\|_2$ and $\|\cdot\|_\infty$ the 2-norm and ∞ -norm of the vector, respectively. Let $A, B \in \mathbb{R}^{d_1 \times d_2}$ 177 be two matrices. Define the Frobenius inner product $\langle A, B \rangle_F = \operatorname{tr}(A^\top B)$. Denote by $\|\cdot\|_F$ 178 the Frobenius norm induced by Frobenius inner product. We denote the ∞ -norm of the matrix as 179 $\|A\|_\infty := \sup_{\|x\|_\infty \leq 1} \|Ax\|_\infty$, which is the maximum 1-norm of the rows of A. The Lipschitz 180 norm of a map f from \mathbb{R}^{d_1} to \mathbb{R}^{d_2} is defined as $\|f\|_{\operatorname{Lip}} := \sup_{x \neq y} \frac{\|f(x) - f(y)\|_2}{\|x - y\|_2}$.

Let $L, N_1, \ldots, N_L \in \mathbb{N}, 0 < B_1 \leq B_2$. A deep ReLU neural network hypothesis space is defined as

$$\mathcal{NN}_{d_1,d_2}(W,L,\mathcal{K},B_1,B_2) := \begin{cases} \phi_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{A}_L \sigma(\boldsymbol{A}_{L-1}\sigma(\cdots\sigma(\boldsymbol{A}_0\boldsymbol{x}+\boldsymbol{b}_0)) + \boldsymbol{b}_{L-1}), \\ W = \max\{N_1,\dots,N_L\}, \ \kappa(\boldsymbol{\theta}) \leq \mathcal{K}, \ B_1 \leq \|\phi_{\boldsymbol{\theta}}\|_2 \leq B_2 \end{cases},$$

 $\kappa(\boldsymbol{ heta}) := \|\boldsymbol{A}_L\|_{\infty} \prod_{l=0}^{L-1} \max\{\|(\boldsymbol{A}_l, \boldsymbol{b}_l)\|_{\infty}, 1\}.$

187 where $\sigma(x) := x \vee 0$ is the ReLU activate function, $N_0 = d_1$, $N_{L+1} = d_2$, $A_i \in \mathbb{R}^{N_{i+1} \times N_i}$ 188 and $b_i \in \mathbb{R}^{N_{i+1}}$. The integers W and L are called the width and depth of the neural network, 189 respectively. $B_1 \leq \|\phi_{\theta}\|_2 \leq B_2$ is used to indicate any $u \in [0,1]^d$, $B_1 \leq \|\phi_{\theta}(u)\|_2 \leq B_2$. 190 The parameters set of the neural network is defined as $\theta := ((A_0, b_0), \dots, (A_{L-1}, b_{L-1}), A_L)$. 191 Further, $\kappa(\theta)$ is defined as

192

185 186

193

194

195 196

197

198

Appendix **B**.1 shows that $\|\phi_{\theta}\|_{\text{Lip}} \leq \mathcal{K}$ for each $\phi_{\theta} \in \mathcal{NN}_{d_1,d_2}(W, L, \mathcal{K}, B_1, B_2)$.

2.2 ADVERSARIAL CONTRASTVE TRAINING

Learning representations from large amounts of unlabeled data has recently gained significant attention, as highly transferable representations offer substantial benefits for downstream tasks. Adversarial contrastve training is driven by two key factors: augmentation invariance and a regularization term to prevent model collapse. Specifically, augmentation invariance aims to make representations of different augmented views of the same sample as similar as possible. However, a trivial representation that maps all augmented views to the same point is ineffective for downstream tasks, making the regularization term essential.

Data augmentation $A : \mathbb{R}^d \to \mathbb{R}^d$ is essentially a transformation of the original sample before training. A commonly-used augmentation is the composition of random transformations, such as RandomCrop, HorizontalFlip, and Color distortion (Chen et al., 2020a). Denote by $\mathcal{A} = \{A_{\gamma}(\cdot) :$ $\gamma \in [m]\}$ the collection of data augmentations, and denote the source domain as $\mathcal{X}_s \subseteq [0, 1]^d$, with its corresponding unknown distribution denoted by P_s . Let $\{x^{(1)}, \ldots, x^{(n_s)}\}$ be n_s i.i.d. unlabeled samples from the source distribution. For each sample $x^{(i)}$, we define the corresponding augmented pair as

$$\tilde{\boldsymbol{x}}^{(i)} = (\boldsymbol{x}_1^{(i)}, \boldsymbol{x}_2^{(i)}) = (A(\boldsymbol{x}^{(i)}), A'(\boldsymbol{x}^{(i)})),$$
(4)

214 215 where A and A' are drawn from the uniform distribution on \mathcal{A} independently. Further, the augmented dataset for ACT is defined as $D_s := {\tilde{x}^{(i)}}_{i \in [n_s]}$. 216 The ACT method can be formulated as a minimax problem 217

$$\hat{f}_{n_s} \in \underset{f \in \mathcal{F}}{\operatorname{arg\,min}} \sup_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{L}}(f, G), \tag{5}$$

where the empirical risk is defined as

218

219 220

221 222

224

225

231

235

240 241 242

247

248 249

250 251

253

254

256 257

$$\widehat{\mathcal{L}}(f,G) := \frac{1}{n_s} \sum_{i=1}^{n_s} \Big[\|f(\boldsymbol{x}_1^{(i)}) - f(\boldsymbol{x}_2^{(i)})\|_2^2 + \lambda \big\langle f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - I_{d^*}, G \big\rangle_F \Big], \tag{6}$$

and $\lambda > 0$ is the regularization parameter, the hypothesis space \mathcal{F} is chosen as the neural network class $\mathcal{NN}_{d,d^*}(W, L, \mathcal{K}, B_1, B_2)$, and the feasible set $\mathcal{G}(f)$ is defined as

$$\widehat{\mathcal{G}}(f) := \left\{ G \in \mathbb{R}^{d^* \times d^*} : \|G\|_F \le \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - I_{d^*} \right\|_F \right\}.$$

The first term of (6) helps the representation to achieve the augmentation invariance while the second term is used to prevent model collapse. It is worth noting that, unlike existing contrastive learning 232 methods (Ye et al., 2019; He et al., 2020; Chen et al., 2020a;; HaoChen et al., 2021), the loss func-233 tion of ACT (6) does not need to construct negative pairs for preventing model collapse, avoiding 234 the issues introduced by negative samples.

We now introduce an alternating algorithm for solving the minimax problem (5). We take the t-th 236 iteration as an example. Observe that the inner maximization problem is linear. Given the previous 237 representation mapping $f_{(t-1)} : \mathbb{R}^d \to \mathbb{R}^{d^*}$, the explicit solution to the maximization problem is 238 given as 239

$$\widehat{G}_{(t)} = \frac{1}{n_s} \sum_{i=1}^{n_s} f_{(t-1)}(\boldsymbol{x}_1^{(i)}) f_{(t-1)}(\boldsymbol{x}_2^{(i)})^\top - I_{d^*}.$$
(7)

Then it suffices to solve the outer minimization problem

$$\hat{f}_{(t)} \in \underset{f \in \mathcal{F}}{\operatorname{arg\,min}} \frac{1}{n_s} \sum_{i=1}^{n_s} \Big[\|f(\boldsymbol{x}_1^{(i)}) - f(\boldsymbol{x}_2^{(i)})\|_2^2 + \lambda \big\langle f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - I_{d^*}, \widehat{G}_{(t)} \big\rangle_F \Big]. \tag{8}$$

Solving the inner problem (7) and the outer problem (8) alternatively yields the desired representation mapping. The detailed algorithm is summarized as Algorithm 1.

Algorithm 1 Adversarial contrastive training (ACT)

Require: Augmented dataset $D_s = {\tilde{x}^{(i)}}_{i \in [n]}$, initial representation $\hat{f}_{(0)}$, iteration horizon T. 1: for $t \in [T]$ do

Update G by solving the inner problem (7). 2:

Update the representation by solving the outer problem (8). 3:

4: end for

5: **return** The learned representation mapping $\hat{f}_{(T)}$.

Remark 2.1. We note that $\hat{G}_{(t)}$ will be detached from the computational graph when solving the outer problem (8) in practice, which means that the gradient of the second term in (8) should be written as $\hat{G}_{(t)}$ will be detached from the computational graph when solving the outer problem (8) in practice, which means that the gradient of the second term in (8) should be written as $\hat{G}_{(t)}$ will be detached from the computational graph when solving the outer problem (8) in practice, which means that the gradient of the second term in (8) should be written as $\hat{G}_{(t)}$ will be detached from the computational graph when solving the outer problem (8) in practice, which means that the gradient of the second term in (8) should be written as $\hat{G}_{(t)}$ will be detached from the computational graph when solving the outer problem (8) in practice, which means that the gradient of the second term in (8) should be written as $\hat{G}_{(t)}$. 258 259 ten as $\langle \nabla_{\boldsymbol{\theta}} \frac{1}{n_s} \sum_{i=1}^{n_s} f_{\boldsymbol{\theta}}(\boldsymbol{x}_1^{(i)}) f_{\boldsymbol{\theta}}(\boldsymbol{x}_2^{(i)})^\top - I_{d^*}, \widehat{G}_{(t)} \rangle$ instead of $\nabla_{\boldsymbol{\theta}} \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f_{\boldsymbol{\theta}}(\boldsymbol{x}_1^{(i)}) f_{\boldsymbol{\theta}}(\boldsymbol{x}_2^{(i)})^\top - I_{d^*} \right\|_F^2$, 260 261 which is a biased discretization of $\left\|\mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1,\boldsymbol{x}_2\in\mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}_1)f(\boldsymbol{x}_2)^{\top}] - I_{d^*}\right\|_F^2$ 262 263

264 2.3 DOWNSTREAM TASK 265

With the help of the representations learned by ACT, we address the downstream classification task 266 in the target domain. Let $\mathcal{X}_t \subseteq [0,1]^d$ represent the target domain, and let P_t be the corresponding 267 unknown distribution. Suppose we have n_t i.i.d. labeled samples $\{(\boldsymbol{z}^{(1)}, y_1), \dots, (\boldsymbol{z}^{(n_t)}, y_{n_t})\} \subseteq$ 268 $\mathcal{X}_t \times [K]$ for the downstream task. We will say that $z \in C_t(k)$ if its label is $k \in [K]$. By a similar 269 process as in obtaining (4), we can construct the augmented dataset in the target domain as follows.

$$D_t = \{ (\tilde{\boldsymbol{z}}^{(i)}, y_i) : \tilde{\boldsymbol{z}}^{(i)} = (\boldsymbol{z}_1^{(i)}, \boldsymbol{z}_2^{(i)}) \}_{i \in [n_t]}, \quad \boldsymbol{z}_1^{(i)} = A(\boldsymbol{z}^{(i)}), \ \boldsymbol{z}_2^{(i)} = A'(\boldsymbol{z}^{(i)})$$

where A and A' are drawn from the uniform distribution on A independently.

Given the representation \hat{f}_{n_e} learned by our self-supervised learning method (5), we adopt following linear probe as the classifier for downstream task:

$$Q_{\hat{f}_{n_s}}(\boldsymbol{z}) = \operatorname*{arg\,max}_{k \in [K]} \left(\widehat{W} \widehat{f}_{n_s}(\boldsymbol{z}) \right)_k,\tag{9}$$

where the k-th row of \widehat{W} is given as

$$\widehat{\mu}_t(k) = \frac{1}{2n_t(k)} \sum_{i=1}^{n_t} (\widehat{f}_{n_s}(\boldsymbol{z}_1^{(i)}) + \widehat{f}_{n_s}(\boldsymbol{z}_2^{(i)})) \mathbb{1}\{y_i = k\}, \quad n_t(k) := \sum_{i=1}^{n_t} \mathbb{1}\{y_i = k\}.$$

This means that we build a template for each class of downstream task through calculating the average representations of each class. Whenever a new sample needs to be classified, simply classify it into the category of the template that it most closely resembles. The algorithm for downstream task can be summarized as Algorithm 2. Finally, the misclassification rate is defined as

$$\operatorname{Err}(Q_{\hat{f}_{n_s}}) = \sum_{k=1}^{K} P_t \big(Q_{\hat{f}_{n_s}}(\boldsymbol{z}) \neq k, \boldsymbol{z} \in C_t(k) \big),$$
(10)

which are used to evaluate the performance of the representation learned by ACT.

Algorithm 2 Downstream classification

Require: Representation mapping \hat{f}_{n_s} , augmented dataset in the target domain D_t = $\{(\tilde{z}^{(i)}, y_i)\}_{i \in [n_t]}$, testing data z.

1: Fit the linear probe according to

$$\widehat{W}(k,:) = \frac{1}{2n_t(k)} \sum_{i=1}^{n_t} (\widehat{f}_{n_s}(\boldsymbol{z}_1^{(i)}) + \widehat{f}_{n_s}(\boldsymbol{z}_2^{(i)})) \mathbb{1}\{y_i = k\}$$

2: Predict the label of testing data by (9).

3: return The predicted label of testing data $Q_{\hat{f}_{n-1}}(z)$.

REAL DATA ANALYSIS

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 1 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.

The experimental details are deferred to Appendix C. The PyTorch code be found in https://anonymous.4open.science/r/Adversarial-Contrastive-Training.

THEORETICAL ANALYSIS

In this section, we will explore an end-to-end theoretical guarantee for ACT. It is crucial to introduce several assumptions while expounding on their rationale in Section 4.1. The main theorem and its proof sketch are presented in Section 4.2. The formal version of the main theorem and further details of the proof can be found in Appendix B.2.

We first define the population ACT risk minimizer as

$$f^* \in \underset{f:B_1 \le \|f\|_2 \le B_2}{\operatorname{arg\,min}} \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f,G),$$
(11)

229	٩,	5	ŋ),	1
205	5	2	-		7
205					
	1	2	r,) [

327 328

347 348 349

350 351

352 353

354

355

356

Table 1: Classification accuracy (top 1) of a linear classifier and a 5-nearest neighbors classifier for different loss functions and datasets. While the results for Barlow Twins are from Bandara et al. (2023), the remains are derived from Ermolov et al. (2021).

Method	CIFAR-10		CIFAR-100		Tiny ImageNet	
	linear	5-NN	linear	5-NN	linear	5-NN
SimCLR (Ermolov et al. (2021))	91.80	88.42	66.83	56.56	48.84	32.86
BYOL (Ermolov et al. (2021))	91.73	89.45	66.60	56.82	51.00	36.24
W-MSE 2 (Ermolov et al. (2021))	91.55	89.69	66.10	56.69	48.20	34.16
W-MSE 4 (Ermolov et al. (2021))	91.99	89.87	67.64	56.45	49.20	35.44
BarlowTwins (Bandara et al. (2023))	87.76	86.66	61.64	55.94	41.80	33.60
VICReg (our repro.)	86.76	83.70	57.13	44.63	40.04	30.46
HaoChen et al. (2022) (our repro.)	86.53	84.20	59.68	49.26	35.80	20.36
ACT (our repro.)	92.11	90.01	68.24	58.35	49.72	36.40

where $\mathcal{L}(\cdot, \cdot)$, the unbiased population counterpart of $\widehat{\mathcal{L}}(\cdot, \cdot)$ (6), is defined as

$$\mathcal{L}(f,G) = \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} \left[\|f(\boldsymbol{x}_1) - f(\boldsymbol{x}_2)\|_2^2 \right] + \lambda \left\langle \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} \left[f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top \right] - I_{d^*}, G \right\rangle_F,$$

and the population feasible set is defined as

$$\mathcal{G}(f) = \left\{ G \in \mathbb{R}^{d^* \times d^*} : \|G\|_F \le \|\mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}_1)f(\boldsymbol{x}_2)^\top] - I_{d^*}\|_F \right\}$$

Here B_1 and B_2 are two positive constant, and we will detail how to set B_1 and B_2 later.

4.1 ASSUMPTIONS

In this subsection, we will put forward certain assumptions that are necessary to establish our main theorem. We first assume that each component of f^* exhibits a certain regularity and smoothness.

Definition 4.1 (Hölder class). Let $d \in \mathbb{N}$ and $\alpha = r + \beta > 0$, where $r \in \mathbb{N}_0$ and $\beta \in (0, 1]$. We denote the Hölder class $\mathcal{H}^{\alpha}(\mathbb{R}^d)$ as

$$\mathcal{H}^{\alpha}(\mathbb{R}^{d}) := \Big\{ f: \mathbb{R}^{d} \to \mathbb{R}, \max_{\|\boldsymbol{s}\|_{1} \leq r} \sup_{\boldsymbol{x} \in \mathbb{R}^{d}} |\partial^{\boldsymbol{s}} f(\boldsymbol{x})| \leq 1, \max_{\|\boldsymbol{s}\|_{1} = r} \sup_{\boldsymbol{x} \neq \boldsymbol{y}} \frac{\partial^{\boldsymbol{s}} f(\boldsymbol{x}) - \partial^{\boldsymbol{s}} f(\boldsymbol{y})}{\|\boldsymbol{x} - \boldsymbol{y}\|_{\infty}^{\beta}} \leq 1 \Big\},$$

where the multi-index $s \in \mathbb{N}_0^d$. Furthermore, we denote $\mathcal{H}^{\alpha} := \{f : [0,1]^d \to \mathbb{R}, f \in \mathcal{H}^{\alpha}(\mathbb{R}^d)\}$ 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.

Assumption 4.1. 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 4.1 is a standard assumption in nonparametric statistics (Tsybakov, 2008; Schmidt-Hieber, 2020), more specifically in studies of neural network approximation capacity (Yarotsky, 2018; Yarotsky & Zhevnerchuk, 2020). It is a pretty mild requirement due to the universality of Hölder class.

372 Next we enumerate the assumptions about the data augmentations A.

Assumption 4.2 (Lipschitz augmentation). Any data augmentation $A_{\gamma} \in \mathcal{A}$ is *M*-Lipschitz, i.e., $\|A_{\gamma}(\boldsymbol{u}_1) - A_{\gamma}(\boldsymbol{u}_2)\|_2 \leq M \|\boldsymbol{u}_1 - \boldsymbol{u}_2\|_2$ for any $\boldsymbol{u}_1, \boldsymbol{u}_2 \in [0, 1]^d$.

A typical example to understand Assumption 4.2 is that the resulting augmented data obtained through cropping would not undergo drastic changes when minor perturbations are applied to the original image. 378 Denote the corresponding latent classes on source domain by $\{C_s(k)\}_{k \in [K]}$. Beyond the general 379 assumption regarding data augmentation \mathcal{A} above, we require a more precise way to describe the 380 intensity of data augmentations A. A more general version of the (σ, δ) -augmentation employed by 381 Huang et al. (2023) is adopted by us to distinguish the efficiency of data augmentations.

382 **Definition 4.2** ($(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -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 $\hat{C}_s(k) \subseteq C_s(k)$ and $\hat{C}_t(k) \subseteq C_t(k)$, such that

$$P_{s}(\boldsymbol{x} \in C_{s}(k)) \geq \sigma_{s}P_{s}(\boldsymbol{x} \in C_{s}(k)), \quad \sup_{\boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \widetilde{C}_{s}(k)} \min_{\boldsymbol{x}_{1}' \in \mathcal{A}(\boldsymbol{x}_{1}), \boldsymbol{x}_{2}' \in \mathcal{A}(\boldsymbol{x}_{2})} \|\boldsymbol{x}_{1}' - \boldsymbol{x}_{2}'\|_{2} \leq \delta_{s},$$

$$P_{t}(\boldsymbol{z} \in \widetilde{C}_{t}(k)) \geq \sigma_{t}P_{t}(\boldsymbol{z} \in C_{t}(k)), \quad \sup_{\boldsymbol{z}_{1}, \boldsymbol{z}_{2} \in \widetilde{C}_{t}(k)} \min_{\boldsymbol{z}_{1}' \in \mathcal{A}(\boldsymbol{z}_{1}), \boldsymbol{z}_{2}' \in \mathcal{A}(\boldsymbol{z}_{2})} \|\boldsymbol{z}_{1}' - \boldsymbol{z}_{2}'\|_{2} \leq \delta_{t},$$

$$K$$

$$P_t\big(\bigcup_{k=1}^K \widetilde{C}_t(k)\big) \ge \sigma_t,$$

393 where $\sigma_s, \sigma_t \in (0, 1]$ and $\delta_s, \delta_t \ge 0$.

394 *Remark* 4.1. The $(\sigma_s, \sigma_t, \delta_s, \delta_t)$ -augmentation methods emphasize that a robust data augmentation 395 should adhere to the principle that when the semantic information of the original images exhibit 396 heightened similarity, augmented views from them should be close according to specific crite-397 ria. Among above requirements, $P_t\left(\bigcup_{k=1}^{n} \widetilde{C}_t(k)\right) \geq \sigma_t$, which is used to replace the assumption 398 399 $\mathcal{A}(C_t(i)) \cap \mathcal{A}(C_t(j)) = \emptyset$ of Huang et al. (2023), implies that the augmentation used should be in-400 telligent enough to recognize objectives aligned with the image labels for the majority of samples in 401 the dataset. For instance, consider a downstream task involving classifying images of cats and dogs, 402 where the dataset includes some images featuring both cats and dogs together. This requirement demands that the data augmentation intelligently selects dog-specific augmentations when the image 403 is labeled as dog, and similarly for cat-specific augmentations when the image is labeled as cat. A 404 simple alternative to this requirement is assuming different class $C_t(k)$ are pairwise disjoint, i.e., 405 406

$$\forall i \neq j, C_t(i) \cap C_t(j) = \emptyset, \text{ which implies } P_t\left(\bigcup_{k=1}^K \widetilde{C}_t(k)\right) = \sum_{k=1}^K P_t(\widetilde{C}_t(k)) \ge \sigma_t \sum_{k=1}^K P_t(C_t(k)) = \sigma_t.$$

Assumption 4.3 (Existence of augmentation sequence). Assume there exists a sequence of 409 $(\sigma_s^{(n_s)}, \sigma_t^{(n_s)}, \delta_s^{(n_s)}, \delta_t^{(n_s)})$ -data augmentations $\mathcal{A}_{n_s} = \{A_{\gamma}^{(n_s)}(\cdot) : \gamma \in [m]\}$ and $\tau > 0$ such that 410

426 427

429

407 408

384

 $\max\{\delta_s^{(n_s)}, \delta_t^{(n_s)}\} \le n_s^{-\frac{\tau+d+1}{2(\alpha+d+1)}}, \quad \min\{\sigma_s^{(n_s)}, \sigma_t^{(n_s)}\} \stackrel{n_s \to \infty}{\to} 1$

413 It is worth mentioning that this assumption essentially aligns with Assumption 3.5 in HaoChen 414 et al. (2021), both stipulating the augmentations must be sufficiently robust so that the internal 415 connections within latent classes are strong enough to prevent instance clusters from being separated. 416 Recently, methods for building stronger data augmentation, as discussed by Jahanian et al. (2022) 417 and Trabucco et al. (2024), are constantly being proposed, making it more feasible to meet the 418 theoretical requirements for data augmentation.

419 Next we are going to introduce the assumption about distribution shift. For simplicity, denote 420 $p_s(k) = P_s(x \in C_s(k))$ and $P_s(k)$ be the conditional distribution of $P_s(x|x \in C_s(k))$ on the 421 upstream data, $p_t(k) = P_t(z \in C_t(k))$ and $P_t(k)$ be the conditional distribution $P_t(z | z \in C_t(k))$ 422 on the downstream task. Following assumption is needed to quantify our requirement on domain 423 shift. 424

Assumption 4.4. Assume there exists $\nu > 0$ and $\varsigma > 0$ such that 425

$$\max_{k \in [K]} \mathcal{W}(P_s(k), P_t(k)) \le n_s^{-\frac{\nu+d+1}{2(\alpha+d+1)}}, \quad \max_{k \in [K]} |p_s(k) - p_t(k)| \le n_s^{-\frac{\varsigma}{2(\alpha+d+1)}}.$$

428 where \mathcal{W} is the Wasserstein-1 distance.

A trivial scenario occurs when there is no gap between the upstream and downstream dis-430 tributions, i.e., when $(\mathcal{X}_s, P_s) = (\mathcal{X}_t, P_t)$, leading to both $\max_{k \in [K]} \mathcal{W}(P_s(k), P_t(k))$ and 431 $\max_{k \in [K]} |p_s(k) - p_t(k)|$ vanishing.

432 433 434 **Assumption 4.5.** Assume there exists a measurable partition $\{\mathcal{P}_1, \ldots, \mathcal{P}_{d^*}\}$ of \mathcal{X}_s , such that $1/B_2^2 \leq P_s(\mathcal{P}_i) \leq 1/B_1^2$ for each $i \in [d^*]$.

Assumption 4.5 is used to construct a measurable function \tilde{f} with $B_1 \leq ||\tilde{f}||_2 \leq B_2$, such that $\mathcal{L}(\tilde{f}) = 0$, tackling one of theoretical challenges introduced in Theorem 4.2 of HaoChen & Ma (2023), further implying that $\mathcal{L}(f^*)$ vanishes (see B.2.6 for more details). It suggests that the data distribution in the source domain should not be overly singular. All common continuous distributions defined on Borel algebra apparently satisfy these requirements, as the measure of any single point is zero.

4.2 END-TO-END THEORETICAL GUARANTEE

Our main theoretical result is stated as follows.

Theorem 4.2. Suppose Assumptions 4.1-4.5 hold. Set the width, depth and the Lipschitz constraint of the deep neural network as

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

Then the following inequality holds

441 442

443

444 445

450 451 452

465

466

$$\mathbb{E}_{D_s}\left[\mathrm{Err}(Q_{\hat{f}_{n_s}})\right] \le (1 - \sigma_t^{(n_s)}) + \mathcal{O}(n_s^{-\frac{\min\{\alpha,\nu,\varsigma\}}{8(\alpha+d+1)}}),$$

453 454 with probability at least $\sigma_s^{(n_s)} - \mathcal{O}(n_s^{-\frac{\min\{\alpha,\nu,\varsigma,\tau\}}{16(\alpha+d+1)}}) - \mathcal{O}(\frac{1}{\sqrt{\min_k n_t(k)}})$ for n_s sufficiently large.

455 *Remark* 4.3. Note that only the probability term depends on the downstream sample size and the 456 failure probability decays rapidly with respect to $\min_k n_t(k)$ with order 1/2, implying that the 457 learned representation via ACT from a large amount of unlabeled data can indeed help capture 458 downstream knowledge, despite a limited downstream sample size. This demonstrates the proven 459 advantage of ACT and provides an explanation for the empirical success of few-shot learning, which aligns with the concept of K-way $\min_k n_t(k)$ -shot learning. Apart from that, note the conditions of 460 Theorem 4.2 only require $W \ge \mathcal{O}(n_s^{\frac{2d+\alpha}{4(\alpha+d+1)}}), L \ge \mathcal{O}(1)$ and $\mathcal{K} = \mathcal{O}(n_s^{\frac{d+1}{2(\alpha+d+1)}})$, which implies 461 462 that the number of network parameters could be arbitrarily large if we control the norm of weight 463 properly, which is coincide with the concept of over-parametrization. 464

4.3 PROOF SKETCH OF THEOREM 4.2

467 **Step 1.** In Appendix B.2.1, we initially investigate the sufficient condition for achieving a low error 468 rate in a downstream task at the population level in Lemma B.1. It reveals that the misclassifica-469 tion rate bounded by the strength of data augmentations $1 - \sigma_s$, and the augmented concentration, 470 represented by $R_t(\varepsilon, f)$. This dependence arises when the divergence between different classes, 471 quantified by $\mu_t(i)^{\top} \mu_t(j)$, is sufficiently dispersed.

472 Step 2. Subsequently in Appendix B.2.2 and B.2.3, we regard $\sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G)$ as the 473 weighted summation of $\mathcal{L}_{align}(f)$ and $\mathcal{L}_{div}(f)$, then attempt to show they are the upper bound 474 of $R_t(\varepsilon, f), \max_{i \neq j} |\mu_t(i)^\top \mu_t(j)|$ respectively in Lemma B.4, which implies that optimizing our 475 adversarial self-supervised learning loss is equivalent to optimize the upper bound of $R_t(\varepsilon, f)$ and 476 $\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)|$ simultaneously, because $\mathcal{L}_{align}(f)$ and $\mathcal{L}_{div}(f)$ are positive. Finally, apply 477 Lemma B.1 and Lemma B.4 to \hat{f}_{n_s} , combining with the Markov inequality, to conclude Theorem 478 B.1, which is population version of Theorem 4.2.

Step 3. To further obtain an end-to-end theoretical guarantee, we subsequently decompose $\mathcal{E}(\hat{f}_{n_s})$, the excess risk defined in Definition B.3, into three parts: statistical error: \mathcal{E}_{sta} , approximation error introduced by neural network class: $\mathcal{E}_{\mathcal{F}}$, and the error brought by $\hat{\mathcal{G}}$: $\mathcal{E}_{\hat{\mathcal{G}}}$ in Appendix B.2.7. Note that the unbiased design of ACT plays a key role in such misspecified decomposition. We successively deal each produced term. For $\mathbb{E}_{D_s}[\mathcal{E}_{sta}]$, we claim it can be bounded by $\frac{\mathcal{K}\sqrt{L}}{\sqrt{n_s}}$ by adopting some typical techniques of empirical process and the result claimed by Golowich et al. (2018) in Appendix B.2.8. For $\mathcal{E}_{\mathcal{F}}$, according to the existing conclusion of Jiao et al. (2023), we can show $\mathcal{E}_{\mathcal{F}}$ can be bounded by $\mathcal{K}^{-\alpha/(d+1)}$ in Appendix B.2.9. By leveraging the unbiased property of ACT, the problem bounding $\mathbb{E}_{D_s}[\mathcal{E}_{\widehat{G}}]$ can be transformed into a common problem of mean convergence rate, so that it can be controlled by $\frac{1}{n_s^{1/4}}$ with high probability, shown as Appendix B.2.10. Trading off over three errors helps us determine a appropriate \mathcal{K} to bound $\mathbb{E}_{D_s}[\mathcal{E}(\widehat{f}_{n_s})]$, more details is showed in Appendix B.2.11.

492 **Step 4.** However, $\mathcal{L}(f^*)$, the difference between the excess risk and the term $\mathcal{L}(\hat{f}_{n_s})$ involving 493 in Theorem B.1, still impedes us from building an end-to-end theoretical guarantee for ACT. To 494 address this issue, in Appendix B.2.6, we construct a representation making this term vanishing 495 under Assumption 4.5. Finally, just set appropriate parameters of Theorem B.1 to conclude Lemma 496 B.12, whose direct corollary is Theorem 4.2, and proof is presented in Appendix B.12. The bridge 497 between Lemma B.12 and Theorem 4.2 is shown in Appendix B.2.12.

498 499

500

5 CONCLUSIONS

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 finetuned 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.

Minimax rates for supervised transfer learning are established in Cai & Wei (2019); Kpotufe &
 Martinet (2021); Cai & Pu (2024). However, the minimax rate for unsupervised transfer learning
 remains unclear. Establishing a lower bound to gain a deeper understanding of our ACT model
 presents an interesting and challenging problem for future research.

513 514

References

- Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Dipendra Misra. Investigating the role of negatives in contrastive representation learning. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (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.
- Pranjal Awasthi, Nishanth Dikkala, and Pritish Kamath. Do more negative samples necessarily hurt in contrastive learning? In *International conference on machine learning*, pp. 1101–1116. PMLR, 2022.
- Wele Gedara Chaminda Bandara, Celso M. De Melo, and Vishal M. Patel. Guarding barlow twins against overfitting with mixed samples, 2023. URL https://arxiv.org/abs/2312.02151.
- Adrien Bardes, Jean Ponce, and Yann LeCun. 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.
- T. Tony Cai and Hongming Pu. Transfer learning for nonparametric regression: Non-asymptotic minimax analysis and adaptive procedure, 2024. URL https://arxiv.org/abs/2401.
 12272.
- T. Tony Cai and Hongji Wei. Transfer learning for nonparametric classification: Minimax rate and adaptive classifier, 2019. URL https://arxiv.org/abs/1906.02903.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.
 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020.

550

551

555

556

- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
 contrastive learning of visual representations. In *International conference on machine learning*,
 pp. 1597–1607. PMLR, 2020a.
- 547 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of* 548 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.
 - Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.
- Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self supervised representation learning. In *International conference on machine learning*, pp. 3015–
 3024. PMLR, 2021.
 - Quentin Garrido, Yubei Chen, Adrien Bardes, Laurent Najman, and Yann Lecun. On the duality between contrastive and non-contrastive self-supervised learning. *arXiv preprint arXiv:2206.02574*, 2022.
- E. Giné and R. Nickl. 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.
- Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet (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.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- Jeff Z. HaoChen and Tengyu Ma. 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.
- Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-supervised
 deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*,
 34:5000–5011, 2021.
- Jeff Z. HaoChen, Colin Wei, Ananya Kumar, and Tengyu Ma. Beyond separability: Analyzing the linear transferability of contrastive representations to related subpopulations. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (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/ ac112e8ffc4e5b9ece32070440a8ca43-Abstract-Conference.html.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Weiran Huang, Mingyang Yi, Xuyang Zhao, and Zihao Jiang. 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.

594 Ali Jahanian, Xavier Puig, Yonglong Tian, and Phillip Isola. Generative models as a data source 595 for multiview representation learning. In The Tenth International Conference on Learning 596 Representations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net, 2022. URL 597 https://openreview.net/forum?id=qhAeZjs7dCL. 598 Yuling Jiao, Yang Wang, and Yunfei Yang. Approximation bounds for norm constrained neural networks with applications to regression and gans. Applied and Computational Harmonic Analysis, 600 65:249-278, 2023. 601 602 Samory Kpotufe and Guillaume Martinet. Marginal singularity and the benefits of labels in 603 covariate-shift. The Annals of Statistics, 49(6):3299-3323, 2021. 604 Yunwen Lei, Tianbao Yang, Yiming Ying, and Ding-Xuan Zhou. Generalization analysis for con-605 trastive representation learning. In International Conference on Machine Learning, pp. 19200– 606 19227. PMLR, 2023. 607 608 Jit Yan Lim, Kian Ming Lim, Chin Poo Lee, and Yong Xuan Tan. Scl: Self-supervised contrastive 609 learning for few-shot image classification. *Neural Networks*, 165:19–30, 2023. 610 Chen Liu, Yanwei Fu, C. Xu, Siqian Yang, Jilin Li, Chengjie Wang, and Li Zhang. Learning a few-611 shot embedding model with contrastive learning. In AAAI Conference on Artificial Intelligence, 612 2021. URL https://api.semanticscholar.org/CorpusID:235349153. 613 614 Andreas Maurer. A vector-contraction inequality for rademacher complexities. In Algorithmic 615 Learning Theory: 27th International Conference, ALT 2016, Bari, Italy, October 19-21, 2016, 616 Proceedings 27, pp. 3–17. Springer, 2016. 617 Serdar Ozsoy, Shadi Hamdan, Sercan ö. Arik, Deniz Yuret, and Alper T. Erdogan. Self-supervised 618 learning with an information maximization criterion, 2022. URL https://arxiv.org/ 619 abs/2209.07999. 620 621 Mamshad Nayeem Rizve, Salman Khan, Fahad Shahbaz Khan, and Mubarak Shah. Exploring 622 complementary strengths of invariant and equivariant representations for few-shot learning. In 623 Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 10836– 624 10846, 2021. 625 Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khandeparkar. 626 A theoretical analysis of contrastive unsupervised representation learning. In Kamalika Chaudhuri 627 and Ruslan Salakhutdinov (eds.), Proceedings of the 36th International Conference on Machine 628 Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA, volume 97 of Proceedings 629 of Machine Learning Research, pp. 5628–5637. PMLR, 2019. URL http://proceedings. 630 mlr.press/v97/saunshi19a.html. 631 Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU acti-632 vation function. The Annals of Statistics, 48(4):1875 - 1897, 2020. doi: 10.1214/19-AOS1875. 633 URL https://doi.org/10.1214/19-AOS1875. 634 635 Brandon Trabucco, Kyle Doherty, Max Gurinas, and Ruslan Salakhutdinov. Effective data aug-636 mentation with diffusion models. In The Twelfth International Conference on Learning Rep-637 resentations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net, 2024. URL 638 https://openreview.net/forum?id=ZWzUA9zeAq. 639 A.B. Tsybakov. Introduction to Nonparametric Estimation. Springer Series in Statistics. Springer 640 New York, 2008. ISBN 9780387790527. URL https://books.google.com/books? 641 id=mwB8rUBsbqoC. 642 643 Roman Vershynin. High-Dimensional Probability: An Introduction with Applications in Data Sci-644 ence. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 645 2018. 646 Yifei Wang, Qi Zhang, Yaoyu Guo, and Yisen Wang. Non-negative contrastive learning. In ICLR, 647 2024.

- 648 Zhanyuan Yang, Jinghua Wang, and Yingying Zhu. Few-shot classification with contrastive learning. 649 In European conference on computer vision, pp. 293–309. Springer, 2022. 650
- Dmitry Yarotsky. Optimal approximation of continuous functions by very deep relu networks. 651 In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet (eds.), Proceedings of the 31st 652 Conference On Learning Theory, volume 75 of Proceedings of Machine Learning Research, 653 pp. 639-649. PMLR, 06-09 Jul 2018. URL https://proceedings.mlr.press/v75/ 654 yarotsky18a.html. 655
- Dmitry Yarotsky and Anton Zhevnerchuk. The phase diagram of approximation rates for deep neural 656 networks. Advances in neural information processing systems, 33:13005–13015, 2020. 657
- 658 Mang Ye, Xu Zhang, Pong C Yuen, and Shih-Fu Chang. Unsupervised embedding learning via 659 invariant and spreading instance feature. In Proceedings of the IEEE/CVF conference on computer 660 vision and pattern recognition, pp. 6210–6219, 2019.
 - Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In International conference on machine learning, pp. 12310-12320. PMLR, 2021.
 - Qi Zhang, Yifei Wang, and Yisen Wang. Identifiable contrastive learning with automatic feature importance discovery. In NeurIPS, 2023.
- 666 667 668

678 679

681 682

686 687

689

693 694

661

662

663

664

665

EXPLANATION OF THE REGULARIZATION TERM А

670 In brief, contrastive learning utilizes data augmentation to construct the loss function (specifically, 671 the first term in our loss) that aligns representations of the same class. However, to avoid trivial 672 solutions, an additional regularization term is necessary to ensure that clusters representing different classes are well-separated. We measure this separation using the angles between the centroids of 673 different classes. While these angles are ideal for quantifying separation, they cannot be directly 674 optimized because the latent class annotations are unavailable in the upstream task. As an alternative, 675 we propose finding an appropriate computable loss function that serves as an upper bound for these 676 angles, effectively achieving the desired separation. Denote 677

$$\mathcal{L}_{\text{div}}(f) = \left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*} \right\|_F^2.$$

It can severs as a regularization term since in Lemma B.4, we can show 680

$$\mu_{s}(i)^{\top}\mu_{s}(j) \lesssim \left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_{1},\boldsymbol{x}_{2} \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_{1})f(\boldsymbol{x}_{2})^{\top}] - I_{d^{*}} \right\|_{F},$$
(12)

where $\mu_s(i) = \mathbb{E}_{\boldsymbol{x} \in C_s(i)} \mathbb{E}_{\boldsymbol{x}' \in \mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}')]$ is the center of the latent class *i*. (12) implies that a lower 683 value of the regularization term leads the separation between different categories' center, thereby 684 benefits classification in downstream tasks. 685

At the sample level, one can use
$$\widehat{\mathcal{L}}_{\text{div}}(f) = \left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - I_{d^*} \right\|_F^2$$
 to estimate $\mathcal{L}_{\text{div}}(f)$.

However, this lead to a bias loss, i.e. 688

 $\mathbb{E}_{D_s}[\widehat{\mathcal{L}}_{\operatorname{div}}(f)] \neq \mathcal{L}_{\operatorname{div}}(f),$

690 where D_s is augmented dataset. This bias is caused by the non-commutativity of the expectation 691 and the Frobenius norm. To overcome this we can reformulate it as an equivalent form 692

$$\mathcal{L}_{\operatorname{div}}(f) = \sup_{G \in \mathcal{G}(f)} \mathcal{R}(f, G) := \langle \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*}, G \rangle_F$$

The counterpart of $\mathcal{R}(f, G)$ at the sample level is

$$\widehat{\mathcal{R}}(f,G) = \langle \frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - I_{d^*}, G \rangle_F.$$

We can see that $\mathbb{E}_{D_{e}}[\mathcal{R}(f,G)] = \mathcal{R}(f,G)$ for any fixed G due to the linearity of Frobenius inner 699 product, combining this property with the new decomposition method proposed by us, we build an 700 end-to-end theoretical guarantee in the transfer learning setting to provide an explanation for few 701 shot learning. And using an alternative optimization method to optimize this loss is natural.

702 B DEFERRED PROOF

704

705

706 707

708 709

710

711 712

715 716 717

727 728 729 The Section B will be divided into two parts. The first part B.1 is used to prove $\|\phi_{\theta}\|_{\text{Lip}} \leq \mathcal{K}$ for any $\phi_{\theta} \in \mathcal{NN}_{d_1,d_2}(W,L,\mathcal{K},B_1,B_2)$. The proof of Theorem 4.2 is shown in the second part B.2.

B.1 *K*-LIPSCHITZ PROPERTY OF $\mathcal{NN}_{d_1,d_2}(W,L,\mathcal{K},B_1,B_2)$

Proof. To claim any $\phi_{\theta} \in \mathcal{NN}_{d_1,d_2}(W, L, \mathcal{K}, B_1, B_2)$ is \mathcal{K} -Lipschitz function, we need to define two special classes of neural network functions, the first is

$$\mathcal{NN}_{d_1,d_2}(W,L,\mathcal{K}) := \{ \phi_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{A}_L \sigma(\boldsymbol{A}_{L-1}\sigma(\cdots\sigma(\boldsymbol{A}_0\boldsymbol{x})) : \kappa(\boldsymbol{\theta}) \le \mathcal{K} \},$$
(13)

which equivalent to $\mathcal{NN}_{d_1,d_2}(W,L,\mathcal{K},B_1,B_2)$ ignoring the condition $\|\phi_{\theta}\|_2 \in [B_1,B_2]$, and the second one

$$\mathcal{SNN}_{d_1,d_2}(W,L,\mathcal{K}) := \{ \breve{\phi}(\boldsymbol{x}) = \breve{\boldsymbol{A}}_L \sigma(\breve{\boldsymbol{A}}_{L-1}\sigma(\cdots\sigma(\breve{\boldsymbol{A}}_0\breve{\boldsymbol{x}})) : \prod_{l=1}^L \|\breve{\boldsymbol{A}}_l\|_{\infty} \le \mathcal{K} \}, \qquad \breve{\boldsymbol{x}} := \begin{pmatrix} \boldsymbol{x} \\ 1 \end{pmatrix}$$

718 719 where $\breve{A}_l \in \mathbb{R}^{N_{l+1} \times N_l}$ with $N_0 = d_1 + 1$.

720 It is obvious 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 721 $\mathcal{SNN}_{d_1,d_2}(W,L,\mathcal{K})$ is \mathcal{K} -Lipschitz function as the 1-Lipschitz property of ReLU, thus it suffices 722 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 yield what 723 we desired.

724 725 726 In fact, any $\phi_{\theta}(\boldsymbol{x}) = \boldsymbol{A}_{L}\sigma(\boldsymbol{A}_{L-1}\sigma(\cdots\sigma(\boldsymbol{A}_{0}\boldsymbol{x}+\boldsymbol{b}_{0})) + \boldsymbol{b}_{L-1}) \in \mathcal{NN}_{d_{1},d_{2}}(W,L,\mathcal{K})$ can be rewritten as $\check{\phi}(\boldsymbol{x}) = \check{\boldsymbol{A}}_{L}\sigma(\check{\boldsymbol{A}}_{L-1}\sigma(\cdots\sigma(\check{\boldsymbol{A}}_{0}\check{\boldsymbol{x}})))$, where

$$\breve{\boldsymbol{x}} := \begin{pmatrix} \boldsymbol{x} \\ 1 \end{pmatrix}, \breve{\boldsymbol{A}} = (\boldsymbol{A}_L, \boldsymbol{0}), \breve{\boldsymbol{A}}_l = \begin{pmatrix} \boldsymbol{A}_l & \boldsymbol{b}_l \\ \boldsymbol{0} & 1 \end{pmatrix}, l = 0, \dots, L-1.$$

730 Notice that $\prod_{l=0}^{L} \|\breve{\boldsymbol{A}}\|_{\infty} = \|\boldsymbol{A}_L\|_{\infty} \prod_{l=0}^{L-1} \max\{\|(\boldsymbol{A}_l, \boldsymbol{b}_l)\|_{\infty}, 1\} = \kappa(\boldsymbol{\theta}) \leq \mathcal{K}$, which implies that 731 $\phi_{\boldsymbol{\theta}} \in SNN_{d_1, d_2}(W+1, L, \mathcal{K}).$

Conversely, since any $\check{\phi} \in \mathcal{SNN}(W, L, \mathcal{K})$ can also be parameterized in the form of $A_L \sigma(A_{L-1}\sigma(\cdots\sigma(A_0x + b_0)) + b_{L-1})$ with $\theta = (\check{A}_0, (\check{A}_1, \mathbf{0}), \dots, (\check{A}_{L-1}, \mathbf{0}), \check{A}_L)$, and by the absolute homogeneity of the ReLU function, we can always rescale \check{A}_l such that $||\check{A}_L||_{\infty} \leq \mathcal{K}$ and $||\check{A}_l||_{\infty} = 1$ for $l \neq L$. Hence $\kappa(\theta) = \prod_{l=0}^L ||\check{A}_l||_{\infty} \leq \mathcal{K}$, which yields $\check{\phi} \in \mathcal{NN}(W, L, \mathcal{K})$. \Box

737 738 E

B.2 PROOF OF THEOREM 4.2

739 We will begin by exploring the sufficient condition for achieving small $Err(Q_f)$ in B.2.1. Follow-740 ing that, we build the connection between the required condition and optimizing our adversarial 741 self-supervised learning loss in Theorem B.1 of B.2.3, it reveals that small quantity of our loss 742 function may induce significant class divergence and high augmented concentration. Although this 743 theorem can explain the essential factors behind the success of our method to some extent, its anal-744 ysis still stay at population level, the impact of sample size on $Err(Q_f)$ remains unresolved. To 745 obtain an end-to-end theoretical guarantee as Theorem 4.2, we first decompose $\mathcal{E}(f_{n_0})$, which is the 746 excess risk defined in the Definition B.3, into three parts: statistical error: \mathcal{E}_{sta} , approximation error 747 brought by \mathcal{F} : $\mathcal{E}_{\mathcal{F}}$ and the error introduced by using $\mathcal{G}(f)$ to approximate $\mathcal{G}(f)$: $\mathcal{E}_{\widehat{G}}$ in B.2.7, then 748 successively deal each produced term. For $\mathbb{E}_{D_s}[\mathcal{E}_{sta}]$, we adopt some typical techniques of empirical 749 process and the result provided by Golowich et al. (2018) in B.2.8 for bounding it by $\frac{\mathcal{K}\sqrt{L}}{\sqrt{n_e}}$. Regard-750 ing bounding $\mathcal{E}_{\mathcal{F}}$, we first convert the problem to a function approximation problem and adopt the 751 existing conclusion proposed by Jiao et al. (2023), yielding $\mathcal{E}_{\mathcal{F}}$ can be bounded by $\mathcal{K}^{-\alpha/(d+1)}$ in 752 B.2.9. By leveraging the property $\mathbb{E}_{D_s}[\hat{\mathcal{L}}(f,G)] = \mathcal{L}(f,G)$, we find that the problem of bounding 753 $\mathbb{E}_{D_s}[\mathcal{E}_{\widehat{G}}]$ can be transformed into a common problem of mean convergence rate and further control 754 it by $\frac{1}{n_s^{1/4}}$ in B.2.10. After finishing these preliminaries, trade off between these errors to determine 755 a appropriate Lipschitz constant \mathcal{K} of neural network, while bound the expectation of excess risk

 $\mathbb{E}_{D_s}[\mathcal{E}(\hat{f}_{n_s})]$, more details are deferred to B.2.11. However, $\sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G)$, the difference between the excess risk and the term $\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)$ involving in Theorem B.1, still im-pedes us from building an end-to-end theoretical guarantee for ACT. To address this issue, in B.2.6, we construct a representation making this term vanishing under Assumption 4.5. Finally, just set appropriate parameters of Theorem B.1 to conclude Lemma B.12, and the bridge between Lemma B.12 and Theorem 4.2 is built in B.2.12.

B.2.1 SUFFICIENT CONDITION OF SMALL MISCLASSIFICATION RATE

Lemma B.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 K-Lipschitz and

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

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

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

where $\varepsilon > 0$, $\mu_t(k) = \mathbb{E}_{\boldsymbol{z} \in C_t(k)} \mathbb{E}_{\boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})}[f(\boldsymbol{z}')]$ for any $k \in [K]$, $\Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, f) = (\sigma_t - \sigma_t)$ $\frac{R_t(\varepsilon,f)}{\min_i p_t(i)} \Big) \Big(1 + \Big(\frac{B_1}{B_2} \Big)^2 - \frac{K\delta_t}{B_2} - \frac{2\varepsilon}{B_2} \Big) - 1, \ \Delta_{\hat{\mu}_t} = 1 - \frac{\min_{k \in [K]} \|\hat{\mu}_t(k)\|_2^2}{B_2^2}, \ R_t(\varepsilon,f) = P_t \Big(\boldsymbol{z} \in \bigcup_{k=1}^K C_t(k) : \sup_{\boldsymbol{z}_1, \boldsymbol{z}_2 \in \mathcal{A}(\boldsymbol{z})} \|f(\boldsymbol{z}_1) - f(\boldsymbol{z}_2)\|_2 > \varepsilon \Big) \text{ and } \Theta(\sigma_t, \delta_t, \varepsilon, f) = \Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, f) - \sqrt{2 - 2\Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, f)} - \frac{\Delta_{\hat{\mu}_t}}{2} - \frac{2\max_{k \in [K]} \|\hat{\mu}_t(k) - \mu_t(k)\|_2}{B_2}.$

 $\textit{Proof. For any encoder } f, \text{ let } S_t(\varepsilon, f) := \{ \boldsymbol{z} \in \cup_{k=1}^K C_t(k) : \sup_{\boldsymbol{z}_1, \boldsymbol{z}_2 \in \mathcal{A}(\boldsymbol{z})} \| f(\boldsymbol{z}_1) - f(\boldsymbol{z}_2) \|_2 \leq C_t(\varepsilon, f) \}$ ε }, if any $z \in (\widetilde{C}_t(1) \cup \cdots \cup \widetilde{C}_t(K)) \cap S_t(\varepsilon, f)$ can be correctly classified by Q_f , it turns out that $\operatorname{Err}(Q_f)$ can be bounded by $(1 - \sigma_t) + R_t(\varepsilon, f)$. In fact,

$$\operatorname{Err}(Q_f) = \sum_{k=1}^{K} P_t \left(Q_f(\boldsymbol{z}) \neq k, \forall \boldsymbol{z} \in C_t(k) \right)$$
$$\leq P_t \left(\left(\widetilde{C}_t(1) \cup \dots \cup \widetilde{C}_t(K) \right) \cap S_t(\varepsilon, f) \right)^c \right)$$
$$= P_t \left(\left(\widetilde{C}_t(1) \cup \dots \cup \widetilde{C}_t(K) \right)^c \cup \left(S_t(\varepsilon, f) \right)^c \right)$$
$$\leq (1 - \sigma_t) + P_t \left((S_t(\varepsilon, f))^c \right)$$
$$= (1 - \sigma_t) + R_t(\varepsilon, f).$$

The first row is derived according to the definition of $\operatorname{Err}(Q_f)$. Since any $z \in (\widetilde{C}_t(1) \cup \cdots \cup$ $\widetilde{C}_t(K)) \cap S_t(\varepsilon, f)$ can be correctly classified by Q_f , we yields the second row. De Morgan's laws implies the third row. The fourth row stems from the Definition 4.2. Finally, just note $R_t(\varepsilon, f) =$ $(S_t(\varepsilon, f))^c$ to obtain the last line.

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

$$\mu_{t}(i)^{\top}\mu_{t}(j) < B_{2}^{2} \Big(\Gamma_{i}(\sigma_{t}, \delta_{t}, \varepsilon, f) - \sqrt{2 - 2\Gamma_{i}(\sigma_{t}, \delta_{t}, \varepsilon, f)} - \frac{\Delta_{\hat{\mu}_{t}}}{2} - \frac{\|\hat{\mu}_{t}(i) - \mu_{t}(i)\|_{2}}{B_{2}} - \frac{\|\hat{\mu}_{t}(j) - \mu_{t}(j)\|_{2}}{B_{2}} \Big),$$

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

To this end, without losing generality, consider the case i = 1. To turn out $z_0 \in \widetilde{C}_t(1) \cap S_t(\varepsilon, f)$ can be correctly classified 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 is equivalent to

$$f(\boldsymbol{z}_0)^{\top} \hat{\mu}_t(1) - f(\boldsymbol{z}_0)^{\top} \hat{\mu}_t(k) - \left(\frac{1}{2} \| \hat{\mu}_t(1) \|_2^2 - \frac{1}{2} \| \hat{\mu}_t(k) \|_2^2\right) > 0.$$

We will firstly deal with the term
$$f(z_0)^{\top} \hat{\mu}_t(1)$$
,
 $f(z_0)^{\top} \hat{\mu}_t(1) = f(z_0)^{\top} \mu_t(1) + f(z_0)^{\top} (\hat{\mu}_t(1) - \mu_t(1))$
 $\geq f(z_0)^{\top} \mathbb{E}_{z \in C_t(1)} \mathbb{E}_{z' \in A(z)} [f(z')] = ||f(z_0)||_2 ||\hat{\mu}_t(1) - \mu_t(1)||_2$
 $\geq \frac{1}{p_t(1)} f(z_0)^{\top} \mathbb{E}_z \mathbb{E}_{z' \in A(z)} [f(z')] = C_t(1) \cap C_t(1) - \mu_t(1)||_2$
 $= \frac{1}{p_t(1)} f(z_0)^{\top} \mathbb{E}_z \mathbb{E}_{z' \in A(z)} [f(z')] = C_t(1) \cap C_t(1) \cap S_t(\varepsilon, f)]$
 $+ \frac{1}{p_t(1)} f(z_0)^{\top} \mathbb{E}_z \mathbb{E}_{z' \in A(z)} [f(z')] = C_t(1) \cap (C_t(1) \cap S_t(\varepsilon, f))^c]$
 $- B_2 ||\hat{\mu}_t(1) - \mu_t(1)||_2$
 $= \frac{P_t(\tilde{C}_t(1) \cap S_t(\varepsilon, f))}{p_t(1)} f(z_0)^{\top} \mathbb{E}_{z \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{z' \in A(z)} [f(z')]$
 $+ \frac{1}{p_t(1)} \mathbb{E}_z [\mathbb{E}_{z' \in A(z)} [f(z_0)^{\top} f(z')] = C_t(1) \setminus (\tilde{C}_t(1) \cap S_t(\varepsilon, f))]$
 $- B_2 ||\hat{\mu}_t(1) - \mu_t(1)||_2$
 $\geq \frac{P_t(\tilde{C}_t(1) \cap S_t(\varepsilon, f))}{p_t(1)} f(z_0)^{\top} \mathbb{E}_{z \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{z' \in A(z)} [f(z')]$
 $- B_2 ||\hat{\mu}_t(1) - \mu_t(1)||_2$
 $\geq \frac{P_t(\tilde{C}_t(1) \cap S_t(\varepsilon, f))}{p_t(1)} f(z_0)^{\top} \mathbb{E}_{z \in \tilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{z' \in A(z)} [f(z')]$
 $- \frac{B_2^2}{p_t(1)} P_t(C_t(1) \setminus (\tilde{C}_t(1) \cap S_t(\varepsilon, f))) - B_2 ||\hat{\mu}_t(1) - \mu_t(1)||_2, \quad (14)$

where the second row stems from Cauchy-Schwarz inequality. The third and the last rows are according to the condition $||f||_2 \leq B_2$.

Note that

$$P_t(C_t(1) \setminus (\widetilde{C}_t(1) \cap S_t(\varepsilon, f))) = P_t((C_t(1) \setminus \widetilde{C}_t(1)) \cup (\widetilde{C}_t(1) \cap (S_t(\varepsilon, f))^c))$$

$$\leq (1 - \sigma_t) p_t(1) + R_t(\varepsilon, f),$$
(15)

and

$$P_t(\widetilde{C}_t(1) \cap S_t(\varepsilon, f)) = P_t(C_t(1)) - P_t(C_t(1) \setminus (\widetilde{C}_t(1) \cap S_t(\varepsilon, f)))$$

$$\geq p_t(1) - ((1 - \sigma_t)p_t(1) + R_t(\varepsilon, f))$$

$$= \sigma_t p_t(1) - R_t(\varepsilon, f).$$
(16)

Plugging (15), (16) into (14) yields

$$f(\boldsymbol{z}_{0})^{\top}\hat{\mu}_{t}(1) \geq \left(\sigma_{t} - \frac{R_{t}(\varepsilon, f)}{p_{t}(1)}\right)f(\boldsymbol{z}_{0})^{\top}\mathbb{E}_{\boldsymbol{z}\in\widetilde{C}_{t}(1)\cap S_{t}(\varepsilon, f)}\mathbb{E}_{\boldsymbol{z}'\in\mathcal{A}(\boldsymbol{z})}[f(\boldsymbol{z}')] - B_{2}^{2}\left(1 - \sigma_{t} + \frac{R_{t}(\varepsilon, f)}{p_{t}(1)}\right) - B_{2}\|\hat{\mu}_{t}(1) - \mu_{t}(1)\|_{2}.$$
(17)

Notice that $z_0 \in \widetilde{C}_t(1) \cap S_t(\varepsilon, f)$. Thus for any $z \in \widetilde{C}_t(1) \cap S_t(\varepsilon, f)$, by the defini-tion of $\widetilde{C}_t(1)$, we have $\min_{\boldsymbol{z}_0' \in \mathcal{A}(\boldsymbol{z}_0), \boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})} \|\boldsymbol{z}_0' - \boldsymbol{z}'\|_2 \leq \delta_t$. Further denote $(\boldsymbol{z}_0^*, \boldsymbol{z}^*) = \delta_t$ $\arg\min_{\boldsymbol{z}_0' \in \mathcal{A}(\boldsymbol{z}_0), \boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})} \|\boldsymbol{z}_0' - \boldsymbol{z}'\|_2$, then $\|\boldsymbol{z}_0^* - \boldsymbol{z}^*\|_2 \leq \delta_t$, combining \mathcal{K} -Lipschitz property of f to yield $||f(z_0^*) - f(z^*)||_2 \leq \mathcal{K} ||z_0^* - z^*||_2 \leq \mathcal{K}\delta_t$. Besides that, since $z \in S_t(\varepsilon, f), \forall z' \in \mathcal{A}(z), ||f(z') - f(z^*)||_2 \leq \varepsilon$. Similarly, as $z_0 \in S_t(\varepsilon, f)$ and $z_0, z_0^* \in \mathcal{A}(z_0)$, we know $\|f(\boldsymbol{z}_0) - f(\boldsymbol{z}_0^*)\|_2 \le \varepsilon.$

Therefore,

$$\begin{array}{ll} \mathbf{859} & f(\boldsymbol{z}_0)^\top \mathbb{E}_{\boldsymbol{z} \in \widetilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})}[f(\boldsymbol{z}')] \\ \mathbf{860} & = \mathbb{E}_{\boldsymbol{z} \in \widetilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})}[f(\boldsymbol{z}_0)^\top f(\boldsymbol{z}_0)^\top f(\boldsymbol{z}_0)] \\ \mathbf{862} & = \mathbb{E}_{\boldsymbol{z} \in \widetilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})}[f(\boldsymbol{z}_0)^\top (f(\boldsymbol{z}_0)^\top f(\boldsymbol{z}_0)] \\ \mathbf{863} & = \mathbb{E}_{\boldsymbol{z} \in \widetilde{C}_t(1) \cap S_t(\varepsilon, f)} \mathbb{E}_{\boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})}[f(\boldsymbol{z}_0)^\top f(\boldsymbol{z}_0)] \\ \end{array}$$

$$= \mathbb{E}_{\boldsymbol{z}\in \widetilde{C}_{t}(1)\cap S_{t}(\varepsilon,f)} \mathbb{E}_{\boldsymbol{z}'\in\mathcal{A}(\boldsymbol{z})} [f(\boldsymbol{z}_{0})^{\top}f(\boldsymbol{z}')]$$

$$= \mathbb{E}_{\boldsymbol{z}\in \widetilde{C}_{t}(1)\cap S_{t}(\varepsilon,f)} \mathbb{E}_{\boldsymbol{z}'\in\mathcal{A}(\boldsymbol{z})} [f(\boldsymbol{z}_{0})^{\top}(f(\boldsymbol{z}') - f(\boldsymbol{z}_{0}) + f(\boldsymbol{z}_{0}))]$$

$$\geq B_{1}^{2} + \mathbb{E}_{\boldsymbol{z}\in \widetilde{C}_{t}(1)\cap S_{t}(\varepsilon,f)} \mathbb{E}_{\boldsymbol{z}'\in\mathcal{A}(\boldsymbol{z})} [f(\boldsymbol{z}_{0})^{\top}(f(\boldsymbol{z}') - f(\boldsymbol{z}_{0}))]$$

$$\begin{split} &=B_1^2 + \mathbb{E}_{z\in\tilde{C}_1(1)\cap S_1(\varepsilon,f)} \mathbb{E}_{z'\in A(z)} [f(z_0)^\top (f(z') - f(z') + f(z') - f(z_0) + f(z_0) + f(z_0) - f(z_0))] \\ &\geq B_1^2 - B_2(\mathcal{K}\delta_t + B_2 \mathcal{K}) \\ &\geq B_1^2 - B_2(\mathcal{K}\delta_t + 2\varepsilon), \end{split} (18) \\ &\text{where the fourth row is derived from $\|f\|_2 \geq B_1. \\ &\text{Plugging (18) to the inequality (17) knows} \\ &f(z_0)^\top \hat{\mu}_t(1) \geq \left(\sigma_t - \frac{R_t(\varepsilon,f)}{p_t(1)}\right) f(z_0)^\top \sum_{z\in\tilde{C}_t(1)\cap S_t(\varepsilon,f)} \sum_{z'\in\tilde{\mathcal{K}}(z_0)} [f(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 + (\frac{B_1}{D_2})^2\right) \left(\sigma_t - \frac{R_t(\varepsilon,f)}{B_2}\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) \\ &- D_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 + (\frac{B_1}{B_2})^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 \left(\left(\sigma_t - \frac{R_t(\varepsilon,f)}{p_t(1)}\right) \left(1 + (\frac{B_1}{B_2})^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 \left(\left(\sigma_t - \frac{R_t(\varepsilon,f)}{p_t(1)}\right) \left(1 + (\frac{B_1}{B_2})^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 \left(\left(\sigma_t - \frac{R_t(\varepsilon,f)}{p_t(1)}\right) \left(1 + (\frac{B_1}{B_2})^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 \left(\int (\sigma_t - \frac{R_t(\varepsilon,f)}{p_t(1)}\right) \left(1 + (\frac{B_1}{B_2})^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 \left(\int (\sigma_t - \frac{R_t(\varepsilon,f)}{p_t(1)}\right) \left(1 + (\frac{R_1}{B_2})^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 \left(\int (\sigma_t - \frac{R_t(\varepsilon,f)}{p_t(1)}\right) \left(1 + (\frac{R_1}{B_2})^2 - \frac{\mathcal{K}\delta_t}{B_2} - \frac{2\varepsilon}{B_2}\right) - 1\right) \\ (\text{combining the fact that} \\ \|\mu_t(k)\|_2 = \|B_2\varepsilon\hat{C}_t(k)|_E^2 + (k, 1) - \mu_t(k)|_E \\ &\leq f(z_0)^\top \mu_t(k) + f(z_0) \left[\frac{R_t(k)}{p_t(k)}\right] - \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 \\ &\leq f(z_0)^\top \mu_t(k) + B_1^2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 \\ &\leq f(z_0)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 \\ &\leq f(z_0)^\top \mu_t(k) + B_2 \|\hat{\mu}_t(k) + B_2 \|\hat{\mu}_t(k) -$$$

$$\geq B_2^2 \Gamma_1(\sigma_t, \delta_t, \varepsilon, f) - B_2 \|\hat{\mu}_t(1) - \mu_t(1)\|_2 - \sqrt{2} B_2^2 \sqrt{1 - \Gamma_1(\sigma_t, \delta_t, \varepsilon, \varepsilon, f)} \\ - \mu_t(1)^\top \mu_t(k) - B_2 \|\hat{\mu}_t(k) - \mu_t(k)\|_2 - \frac{1}{2} B_2^2 \Delta_{\hat{\mu}_t} > 0,$$

f)

which finishes the proof.

B.2.2 PRELIMINARIES FOR LEMMA B.4

To establish Lemma B.4, we must first prove Lemmas B.2 and B.3 in advance. Following the notations in the target domain, we employ $\mu_s(k) := \mathbb{E}_{\boldsymbol{x} \in C_s(k)} \mathbb{E}_{\boldsymbol{x}' \in \mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}')] = \frac{1}{p_s(k)} \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}' \in \mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}') \mathbb{1}\{\boldsymbol{x} \in C_s(k)\}]$ to denote the centre of k-th latent class in representation space. Apart from that, it is necessary to introduce following assumption, which is the abstract version of Assumption 4.4.

Assumption B.1. Review $P_s(k)$ and $P_t(k)$ are the conditional measures that $P(\boldsymbol{x}|\boldsymbol{x} \in C_s(k))$ and $P(\boldsymbol{z}|\boldsymbol{z} \in C_t(k))$ respectively, assume $\exists \rho > 0$ and $\eta > 0$, $\max_{k \in [K]} \mathcal{W}(P_s(k), P_t(k)) \leq \rho$ and $\max_{k \in [K]} |p_s(k) - p_t(k)| \leq \eta$.

 $\lim_{k \in [K]} |p_s(k) - p_t(k)| \le \eta.$

Lemma B.2. If the encoder f is \mathcal{K} -Lipschitz and Assumption B.1 holds, for any $k \in [K]$, we have:

$$\|\mu_s(k) - \mu_t(k)\|_2 \le \sqrt{d^*} M \mathcal{K} \rho$$

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

$$\begin{aligned} \|\mu_{s}(k) - \mu_{t}(k)\|_{2}^{2} &= \sum_{l=1}^{d^{*}} \left((\mu_{s}(k))_{l} - (\mu_{t}(k))_{l} \right)^{2} \\ &= \sum_{l=1}^{d^{*}} (\mathbb{E}_{\boldsymbol{x} \in C_{s}(k)} \mathbb{E}_{\boldsymbol{x}' \in \mathcal{A}(\boldsymbol{x})} [f_{l}(\boldsymbol{x}')] - \mathbb{E}_{\boldsymbol{z} \in C_{t}(k)} \mathbb{E}_{\boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})} [f_{l}(\boldsymbol{z}')])^{2} \\ &= \sum_{l=1}^{d^{*}} \left[\frac{1}{m} \sum_{\gamma=1}^{m} \left(\mathbb{E}_{\boldsymbol{x} \in C_{s}(k)} [f_{l}(A_{\gamma}(\boldsymbol{x}))] - \mathbb{E}_{\boldsymbol{z} \in C_{t}(k)} [f_{l}(A_{\gamma}(\boldsymbol{z}))] \right) \right]^{2} \\ &\leq d^{*} M^{2} \mathcal{K}^{2} \rho^{2} \end{aligned}$$

The final inequality is obtained by Assumption B.1 along with the fact that $f(A_{\gamma}(\cdot))$ is $M\mathcal{K}$ -Lipschitz continuous. In fact, as $f \in \operatorname{Lip}(\mathcal{K})$, then for every $l \in [d^*], f_l \in \operatorname{Lip}(\mathcal{K})$, combining the property that $A_{\gamma}(\cdot) \in \operatorname{Lip}(M)$ stated in Assumption 4.2, we can turn out $f(A_{\gamma}(\cdot))$ is $M\mathcal{K}$ -Lipschitz continuous.

So that

$$\|\mu_s(k) - \mu_t(k)\|_2 \le \sqrt{d^*} M \mathcal{K} \rho.$$

Lemma B.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

$$\begin{split} & \underset{\boldsymbol{x}\in C_{s}(k)}{\mathbb{E}} \underset{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}{\mathbb{E}} \|f(\boldsymbol{x}_{1})-\mu_{s}(k)\|_{2}^{2} \leq 4B_{2}^{2} \Big[\Big(1-\sigma_{s}+\frac{K\delta_{s}+2\varepsilon}{2B_{2}}+\frac{R_{s}(\varepsilon,f)}{p_{s}(k)}\Big)^{2} + \Big(1-\sigma_{s}+\frac{R_{s}(\varepsilon,f)}{p_{s}(k)}\Big) \Big],\\ & \text{where } R_{s}(\varepsilon,f) = P_{s} \Big(\boldsymbol{x}\in\bigcup_{k=1}^{K}C_{s}(k): \sup_{\boldsymbol{x}_{1},\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x})} \|f(\boldsymbol{x}_{1})-f(\boldsymbol{x}_{2})\|_{2} > \varepsilon \Big).\\ & \text{Proof. Let } S_{s}(\varepsilon,f) := \{\boldsymbol{x}\in\bigcup_{k=1}^{K}C_{s}(k): \sup_{\boldsymbol{x}_{1},\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x})} \|f(\boldsymbol{x}_{1})-f(\boldsymbol{x}_{2})\|_{2} \leq \varepsilon \}, \text{ for each } k\in[K], \end{split}$$

970
$$\mathbb{E}_{x \in C_{s}(k)} \mathbb{E}_{x_{1} \in \mathcal{A}(x)} \|f(x_{1}) - \mu_{s}(k)\|_{2}^{2}$$
971
$$= \frac{1}{p_{s}(k)} \mathbb{E}_{x} \mathbb{E}_{x_{1} \in \mathcal{A}(x)} [\mathbb{1}\{x \in C_{s}(k)\}\|f(x_{1}) - \mu_{s}(k)\|_{2}^{2}]$$

the second inequality is due to

$$P_s(C_s(k) \setminus ((\widetilde{C}_s(k) \cap S_s(\varepsilon, f)))) = P_s((C_s(k) \setminus \widetilde{C}_s(k)) \cup (C_s(k) \setminus S_s(\varepsilon, f))))$$

$$\leq (1 - \sigma_s) p_s(k) + R_s(\varepsilon, f).$$

Furthermore,

$$\begin{aligned}
\mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\|f(\boldsymbol{x}_{1})-\mu_{s}(k)\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\|f(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in C_{s}(k)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\|f(\boldsymbol{x}_{1})-\frac{P(\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f))}{p_{s}(k)}\mathbb{E}_{\boldsymbol{x}'\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2}) \\
&-\frac{P_{s}\left(C_{s}(k)\setminus(\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f))\right)}{p_{s}(k)}\mathbb{E}_{\boldsymbol{x}'\in C_{s}(k)\setminus(\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f))}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\left\|\frac{P_{s}(\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f))}{p_{s}(k)}\left(f(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right)\right\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\left\|\frac{P_{s}(\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f))}{p_{s}(k)}\left(f(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right)\right\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\left\|\frac{P(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in\mathcal{C}_{s}(k)\setminus(\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f))}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\left[\|f(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right)\right\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\left[\|f(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right)\right\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\left[\|f(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\left[\|f(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right\|_{2}^{2} \\
&= \mathbb{E}_{\boldsymbol{x}\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{1}\in\mathcal{A}(\boldsymbol{x})}\left[\|f(\boldsymbol{x}_{1})-\mathbb{E}_{\boldsymbol{x}'\in\tilde{C}_{s}(k)\cap S_{s}(\varepsilon,f)}\mathbb{E}_{\boldsymbol{x}_{2}\in\mathcal{A}(\boldsymbol{x}')}f(\boldsymbol{x}_{2})\right\|_{2}^{2} \\
&= \mathbb{E}_{$$

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

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

 $\|oldsymbol{x}_1-oldsymbol{x}_2\|_2$, we can turn out $\|oldsymbol{x}_1^*-oldsymbol{x}_2^*\|_2 \leq \delta_s,$ thus if we denote $({m x}_1^*, {m x}_2^*) =$ arg min $\boldsymbol{x}_1 {\in} \mathcal{A}(\boldsymbol{x}), \boldsymbol{x}_2 {\in} \mathcal{A}(\boldsymbol{x}')$ further by \mathcal{K} -Lipschitz continuity of f, we yield $||f(\boldsymbol{x}_1^*) - f(\boldsymbol{x}_2^*)||_2 \leq \mathcal{K} ||\boldsymbol{x}_1^* - \boldsymbol{x}_2^*||_2 \leq \mathcal{K} \delta_s$. In addition, since $x \in S_s(\varepsilon, f)$, we know for any $x_1 \in \mathcal{A}(x), \|f(x_1) - f(x_1^*)\|_2 \leq \varepsilon$. Similarly, $x' \in S_s(\varepsilon, f)$ implies $||f(x_2) - f(x_2^*)||_2 \le \varepsilon$ for any $x_2 \in \mathcal{A}(x')$. Therefore, for any $x, x' \in \mathcal{A}(x')$. $\widetilde{C}_s(1) \cap S_s(\varepsilon, f) \text{ and } \boldsymbol{x}_1 \in \mathcal{A}(\boldsymbol{x}), \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x}'),$ $\|f(\boldsymbol{x}_1) - f(\boldsymbol{x}_2)\|_2 \le \|f(\boldsymbol{x}_1) - f(\boldsymbol{x}_1^*)\|_2 + \|f(\boldsymbol{x}_1^*) - f(\boldsymbol{x}_2^*)\|_2 + \|f(\boldsymbol{x}_2^*) - f(\boldsymbol{x}_2)\|_2 \le 2\varepsilon + \mathcal{K}\delta_s.$ (22)Combining inequalities (20), (21), (22) to conclude

 $\mathbb{E}_{\boldsymbol{x}\in C_s(k)}\mathbb{E}_s$ $\|f(m)\|$ $(1_{2}) || 2$

$$k \mathbb{E}_{\boldsymbol{x}_1 \in \mathcal{A}(\boldsymbol{x})} \| f(\boldsymbol{x}_1) - \mu_s(k) \|_2^2$$

$$\leq \left[2\varepsilon + \mathcal{K}\delta_s + 2B_2\left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right)\right]^2 + 4B_2^2\left(1 - \sigma_s + \frac{R_s(\varepsilon, f)}{p_s(k)}\right)$$

B.2.3 THE EFFECT OF MINIMAXING OUR LOSS

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

 $=4B_2^2\Big[\Big(1-\sigma_s+\frac{\mathcal{K}\delta_s}{2B_2}+\frac{\varepsilon}{B_2}+\frac{R_s(\varepsilon,f)}{p_s(k)}\Big)^2+\Big(1-\sigma_s+\frac{R_s(\varepsilon,f)}{p_s(k)}\Big)\Big]$

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

1041 and

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \le \sqrt{\frac{2}{\min_{i \neq j} p_s(i) p_s(j)} \left(\mathcal{L}_{\operatorname{div}}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f) \right)} + 2\sqrt{d^*} B_2 M \mathcal{K} \rho.$$

 $\begin{array}{ll} \text{1046} & \text{where } R_s(\varepsilon, f) = P_s \left(\mathbf{x} \in \cup_{k=1}^K C_s(k) : \sup_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A}(\mathbf{x})} \| f(\mathbf{x}_1) - f(\mathbf{x}_2) \| > \varepsilon \right) \text{ and } \varphi(\sigma_s, \delta_s, \varepsilon, f) \\ \begin{array}{ll} \text{1047} \\ \text{1048} \end{array} & := 4B_2^2 \Big[\left(1 - \sigma_s + \frac{\kappa \delta_s + 2\varepsilon}{2B_2} \right)^2 + (1 - \sigma_s) + KR_s(\varepsilon, f) \left(3 - 2\sigma_s + \frac{\kappa \delta_s + 2\varepsilon}{B_2} \right) + R_s^2(\varepsilon, f) \left(\sum_{k=1}^K \frac{1}{p_s(k)} \right) \Big] + \\ \begin{array}{ll} \text{1049} \end{array} & B_2(\varepsilon^2 + 4B_2^2R_s(\varepsilon, f))^{\frac{1}{2}}. \end{array}$

Proof. Recall the Assumption 4.2, the measure on A is uniform, thus

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

1056 so that

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

1065
Denote $S := \{ \boldsymbol{z} : \mathbb{E}_{\boldsymbol{z}_1, \boldsymbol{z}_2 \in \mathcal{A}(\boldsymbol{z})} \| f(\boldsymbol{z}_1) - f(\boldsymbol{z}_2) \|_2 > \frac{\varepsilon}{m^2} \}$, by the definition of $R_t(\varepsilon, f)$ along with Markov inequality, we have $P^2(-f) = f(\boldsymbol{z}_2)$

$$R_{t}^{2}(\varepsilon, f) \leq P_{t}^{2}(S)$$

$$\leq \left(\frac{\mathbb{E}_{z}\mathbb{E}_{z_{1}, z_{2} \in \mathcal{A}(z)} \|f(z_{1}) - f(z_{2})\|_{2}}{\frac{\varepsilon}{m^{2}}}\right)^{2}$$

$$\leq \frac{\mathbb{E}_{z}\mathbb{E}_{z_{1}, z_{2} \in \mathcal{A}(z)} \|f(z_{1}) - f(z_{2})\|_{2}^{2}}{\frac{\varepsilon^{2}}{m^{4}}}$$

$$= \frac{m^{4}}{\varepsilon^{2}}\mathbb{E}_{z}\mathbb{E}_{z_{1}, z_{2} \in \mathcal{A}(z)} \|f(z_{1}) - f(z_{2})\|_{2}^{2}$$
(23)

Similar as above process, we can also get the first part stated in Lemma B.4:

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

Besides that, we can turn out

$$\mathbb{E}_{\boldsymbol{z}} \mathbb{E}_{\boldsymbol{z}_1, \boldsymbol{z}_2 \in \mathcal{A}(\boldsymbol{z})} \| f(\boldsymbol{z}_1) - f(\boldsymbol{z}_2) \|_2^2$$

$$\begin{aligned} & = \mathbb{E}_{\substack{\boldsymbol{x} \ \boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \mathcal{A}(\boldsymbol{x})}} \mathbb{E}_{\substack{\boldsymbol{x} \ \boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \mathcal{A}(\boldsymbol{x})}} \|f(\boldsymbol{x}_{1}) - f(\boldsymbol{x}_{2})\|_{2}^{2} + \mathbb{E}_{\substack{\boldsymbol{z} \ \boldsymbol{z}_{1}, \boldsymbol{z}_{2} \in \mathcal{A}(\boldsymbol{z})}} \mathbb{E}_{\substack{\boldsymbol{x} \ \boldsymbol{z}_{1}, \boldsymbol{z}_{2} \in \mathcal{A}(\boldsymbol{z})}} \|f(\boldsymbol{z}_{1}) - f(\boldsymbol{z}_{2})\|_{2}^{2} \\ & - \mathbb{E}_{\substack{\boldsymbol{x} \ \boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \mathcal{A}(\boldsymbol{x})}} \mathbb{E}_{\substack{\boldsymbol{x} \ \boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \mathcal{A}(\boldsymbol{z})}} \|f(\boldsymbol{x}_{1}) - f(\boldsymbol{x}_{2})\|_{2}^{2} \\ & - \mathbb{E}_{\substack{\boldsymbol{x} \ \boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \mathcal{A}(\boldsymbol{x})}} \mathbb{E}_{\substack{\boldsymbol{x} \ \boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \mathcal{A}(\boldsymbol{x})}} \|f(\boldsymbol{x}_{1}) - f(\boldsymbol{x}_{2})\|_{2}^{2} \\ & = \frac{1}{m^{2}} \sum_{\gamma=1}^{m} \sum_{\beta=1}^{m} \mathbb{E}_{\substack{\boldsymbol{x} \ \boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \mathcal{A}(\boldsymbol{x})} - f(\boldsymbol{A}_{\beta}(\boldsymbol{z}))\|_{2}^{2} - \mathbb{E}_{\boldsymbol{x}} \|f(\boldsymbol{A}_{\gamma}(\boldsymbol{x})) - f(\boldsymbol{A}_{\beta}(\boldsymbol{x}))\|_{2}^{2} \end{bmatrix} \end{aligned}$$

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

$$+ \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} \| f(\boldsymbol{x}_1) - f(\boldsymbol{x}_2) \|$$

 $+\mathbb{E}_{x}\mathbb{E}_{x_{1},x_{2}\in\mathcal{A}(x)}\|f(x_{1})-f(x_{2})\|_{2}^{2}$

since for all $\gamma \in [m], \beta \in [m]$ and $l \in [d^*]$, we have

$$\begin{aligned} & \begin{bmatrix} 1094\\ 1095\\ 1096\\ 1097\\ 1096\\ 1097\\ 1098\\ 1098\\ 1099\\ 1009\\ 1009\\ 1009\\ 1000\\ 100$$

It is necessary to claim $g(\mathbf{x}) \in \text{Lip}(8B_2M\mathcal{K})$ at first to obtain the last inequality shown above. In fact, $\forall l \in [d^*], f_l \in \operatorname{Lip}(\mathcal{K})$ as $f \in \operatorname{Lip}(\mathcal{K})$, and review that $A_{\gamma}(\cdot)$ and $A_{\beta}(\cdot)$ are both M-Lipschitz continuous according to Assumption 4.2, therefore we can turn out $f_l(A_{\gamma}(\cdot)) - f_l(A_{\beta}(\cdot)) \in$ $\operatorname{Lip}(2M\mathcal{K})$. In addition, note that $|f_l(A_{\gamma}(\cdot)) - f_l(A_{\beta}(\cdot))| \leq 2B_2$ as $||f||_2 \leq B_2$, hence the out-ermost quadratic function remains locally $4B_2$ -Lipschitz continuity in $[-2B_2, 2B_2]$, which implies that $g \in \operatorname{Lip}(8B_2M\mathcal{K})$.

Now let's separately derive the two terms of the last inequality, combine the conclusion that $g \in$ $\operatorname{Lip}(8B_2M\mathcal{K})$, the definition of Wasserstein distance and Assumption B.1 can obtain

$$\sum_{k=1}^{K} \left[p_t(k) \Big(\mathbb{E}_{\boldsymbol{z} \in C_t(k)} \big[f_l(A_{\gamma}(\boldsymbol{z})) - f_l(A_{\beta}(\boldsymbol{z})) \big]^2 - \mathbb{E}_{\boldsymbol{x} \in C_s(k)} \big[f_l(A_{\gamma}(\boldsymbol{x})) - f_l(A_{\beta}(\boldsymbol{x})) \big]^2 \Big) \right]$$

$$\leq 8B_2 M \mathcal{K} \rho \sum_{k=1}^{K} p_t(k)$$

$$= 8B_2 M \mathcal{K} \rho,$$

For the second term in the last inequality, just need to notice that $f_l(A_{\gamma}(\boldsymbol{x})) - f_l(A_{\beta}(\boldsymbol{x})) \leq 2B_2$, and then apply Assumption **B**.1 to yield

$$\sum_{k=1}^{K} \left[\left(p_t(k) - p_s(k) \right) \mathbb{E}_{\boldsymbol{x} \in C_s(k)} \left[f_l(A_{\gamma}(\boldsymbol{x})) - f_l(A_{\beta}(\boldsymbol{x})) \right]^2 \right] \le 4B_2^2 K \eta.$$

Hence we have

1129
$$\mathbb{E}_{\boldsymbol{z}} \mathbb{E}_{\boldsymbol{z}_1, \boldsymbol{z}_2 \in \mathcal{A}(\boldsymbol{z})} \| f(\boldsymbol{z}_1) - f(\boldsymbol{z}_2) \|_2^2 \le \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} \| f(\boldsymbol{x}_1) - f(\boldsymbol{x}_2) \|_2^2 + 8B_2 d^* M \mathcal{K} \rho + 4B_2^2 d^* \mathcal{K} \eta$$
1130 (24)

Combining (23) and (24) turn out the second inequality of Lemma B.4.

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

To prove the third part of this Lemma, first recall Lemma B.2 that $\forall k \in [K]$, $\forall k \in [K]$,

$$\|\mu_s(k) - \mu_t(k)\|_2 \le \sqrt{d^* M \mathcal{K} \rho}.$$

1137 Hence, $\forall i \neq j$, we have

$$\begin{aligned} |\mu_t(i)^\top \mu_t(j) - \mu_s(i)^\top \mu_s(j)| &= |\mu_t(i)^\top \mu_t(j) - \mu_t(i)^\top \mu_s(j) + \mu_t(i)^\top \mu_s(j) - \mu_s(i)^\top \mu_s(j)| \\ &\leq \|\mu_t(i)\|_2 \|\mu_t(j) - \mu_s(j)\|_2 + \|\mu_s(j)\|_2 \|\mu_t(i) - \mu_s(i)\|_2 \\ &\leq 2\sqrt{d^*} B_2 M \mathcal{K} \rho, \end{aligned}$$

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

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \le \max_{i \neq j} |\mu_s(i)^\top \mu_s(j)| + 2\sqrt{d^*} B_2 M \mathcal{K} \rho.$$
(25)

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

1163 Therefore,

$$\begin{aligned}
& 1164 \\
& 1165 \\
& 1166 \\
& 1166 \\
& 1167 \\
& 1168 \\
& 1169 \\
& 1169 \\
& 1169 \\
& 1169 \\
& 1169 \\
& 1169 \\
& 1170 \\
& 1170 \\
& 1171 \\
& 1172 \\
& & \leq \frac{2 \left\| \sum_{\boldsymbol{x} \, \boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*} + \sum_{k=1}^{K} p_s(k) \mu_s(k) \mu_s(k)^\top - \sum_{\boldsymbol{x} \, \boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] \right\|_F^2}{p_s(i) p_s(j)} \\
& & 1171 \\
& & & \frac{2 \left\| \sum_{\boldsymbol{x} \, \boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*} \right\|_F^2 + 2 \left\| \sum_{k=1}^{K} p_s(k) \mu_s(k) \mu_s(k)^\top - \sum_{\boldsymbol{x} \, \boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] \right\|_F^2}{p_s(i) p_s(j)} \end{aligned}$$
(26)

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

$$\mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top$$
$$= \sum_{k=1}^K p_s(k) \mathbb{E}_{\boldsymbol{x} \in C_s(k)} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top$$
$$= \sum_{k=1}^K p_s(k) \mathbb{E}_{\boldsymbol{x} \in C_s(k)} \mathbb{E}_{\boldsymbol{x}_1 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_1)^\top] - \sum_{k=1}^K p_s(k) \mu_s(k) \mu_s(k)^\top$$

1185
1186
1187

$$\sum_{k=1}^{K} p_{s}(k) = x_{1} \in \mathcal{A}(x) [f(x_{1})(f(x_{2}) - f(x_{1}))^{\top}]$$

$$\sum_{k=1}^{K} p_{s}(k) = x_{1} \in \mathcal{A}(x) [f(x_{1})(f(x_{2}) - f(x_{1}))^{\top}]$$

1192 1193

1194

1202

1203

where the last equation is derived from

$$\begin{split} & \begin{array}{l} & \begin{array}{l} & 1195 \\ & 1196 \\ & 1196 \\ & 1197 \\ & 1197 \\ & 1197 \\ & 1198 \\ & 1199 \\ & 1199 \\ & 1199 \\ & 1199 \\ & 1199 \\ & 1190 \\ & 1190 \\ & 1190 \\ & 1190 \\ & 1190 \\ & 1190 \\ & 1190 \\ & 1100 \\ &$$

 $=\sum_{k=1}^{K} p_s(k) \mathbb{E}_{\boldsymbol{x} \in C_s(k)} \mathbb{E}_{\boldsymbol{x}_1 \in \mathcal{A}(\boldsymbol{x})} [(f(\boldsymbol{x}_1) - \mu_s(k))(f(\boldsymbol{x}_1) - \mu_s(k))^\top]$

(27)

+ $\mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1)(f(\boldsymbol{x}_2) - f(\boldsymbol{x}_1))^\top],$

So its norm is

1242
1243 If we define
$$\varphi(\sigma_s, \delta_s, \varepsilon, f) := 4B_2^2 \Big[\Big(1 - \sigma_s + \frac{\kappa \delta_s + 2\varepsilon}{2B_2} \Big)^2 + (1 - \sigma_s) + KR_s(\varepsilon, f) \Big(3 - 2\sigma_s + \frac{\kappa \delta_s + 2\varepsilon}{2B_2} \Big) + B_2^2(\varepsilon, f) \Big(\sum_{s=1}^K \frac{1}{2\pi s} \Big) \Big] + B_2(\varepsilon^2 + 4B_2^2R_s(\varepsilon, f))^{\frac{1}{2}}$$
 above derivation implies

 $\overline{B_2}$ $\int + \kappa_s^-(\varepsilon, f) \left(\sum_{k=1} \frac{1}{p_s(k)} \right) \right] + B_2(\varepsilon^2 + 4B_2^2 R_s(\varepsilon, f))^{\frac{1}{2}}$, above derivation implies Kш

$$\left\|\sum_{k=1} p_s(k)\mu_s(k)\mu_s(k)^{\top} - \mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1,\boldsymbol{x}_2\in\mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}_1)f(\boldsymbol{x}_2)^{\top}]\right\|_F \le \varphi(\sigma_s,\delta_s,\varepsilon,f).$$
(28)

Besides that, Note that

$$\mathcal{L}_{\text{div}}(f) = \sup_{G \in \mathcal{G}(f)} \langle \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*}, G \rangle_F$$
$$= \left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*} \right\|_F^2,$$
(29)

which is from the facts that $\mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1,\boldsymbol{x}_2\in\mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}_1)f(\boldsymbol{x}_2)^{\top}] - I_{d^*}\in\mathcal{G}(f)$ and

$$\langle \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*}, G \rangle_F \leq \left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*} \right\|_F \cdot \|G\|_F$$

Combining (26), (27), (28), (29) yields for any $i \neq j$

$$(\mu_s(i)^{\top}\mu_s(j))^2 \le \frac{2}{p_s(i)p_s(j)} \Big(\mathcal{L}_{\operatorname{div}}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f) \Big),$$

which implies that

$$\max_{i \neq j} |\mu_s(i)^\top \mu_s(j)| \le \sqrt{\frac{2}{\min_{i \neq j} p_s(i) p_s(j)} \Big(\mathcal{L}_{\operatorname{div}}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f) \Big)}.$$

So we can get what we desired according to (25)

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \le \sqrt{\frac{2}{\min_{i \neq j} p_s(i) p_s(j)} \left(\mathcal{L}_{\text{div}}(f) + \varphi(\sigma_s, \delta_s, \varepsilon, f) \right)} + 2\sqrt{d^*} B_2 M \mathcal{K} \rho.$$

B.2.4 CONNECTION BETWEEN PRETRAINING AND DOWNSTREAM TASK

Following theorem reveals that minimaxing our loss may achieve a small misclassification rate in downstream task.

1278
1279
1279
1280
1281
1282
we have
Theorem B.1. Given a
$$(\sigma_s, \sigma_t, \delta_s, \delta_t)$$
-augmentation, for any $\varepsilon > 0$, if $\Theta(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s}) > 0$, then
 $(\int_{i \neq j}^{2} \frac{2}{i \neq j^{n_s(i)p_s(j)}} \left(\frac{1}{\lambda} \mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(f_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + \psi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}) \right) + 2\sqrt{d^*} B_2 M \mathcal{K} \rho$
 $B_2^2 \Theta(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})$,

$$\mathbb{E}_{D_s}[\operatorname{Err}(Q_{\hat{f}_{n_s}})] \le (1 - \sigma_t) + \frac{m^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + 8B_2 d^* M \mathcal{K}\rho + 4B_2^2 d^* K \eta}$$

where

$$\psi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}) := B_2 \Big(\varepsilon^2 + 4B_2^2 \frac{m^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s} [\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)]} \Big)^{\frac{1}{2}} + 4B_2^2 \Big[\Big(1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2} \Big)^2 \Big]$$

$$+ (1 - \sigma_s) + \frac{Km^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)]} \left(3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2}\right)$$

$$+ (1 - \sigma_s) + \frac{Km^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)]} \left(3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2}\right)$$

1294
1295
$$+ \frac{m^4}{\varepsilon^2} \mathbb{E}_{D_s} [\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] \Big(\sum_{k=1}^K \frac{1}{p_s(k)} \Big) \Big],$$

1296
1297
$$\Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, f) = \left(\sigma_t - \frac{R_t(\varepsilon, f)}{\min_i p_t(i)}\right) \left(1 + \left(\frac{B_1}{B_2}\right)^2 - \frac{\mathcal{K}\delta_t}{B_2} - \frac{2\varepsilon}{B_2}\right) - 1, \ \Delta_{\hat{\mu}_t} = 1 - \frac{\min_{k \in [K]} \|\hat{\mu}_t(k)\|^2}{B_2^2}$$

 $R_t(\varepsilon, f) = P_t \left(\boldsymbol{z} \in \bigcup_{k=1}^K \widetilde{C}_t(k) : \sup_{\boldsymbol{z}_1, \boldsymbol{z}_2 \in \mathcal{A}(\boldsymbol{z})} \| f(\boldsymbol{z}_1) - f(\boldsymbol{z}_2) \| > \varepsilon \right) \text{ and } \Theta(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s}) = C_s \left(\sum_{k=1}^K \widetilde{C}_t(k) : \sum_{k=1}^K \widetilde{C}_t$ $\Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s}) - \sqrt{2 - 2\Gamma_{\min}(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})} - \frac{\Delta_{\hat{\mu}_t}}{2} - \frac{2\max_{k \in [K]} \|\hat{\mu}_t(k) - \mu_t(k)\|_2}{B_2}.$

In addition, the following inequalities always hold

$$\mathbb{E}_{D_s}[R_t^2(\varepsilon, \hat{f}_{n_s})] \le \frac{m^4}{\varepsilon^2} \Big(\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + 8B_2 d^* M \mathcal{K}\rho + 4B_2^2 d^* \mathcal{K}\eta \Big)$$

Proof. Note the facts that $\sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) \geq \max\{\mathcal{L}_{\operatorname{align}}(f), \lambda \mathcal{L}_{\operatorname{div}}(f)\}, B_1 \leq \|\hat{f}_{n_s}\|_2 \leq B_2$ and \mathcal{K} -Lipschitz continuity of \hat{f}_{n_s} , apply Lemma B.4 to \hat{f}_{n_s} to obtain

$$R_s^2(\varepsilon, \hat{f}_{n_s}) \le \frac{m^4}{\varepsilon^2} \sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s})$$
(30)

$$R_t^2(\varepsilon, \hat{f}_{n_s}) \le \frac{m^4}{\varepsilon^2} \sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G) + \frac{8m^4}{\varepsilon^2} B_2 d^* M \mathcal{K}\rho + \frac{4m^4}{\varepsilon^2} B_2^2 d^* K \eta$$
(31)

and

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| \le \sqrt{\frac{2}{\min_{i \neq j} p_s(i) p_s(j)} \left(\frac{1}{\lambda} \sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G) + \varphi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})\right)} + 2\sqrt{d^*} B_2 M \mathcal{K} \rho$$
(32)

Take expectation w.r.t D_s in the both side of (30), (31), (32) and apply Jensen inequality to yield

$$\mathbb{E}_{D_s}[R_s^2(\varepsilon, \hat{f}_{n_s})] \le \frac{m^4}{\varepsilon^2} \mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)]$$

$$\mathbb{E}_{D_s}[R_t^2(\varepsilon, \hat{f}_{n_s})] \leq \frac{m^4}{\varepsilon^2} \mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + \frac{8m^4}{\varepsilon^2} B_2 d^* M \mathcal{K}\rho + \frac{4m^4}{\varepsilon^2} B_2^2 d^* K \eta$$

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

where $\mathbb{E}_{D_s}[\varphi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})] = 4B_2^2 \Big[\Big(1 - \sigma_s + \frac{\kappa \delta_s + 2\varepsilon}{2B_2}\Big)^2 + (1 - \sigma_s) + K \mathbb{E}_{D_s}[R_s(\varepsilon, \hat{f}_{n_s})] \Big(3 - \frac{\kappa \delta_s + 2\varepsilon}{2B_2}\Big)^2 \Big]$ $2\sigma_s + \frac{\kappa\delta_s + 2\varepsilon}{B_2} + \mathbb{E}_{D_s}[R_s^2(\varepsilon, \hat{f}_{n_s})] \Big(\sum_{k=1}^{\kappa} \frac{1}{p_s(k)}\Big) + B_2 \mathbb{E}_{D_s}[(\varepsilon^2 + 4B_2^2 R_s(\varepsilon, \hat{f}_{n_s}))^{\frac{1}{2}}].$

Therefore, by Jensen inequality, we have

$$\begin{split} & \mathbb{E}_{D_s}[\varphi(\sigma_s, \delta_s, \varepsilon, R_s(\varepsilon, \hat{f}_{n_s}))] \\ & \leq 4B_2^2 \Big[\Big(1 - \sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{2B_2} \Big)^2 + (1 - \sigma_s) + K\mathbb{E}_{D_s}[R_s(\varepsilon, \hat{f}_{n_s})] \Big(3 - 2\sigma_s + \frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2} \Big) \\ & + \mathbb{E}_{D_s}[R_s^2(\varepsilon, \hat{f}_{n_s})] \Big(\sum_{k=1}^K \frac{1}{p_s(k)} \Big) \Big] + B_2(\varepsilon^2 + 4B_2^2 \mathbb{E}_{D_s}[R_s(\varepsilon, \hat{f}_{n_s})])^{\frac{1}{2}} \end{split}$$

$$\frac{\mathcal{K}\delta_s + 2\varepsilon}{B_2} + \frac{m^4}{\varepsilon^2} \mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] \Big(\sum_{k=1}^K \frac{1}{p_s(k)}\Big)\Big]$$

(D) 2

$$+ B_2 \left(\varepsilon^2 + \frac{4B_2^2 m^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s} [\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)]} \right)^{\frac{1}{2}} \\ := \psi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}).$$

Recall Lemma B.1 reveals that we can obtain

$$\operatorname{Err}(Q_{\hat{f}_{n_s}}) \leq (1 - \sigma_t) + R_t(\varepsilon, \hat{f}_{n_s})$$

if $\max_{i \neq j} |(\mu_t(i))^\top \mu_t(j)| < B_2^2 \Theta(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s}).$

So that if $\Theta(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s}) > 0$, apply Markov inequality to know with probability at least

$$1 - \frac{\sqrt{\frac{2}{\min_{i \neq j} p_s(i) p_s(j)} \left(\frac{1}{\lambda} \mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + \psi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})\right)}{B_2^2 \Theta(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})},$$

we have

$$\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)| < B_2^2 \Theta(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})$$

so that we can get what we desired.

$$\mathbb{E}_{D_s}[\operatorname{Err}(Q_{\hat{f}_{n_s}})] \leq (1 - \sigma_t) + R_t(\varepsilon, \hat{f}_{n_s})$$
$$\leq (1 - \sigma_t) + \frac{m^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + 8B_2 d^* M \mathcal{K} \rho + 4B_2^2 d^* K \eta}$$

where the last inequality is due to (31).

B.2.5 PRELIMINARIES FOR ERROR ANALYSIS

To prove Theorem 4.2 based on Theorem B.1, we need to first introduce some related definitions and conclusions, which are going to be used in subsequent contents.

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

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

furthermore, denote $\mathcal{G}_1 := \{ G \in \mathbb{R}^{d^* \times d^*} : \|G\|_F \le B_2^2 + \sqrt{d^*} \}$. It is obvious that both $\mathcal{G}(f)$ for any $f: ||f||_2 \leq B_2$ and $\widehat{\mathcal{G}}(f)$ for any $f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K},B_1,B_2)$ are contained in \mathcal{G}_1 . Apart from that, following Proposition B.1 reveals that $\ell(\boldsymbol{u}, G)$ is a Lipschitz function on the domain $\{\boldsymbol{u} \in \mathbb{R}^{2d^*} : \|\boldsymbol{u}\|_2 \leq \sqrt{2}B_2\} \times \mathcal{G}_1 \subseteq \mathbb{R}^{2d^*+(d^*)^2}.$

Proposition B.1. ℓ is a Lipschitz function on the domain $\{\mathbf{u} \in \mathbb{R}^{2d^*} : \|\mathbf{u}\|_2 \leq \sqrt{2}B_2\} \times \mathcal{G}_1$.

Proof. At first step, we will prove $\|\ell(\cdot,G)\|_{\text{Lip}} < \infty$ for any fixed $G \in \mathcal{G}_1$. To this end, denote $u = (u_1, u_2)$, where $u_1, u_2 \in \mathbb{R}^{d^*}$, we firstly show $J(u) = ||u_1 - u_2||_2^2$ is Lipschtiz function. let $g(u) := u_1 - u_2$, then

$$\begin{split} \|g(\boldsymbol{u}_1, \boldsymbol{u}_2) - g(\boldsymbol{v}_1, \boldsymbol{v}_2)\|_2^2 &= \|\boldsymbol{u}_1 - \boldsymbol{u}_2 - \boldsymbol{v}_1 + \boldsymbol{v}_2\|_2^2 \\ &\leq \left(\|\boldsymbol{u}_1 - \boldsymbol{v}_1\|_2 + \|\boldsymbol{u}_2 - \boldsymbol{v}_2\|_2\right)^2 \\ &= \|\boldsymbol{u}_1 - \boldsymbol{v}_1\|_2^2 + \|\boldsymbol{u}_2 - \boldsymbol{v}_2\|_2^2 + 2\|\boldsymbol{u}_1 - \boldsymbol{v}_1\|_2\|\boldsymbol{u}_2 - \boldsymbol{v}_2\|_2 \\ &\leq 2(\|\boldsymbol{u}_1 - \boldsymbol{v}_1\|_2^2 + \|\boldsymbol{u}_2 - \boldsymbol{v}_2\|_2^2) \\ &= 2\|(\boldsymbol{u}_1, \boldsymbol{u}_2) - (\boldsymbol{v}_1, \boldsymbol{v}_2)\|_2^2, \end{split}$$

1404 which implies that $q(u) \in \operatorname{Lip}(\sqrt{2})$. Apart from that, q also possess the property that $||q(u)||_2 =$ 1405 $\|\boldsymbol{u}_1 - \boldsymbol{u}_2\|_2 \le \|\boldsymbol{u}_1\|_2 + \|\boldsymbol{u}_2\|_2 \le 2\|\boldsymbol{u}\|_2 \le 2\sqrt{2}B_2$. Moreover, let $h(\boldsymbol{v}) := \|\boldsymbol{v}\|_2^2$, we know that 1406 $\left\|\frac{\partial h}{\partial \boldsymbol{u}}(g(\boldsymbol{u}))\right\|_{2} = 2\|g(\boldsymbol{u})\|_{2} \le 4\sqrt{2}B_{2}.$ 1407 1408 Therefore, $J(u) = h(g(u)) = ||u_1 - u_2||_2^2 \in Lip(8B_2)$ 1409 1410 To show $Q(\boldsymbol{u}) = \langle \boldsymbol{u}_1 \boldsymbol{u}_2^\top - I_{d^*}, G \rangle_F$ is also a Lipschtiz function. Define $\tilde{g}(\boldsymbol{u}) := \boldsymbol{u}_1 \boldsymbol{u}_2^\top$, we know 1411 that 1412 $\|\tilde{g}(\boldsymbol{u}) - \tilde{g}(\boldsymbol{v})\|_{F} = \|\boldsymbol{u}_{1}\boldsymbol{u}_{2}^{\top} - \boldsymbol{v}_{1}\boldsymbol{v}_{2}^{\top}\|_{F}$ 1413 $= \| \boldsymbol{u}_1 \boldsymbol{u}_2^\top - \boldsymbol{u}_1 \boldsymbol{v}_2^\top + \boldsymbol{u}_1 \boldsymbol{v}_2^\top - \boldsymbol{v}_1 \boldsymbol{v}_2^\top \|_F$ 1414 1415 $= \| \boldsymbol{u}_1 (\boldsymbol{u}_2 - \boldsymbol{v}_2)^\top + (\boldsymbol{u}_1 - \boldsymbol{v}_1) \boldsymbol{v}_2^\top \|_F$ 1416 $\leq \|m{u}_1\|_F \|m{u}_2 - m{v}_2\|_F + \|m{u}_1 - m{v}_1\|_F \|m{v}_2\|_F$ 1417 $\leq (\|m{u}_1\|_2 + \|m{v}_2\|_2)\|m{u} - m{v}\|_2$ 1418 $< 2\sqrt{2}B_2 \|\boldsymbol{u} - \boldsymbol{v}\|_2.$ 1419 1420 Furthermore, denote $\tilde{h}(A) := \langle A - I_{d^*}, G \rangle_F$, then $\|\nabla \tilde{h}(A)\|_F = \|G\|_F \leq B_2^2 + \sqrt{d^*}$. So that 1421 $Q(\boldsymbol{u}) = \tilde{h}(\tilde{g}(\boldsymbol{u})) \in \operatorname{Lip}(2\sqrt{2}B_2(B_2^2 + \sqrt{d^*})).$ 1422 Combining above conclusions knows that for any $G \in \mathcal{G}_1$, we have $\|\ell(\cdot, G)\|_{\text{Lip}} < \infty$ on the domain 1423 1424 $\{ \boldsymbol{u} : \| \boldsymbol{u} \|_2 \leq \sqrt{2}B_2 \}.$ 1425 Next, fixed $\boldsymbol{u} \in \mathbb{R}^{2d^*}$ such that $\|\boldsymbol{u}\|_2 \leq \sqrt{2}B_2$, we have 1426 1427 $|\ell(\boldsymbol{u}, G_1) - \ell(\boldsymbol{u}, G_2)| = |\langle \boldsymbol{u}, G_1 - G_2 \rangle_F| \le \|\boldsymbol{u}\|_2 \|G_1 - G_2\|_F = \sqrt{2}B_2 \|G_1 - G_2\|_F,$ 1428 which implies that $\ell(\boldsymbol{u}, \cdot) \in \operatorname{Lip}(\sqrt{2}B_2)$. 1429 1430 Finally, note that 1431 $|\ell(\boldsymbol{u}_1, G_1) - \ell(\boldsymbol{u}_2, G_2)|^2 \le (|\ell(\boldsymbol{u}_1, G_1) - \ell(\boldsymbol{u}_2, G_1)| + |\ell(\boldsymbol{u}_2, G_1) - \ell(\boldsymbol{u}_2, G_2)|)^2$ 1432 $\leq \left(\left(\sqrt{2} + 2\sqrt{2}B_2(B_2^2 + \sqrt{d^*})\right)\|\boldsymbol{u}_1 - \boldsymbol{u}_2\|_2 + \sqrt{2}B_2\|G_1 - G_2\|_F\right)^2$ 1433 1434 $< 2(\sqrt{2} + 2\sqrt{2}B_2(B_2^2 + \sqrt{d^*}))^2 \|u_1 - u_2\|_2^2 + 4B_2^2 \|G_1 - G_2\|_F^2$ 1435 1436 $< C \| \operatorname{vec}(\boldsymbol{u}_1, G_1) - \operatorname{vec}(\boldsymbol{u}_2, G_2) \|_2^2$ 1437 where C is a constant s.t $C \ge \max\{2(\sqrt{2} + 2\sqrt{2}B_2(B_2^2 + \sqrt{d^*}))^2, 4B_2^2\}$, which yields what we 1438 desired. 1439 1440 We summary the Lipschitz constants of $\ell(u,G)$ with respect to both $u \in \{u \in \mathbb{R}^{2d^*} : ||u||_2 \le$ 1441 $\sqrt{2}B_2$ and $G \in \mathcal{G}_1$ in Table 2. 1442 1443 1444 Table 2: Lipschitz constant of ℓ with respect to each component 1445 **Lipschitz Constant** Function 1446 $\frac{\sqrt{2}B_2}{2\sqrt{2}B_2(B_2^2+\sqrt{d^*})}$ 1447 $\ell(\boldsymbol{u},\cdot)$ 1448 $\ell(\cdot, G)$ $\max\left\{\sqrt{2}B_2, 2\sqrt{2}B_2(B_2^2 + \sqrt{d^*})\right\}$ 1449 $\ell(\cdot)$ 1450 1451

Definition B.1 (Rademacher complexity). Given a set $S \subseteq \mathbb{R}^n$, the Rademacher complexity of S is denoted by

1454

1455 1456

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

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 1/2.

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

Lemma B.5 (Vector-contraction principle). Let \mathcal{X} be any set, $(x_1, \ldots, x_n) \in \mathcal{X}^n$, let F be a class of functions $f : \mathcal{X} \to \ell_2$ and let $h_i : \ell_2 \to \mathbb{R}$ have Lipschitz norm L. Then

$$\mathbb{E}\sup_{f\in F} \left|\sum_{i} \epsilon_{i} h_{i}(f(x_{i}))\right| \leq 2\sqrt{2}L\mathbb{E}\sup_{f\in F} \left|\sum_{i,k} \varepsilon_{ik} f_{k}(x_{i})\right|,$$

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

Proof. Combining Maurer (2016) and Theorem 3.2.1 of Giné & Nickl (2016) obtains the desired result.

1470

1467

1462

1463 1464

1471 Recall $\mathcal{NN}_{d_1,d_2}(W,L,\mathcal{K}) := \{\phi_\theta(x) = A_L \sigma(A_{L-1}\sigma(\cdots \sigma(A_0x)) : \kappa(\theta) \leq \mathcal{K}\}$, which is de-1472 fined in (13). The second lemma we will employed is related to the upper bound for Rademacher 1473 complexity of hypothesis space consisting of norm-constrained neural networks, which was pro-1474 vided by Golowich et al. (2018).

1475 Lemma B.6 (Theorem 3.2 of Golowich et al. (2018)). $\forall n \in \mathbb{N}^+, \forall x_1, \dots, x_n \in [-B, B]^d$ with 1476 $B \ge 1, S := \{(\phi(x_1), \dots, \phi(x_n)) : \phi \in \mathcal{NN}_{d,1}(W, L, \mathcal{K})\} \subseteq \mathbb{R}^n$, then

1477 1478

1479 1480

1486 1487

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

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

Definition B.2 (Covering number). $\forall n \in \mathbb{N}^+$, Fix $S \subseteq \mathbb{R}^n$ and $\varrho > 0$, the set \mathcal{N} is called an ϱ -net of S with respect to a norm $\|\cdot\|$ on \mathbb{R}^n , if $\mathcal{N} \subseteq S$ and for any $u \in S$, there exists $v \in \mathcal{N}$ such that $\|u - v\| \le \varrho$. The covering number of S is defined as

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

where |Q| is the cardinality of the set Q.

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

Apart from that, we need to employ following finite maximum inequality, which is stated in Lemma
2.3.4 of Giné & Nickl (2016), in later deduction.

Lemma B.7 (Finite maximum inequality). For any $N \ge 1$, if $X_i, i \le N$, are sub-Gaussian random variables admitting constants σ_i , then

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

1501 1502 **Definition B.3** (Excess risk). The difference between $\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)$ and $\sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G)$ is 1503 called excess risk i.e.

¹⁵⁰³ called excess risk, i.e., 1504

$$\mathcal{E}(\hat{f}_{n_s}) = \sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G) - \sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G).$$

1506 1507

1499

1500

1508 B.2.6 DEAL WITH $\sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*)$

We aim to claim $\sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*) = 0$ in two step. At first, we assert that if there exists a measurable map f satisfying $\Sigma = \mathbb{E}_{\boldsymbol{x} \sim P_s}[f(\boldsymbol{x})f(\boldsymbol{x})^{\top}]$ be positive definite, then we can conduct some minor rectification on it to get \tilde{f} such that $\sup_{G \in \mathcal{G}(\tilde{f})} \mathcal{L}(\tilde{f}) = 0$. At the second step, we are going $\left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*} \right\|_F$

to show the required f does exist under Assumption 4.5 and the rectification \tilde{f} also fulfill the requirement that $B_1 \leq \|\tilde{f}\|_2 \leq B_2$, which implies that $\sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*) = 0$ as the definition of f^* implies $\sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*) \leq \sup_{G \in \mathcal{G}(\tilde{f})} \mathcal{L}(\tilde{f})$.

1516 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 $f : B_1 \leq ||f||_2 \leq B_2$ satisfying both $\mathcal{L}_{\text{align}}(f) = 0$ and $\|\mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1,\boldsymbol{x}_2\in\mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}_1)f(\boldsymbol{x}_2)^{\top}] - I_{d^*}\|_F = 0$. Note that

1522 1523 1524

1527 1528

1529

1531

1534 1535 1536

1540 1541 1542

1546 1547

1549 1550

1552

Above deduction tells us that finding a measurable map $f: B_1 \leq ||f||_2 \leq B_2$ making both $\mathcal{L}_{align}(f)$ and $\left\|\mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}'\in\mathcal{A}(\boldsymbol{x})}[f(\boldsymbol{x}')f(\boldsymbol{x}')^{\top}] - I_{d^*}\right\|_F$ vanished is just enough to achieve our goal.

 $\leq \left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_1)^\top] - I_{d^*} \right\|_{F} + \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2} [\|f(\boldsymbol{x}_1)\|_2 \|f(\boldsymbol{x}_1) - f(\boldsymbol{x}_2)\|_2]$

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

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

Lemma B.8. If there exists a measurable map f making $\Sigma = \mathbb{E}_{\boldsymbol{x} \sim P_s}[f(\boldsymbol{x})f(\boldsymbol{x})^\top]$ positive definite, then there exists a measurable map \tilde{f} making both

$$\mathcal{L}_{\text{align}}(\tilde{f}) = 0 \text{ and } \|\mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}' \in \mathcal{A}(\boldsymbol{x})}[\tilde{f}(\boldsymbol{x}')\tilde{f}(\boldsymbol{x}')^{\top}] - I_{d^*}\|_F = 0.$$

1537 1538 *Proof.* We conduct following revision for given f to obtain \tilde{f} .

1539 For any $x \in \mathcal{X}$, define

$$ilde{f}_{m{x}}(m{x}') = egin{cases} V^{-1}f(m{x}) & ext{ if }m{x}' \in \mathcal{A}(m{x}) \ f(m{x}) & ext{ if }m{x}'
ot\in \mathcal{A}(m{x}) \end{cases}$$

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

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

1548 and

$$orall oldsymbol{x} \in \mathcal{X}, oldsymbol{x}_1, oldsymbol{x}_2 \in \mathcal{A}(oldsymbol{x}), \| ilde{f}(oldsymbol{x}_1) - ilde{f}(oldsymbol{x}_2)\|_2 = \| ilde{f}(oldsymbol{x}) - ilde{f}(oldsymbol{x})\|_2 = 0.$$

1551 That is what we desired.

1553 Remark B.2. If we have a measurable partition $\mathcal{X} = \bigcup_{i=1}^{d^*} \mathcal{P}_i$ stated in Assumption 4.5 such 1554 that $\mathcal{P}_i \cap \mathcal{P}_j = \emptyset$ and $\forall i \in [d^*], \frac{1}{B_2^2} \leq P_s(\mathcal{P}_i) \leq \frac{1}{B_1^2}$, just set the $f(\boldsymbol{x}) = \boldsymbol{e}_i$ if $\boldsymbol{x} \in \mathcal{P}_i$, 1555 where \boldsymbol{e}_i is the standard basis of \mathbb{R}^{d^*} , then $\Sigma = \text{diag}\{P_s(\mathcal{P}_1), \dots, P_s(\mathcal{P}_i), \dots, P_s(\mathcal{P}_{d^*})\}, V^{-1} =$ 1556 $\text{diag}\{\sqrt{\frac{1}{P_s(\mathcal{P}_1)}}, \dots, \sqrt{\frac{1}{P_s(\mathcal{P}_i)}}, \dots, \sqrt{\frac{1}{P_s(\mathcal{P}_{d^*})}}\}, \tilde{f}(\boldsymbol{x}) = \sqrt{\frac{1}{P_s(\mathcal{P}_i)}}\boldsymbol{e}_i$ if $\boldsymbol{x} \in \mathcal{P}_i$, it is obviously that 1558 $B_1 \leq \|\tilde{f}\|_2 \leq B_2$.

1559

1560 B.2.7 RISK DECOMPOSITION 1561

1562 If denote
$$\widehat{G}(f) = \frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - I_{d^*}$$
 and $G^*(f) = \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*}$,

use can decompose $\mathcal{E}(f_{n_s})$ into three terms shown as follow and then deal each term successively. To achieve conciseness in subsequent conclusions, we employ $X \leq Y$ or $Y \geq X$ to indicate the statement that $X \leq CY$ form some C > 0 if X and Y are two quantities.

1566 **Lemma B.9.** The excess risk $\mathcal{E}(\hat{f}_{n_n})$ satisfies 1567 $\mathcal{E}(\widehat{f}_{n_s}) \leq 2 \sup_{\substack{f \in \mathcal{F}, G \in \widehat{\mathcal{G}}(f) \\ \text{statistical error : } \mathcal{E}_{\text{sta}}}} |\mathcal{L}(f, G) - \widehat{\mathcal{L}}(f, G)| + \inf_{\substack{f \in \mathcal{F} \\ G \in \mathcal{G}(f) \\ \text{approximation error of } \mathcal{F} : \mathcal{E}_{\mathcal{F}}}} \sup_{\substack{G \in \mathcal{G}(f^*) \\ \text{statistical error : } \mathcal{E}_{\text{statistical error :$ 1568 1569 1570 1571 $+ \underbrace{\sup_{f \in \mathcal{F}} \{ \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f,G) - \mathcal{L}(f,\widehat{G}(f)) \} + 2(B_2^2 + \sqrt{d^*}) \sup_{f \in \mathcal{F}} \{ \mathbb{E}_{D_s} \left[\|\widehat{G}(f)\|_F \right] - \|G^*(f)\|_F \},}_{f \in \mathcal{F}}$ 1572 1573 approximation error of $\widehat{\mathcal{G}}$: $\mathcal{E}_{\widehat{\mathcal{C}}}$ 1574 1575 That is, 1576 $\mathcal{E}(\hat{f}_{n_{\tau}}) \leq 2\mathcal{E}_{\text{sta}} + \mathcal{E}_{\mathcal{F}} + \mathcal{E}_{\widehat{\mathcal{C}}}.$ 1577 1578 1579 *Proof.* Recall $\mathcal{F} = \mathcal{NN}_{d,d^*}(W, L, \mathcal{K}, B_1, B_2)$, for any $f \in \mathcal{F}$, 1580 $\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G) - \sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G)$ 1581 1582 $= \Big[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G) - \sup_{G \in \widehat{\mathcal{G}}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)\Big] + \Big[\sup_{G \in \widehat{\mathcal{G}}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G) - \sup_{G \in \widehat{\mathcal{G}}(\hat{f}_{n_s})} \widehat{\mathcal{L}}(\hat{f}_{n_s}, G)\Big]$ 1585 $+ \Big[\sup_{G \in \widehat{\mathcal{G}}(\widehat{f}_{n_s})} \widehat{\mathcal{L}}(\widehat{f}_{n_s}, G) - \sup_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{L}}(f, G) \Big] + \Big[\sup_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{L}}(f, G) - \sup_{G \in \widehat{\mathcal{G}}(f)} \mathcal{L}(f, G) \Big]$ 1586 1587 $+ \Big[\sup_{G \in \widehat{\mathcal{G}}(f)} \mathcal{L}(f,G) - \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f,G) \Big] + \Big[\sup_{G \in \mathcal{G}(f)} \mathcal{L}(f,G) - \sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*,G) \Big],$ 1589 1590 where the second and fourth terms can be bounded by \mathcal{E}_{sta} . In fact, regarding to the fourth term, we 1591 have 1592 $\sup_{G\in\widehat{\mathcal{G}}(f)}\widehat{\mathcal{L}}(f,G) - \sup_{G\in\widehat{\mathcal{G}}(f)}\mathcal{L}(f,G) \le \sup_{G\in\widehat{\mathcal{G}}(f)}\{\widehat{\mathcal{L}}(f,G) - \mathcal{L}(f,G)\}$ 1593 1594 $G \in \widehat{\mathcal{G}}(f)$ $\leq \sup_{G \in \widehat{\mathcal{G}}(f)} |\widehat{\mathcal{L}}(f,G) - \mathcal{L}(f,G)|$ 1596 $\leq \sup_{f \in \mathcal{F}, G \in \widehat{\mathcal{G}}(f)} |\widehat{\mathcal{L}}(f, G) - \mathcal{L}(f, G)|,$ 1597 1598 and the same conclusion holds for the second term. The addition of first term and fifth term can be bounded by $\mathcal{E}_{\widehat{G}}$. Actually, for the first term 1602 $\sup_{G\in\mathcal{G}(\widehat{f}_{n_s})}\mathcal{L}(\widehat{f}_{n_s},G) - \sup_{G\in\widehat{\mathcal{G}}(\widehat{f}_{n_s})}\mathcal{L}(\widehat{f}_{n_s},G) \leq \sup_{f\in\mathcal{F}} \{\sup_{G\in\mathcal{G}(f)}\mathcal{L}(f,G) - \sup_{G\in\widehat{\mathcal{G}}(f)}\mathcal{L}(f,G)\}$ 1603 1604 $\leq \sup\{ \sup \mathcal{L}(f,G) - \mathcal{L}(f,\widehat{G}(f))\},\$ $f \in \mathcal{F} \ G \in \mathcal{G}(f)$ 1607 (As $\widehat{G}(f) \in \widehat{\mathcal{G}}(f)$) 1608 and for the fifth term, we have 1609 1610 $\sup_{G \in \widehat{\mathcal{G}}(f)} \mathcal{L}(f,G) - \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f,G)$ 1611 $G \in \widehat{\mathcal{G}}(f)$ 1612 $= \sup_{G \in \widehat{\mathcal{G}}(f)} \mathbb{E}_{D_s} \left[\langle \widehat{G}(f), G \rangle_F \right] - \sup_{G \in \mathcal{G}(f)} \langle G^*(f), G \rangle_F \qquad (\langle G^*(f), G \rangle_F = \mathbb{E}_{D_s} \left[\langle \widehat{G}(f), G \rangle_F \right])$ 1613 $G \in \widehat{\mathcal{G}}(f)$ 1614 $\leq \mathbb{E}_{D_s} \big[\sup_{G \in \widehat{\mathcal{G}}(f)} \langle \widehat{G}(f), G \rangle_F \big] - \sup_{G \in \mathcal{G}(f)} \langle G^*(f), G \rangle_F$ 1615 1616 $= \mathbb{E}_{D_s} \left[\| \widehat{G}(f) \|_F^2 \right] - \| G^*(f) \|_F^2$ 1617 1618

1619
$$\leq 2(B_2^2 + \sqrt{d^*}) \left(\mathbb{E}_{D_s} \left[\| \widehat{G}(f) \|_F \right] - \| G^*(f) \|_F \right) \\ (\text{Both } \| \widehat{G}(f) \|_F \leq B_2^2 + \sqrt{d^*} \text{ and } \| G^*(f) \|_F \leq B_2^2 + \sqrt{d^*} \text{ hold})$$

which yields what we desired.

1624 Apart from that, the third term $\sup_{G \in \widehat{\mathcal{G}}(\widehat{f}_{n_s})} \widehat{\mathcal{L}}(\widehat{f}_{n_s}, G) - \sup_{G \in \widehat{\mathcal{G}}(f)} \widehat{\mathcal{L}}(f, G) \leq 0$ because of the 1625 definition of \widehat{f}_{n_s} . Taking infimum over all $f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K}, B_1, B_2)$ yields

$$\mathcal{E}(\hat{f}_{n_s}) \leq 2\mathcal{E}_{\mathrm{sta}} + \mathcal{E}_{\mathcal{F}} + \mathcal{E}_{\widehat{\mathcal{G}}},$$

which completes the proof.

1630 1631 Β.2.8 BOUND *E*_{sta}

1627

1634 1635 1636

1632 Lemma B.10. Regarding to \mathcal{E}_{sta} , we have

$$\mathbb{E}_{D_s}[\mathcal{E}_{ ext{sta}}] \lesssim rac{\mathcal{K}\sqrt{L}}{\sqrt{n_s}}$$

¹⁶³⁷ *Proof.* We are going to be introducing the relevant notations at first.

For any $f : \mathbb{R}^d \to \mathbb{R}^{d^*}$, let $\tilde{f} : \mathbb{R}^{2d} \to \mathbb{R}^{2d^*}$ such that $\tilde{f}(\tilde{x}) = (f(x_1), f(x_2))$, where $\tilde{x} = (x_1, x_2) \in \mathbb{R}^{2d}$. Furthermore, define $\tilde{\mathcal{F}} := \{\tilde{f} : f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})\}$ and denote $D'_s = \{\tilde{x}'^{(i)}\}_{i=1}^{n_s}$ as an independent identically distributed samples to D_s , which is called as ghost samples of D_s .

¹⁶⁴³ Next, we are attempt to establish the relationship between $\mathbb{E}_{D_s}[\mathcal{E}_{sta}]$ and the Rademacher complexity of $\mathcal{NN}_{d,d^*}(W,L,\mathcal{K})$. By the definition of \mathcal{E}_{sta} , we have

$$\begin{split} & \mathbb{E}_{D_{s}}[\mathcal{E}_{\text{sta}}] = \mathbb{E}_{D_{s}}\Big[\sup_{f \in \mathcal{NN}_{d,d^{*}}(W,L,K,B_{1},B_{2}), G \in \widehat{\mathcal{G}}(f)} |\mathcal{L}(f,G) - \widehat{\mathcal{L}}(f,G)|\Big] \\ & \leq \mathbb{E}_{D_{s}}\Big[\sup_{(f,G) \in \mathcal{NN}_{d,d^{*}}(W,L,K,B_{1},B_{2}) \times G_{1}} |\mathcal{L}(f,G) - \widehat{\mathcal{L}}(f,G)|\Big] \\ & \leq \mathbb{E}_{D_{s}}\Big[\sup_{(f,G) \in \mathcal{NN}_{d,d^{*}}(W,L,K,B_{1},B_{2}) \times G_{1}} |\mathcal{L}(f,G) - \widehat{\mathcal{L}}(f,G)|\Big] \\ & \leq \mathbb{E}_{D_{s}}\Big[\sup_{(f,G) \in \mathcal{NN}_{d,d^{*}}(W,L,K,N) \times G_{1}} |\mathcal{L}(f,G) - \widehat{\mathcal{L}}(f,G)|\Big] \\ & \leq \mathbb{E}_{D_{s}}\Big[\sup_{(f,G) \in \widehat{\mathcal{F}} \times G_{1}} \left|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \mathbb{E}_{D_{s}'}[\ell(\widetilde{f}(\widetilde{x}'^{(i)}),G)] - \frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \ell(\widetilde{f}(\widetilde{x}^{(i)}),G)\Big|\Big] \\ & \leq \mathbb{E}_{D_{s},D_{s}'}\Big[\sup_{(\widetilde{f},G) \in \widehat{\mathcal{F}} \times G_{1}} \left|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \ell(\widetilde{f}(\widetilde{x}'^{(i)}),G) - \frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \ell(\widetilde{f}(\widetilde{x}^{(i)}),G)\Big|\Big] \\ & \leq \mathbb{E}_{D_{s},D_{s}',\xi}\Big[\sup_{(\widetilde{f},G) \in \widehat{\mathcal{F}} \times G_{1}} \left|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \ell(\widetilde{f}(\widetilde{x}'^{(i)}),G) - \frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \ell(\widetilde{f}(\widetilde{x}^{(i)}),G)\Big|\Big] \\ & \leq \mathbb{E}_{D_{s},D_{s}',\xi}\Big[\sup_{(\widetilde{f},G) \in \widehat{\mathcal{F}} \times G_{1}} \left|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \ell(\widetilde{f}(\widetilde{x}'^{(i)}),G) - \ell(\widetilde{f}(\widetilde{x}^{(i)}),G)\Big)\Big|\Big] \\ & \leq 2\mathbb{E}_{D_{s},\mathcal{K}_{s}}\Big[\sup_{(\widetilde{f},G) \in \widehat{\mathcal{F}} \times G_{1}} \left|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \xi_{i}\ell(\widetilde{f}(\widetilde{x}^{(i)}),G)\Big|\Big] \\ & \leq 4\sqrt{2}\|\ell\|_{\mathrm{Lip}}\Big(\mathbb{E}_{D_{s},\xi}\Big[\sup_{f \in \mathcal{NN}_{d,d^{*}}(W,L,K)\Big| \frac{1}{n_{s}}\sum_{i=1}^{n_{s}}\sum_{j=1}^{d^{*}} \xi_{i,j,1}f_{j}(x_{1}^{(i)}) + \xi_{i,j,2}f_{j}(x_{2}^{(i)})\Big|\Big] \\ & + \mathbb{E}_{\xi}\Big[\sup_{G \in \mathcal{G}_{1}} \left|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}}\sum_{j=1}^{d^{*}} \xi_{i,j,k}G_{jk}\Big|\Big]\Big) \\ & \leq 8\sqrt{2}\|\ell\|_{\mathrm{Lip}}\mathbb{E}_{D_{s},\xi}\Big[\sup_{f \in \mathcal{NN}_{d,d^{*}}(W,L,K)\Big|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}}\sum_{j=1}^{d^{*}} \xi_{i,j,1}f_{j}(x_{1}^{(i)})\Big|\Big] + 4\sqrt{2}d^{*}\|\ell\|_{\mathrm{Lip}}\mathcal{P} \\ & \leq 8\sqrt{2}\|\ell\|_{\mathrm{Lip}}\mathbb{E}_{D_{s},\xi}\Big[\sup_{f \in \mathcal{NN}_{d,d^{*}}(W,L,K)\Big|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}}\sum_{j=1}^{d^{*}} \xi_{i,j,1}f_{j}(x_{1}^{(i)})\Big|\Big] + 4\sqrt{2}d^{*}\|\ell\|_{\mathrm{Lip}}\mathcal{P} \\ & \leq 8\sqrt{2}\|\ell\|_{\mathrm{Lip}}\mathbb{E}_{D_{s},\xi}\Big[\sup_{f \in \mathcal{NN}_{d,d^{*}}(W,L,K)\Big|\frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \sum_{j=1}^{d^{*}} \xi_{i,j,1}f_{j}(x_{1}^{(i)})\Big|\Big] + 4\sqrt{2}d^{*}\|\ell\|_{\mathrm{Lip}}\mathcal{P} \\ & \leq 8\sqrt{2}\|\ell\|_{\mathrm{Lip}}\mathbb{E}_{D_{s},\xi}\Big[\sup_{f$$

$$+ 4\sqrt{2} \|\ell\|_{\mathrm{Lip}} \mathbb{E}_{\xi} \Big[\max_{G \in \mathcal{N}_{\mathcal{G}_{1}}(\varrho)} \Big| \frac{1}{n_{s}} \sum_{i=1}^{n_{s}} \sum_{j=1}^{d^{*}} \sum_{k=1}^{d^{*}} \xi_{i,j,k} G_{jk} \Big| \Big]$$
(35)

 $\leq 8\sqrt{2} \|\ell\|_{\mathrm{Lip}} \mathbb{E}_{D_s,\xi} \Big[\sup_{f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K})} \Big| \frac{1}{n_s} \sum_{i=1}^{n_s} \sum_{j=1}^{d^*} \xi_{i,j} f_j(\boldsymbol{x}_1^{(i)}) \Big| \Big] + 4\sqrt{2} d^* \|\ell\|_{\mathrm{Lip}} \varrho$

$$+4\sqrt{2}(B_2^2+\sqrt{d^*})\|\ell\|_{\operatorname{Lip}}\sqrt{\frac{2\log\left(2|\mathcal{N}_{\mathcal{G}_1}(\varrho)|\right)}{n_s}}$$
(36)

$$\leq 8\sqrt{2}d^* \|\ell\|_{\mathrm{Lip}} \mathbb{E}_{D_s,\xi} \Big[\sup_{f \in \mathcal{NN}_{d,1}(W,L,\mathcal{K})} \Big| \frac{1}{n_s} \sum_{i=1}^{n_s} \xi_i f(\boldsymbol{x}_1^{(i)}) \Big| \Big] + 4\sqrt{2}d^* \|\ell\|_{\mathrm{Lip}} \varrho$$
$$+ 4\sqrt{2}(B_2^2 + \sqrt{d^*}) \|\ell\|_{\mathrm{Lip}} \sqrt{\frac{2\log\left(2(\frac{3}{(B_2^2 + \sqrt{d^*})\varrho})^{(d^*)^2}\right)}{n_0}}$$

$$(|\mathcal{N}_{\mathcal{G}_1}(\varrho)| \le \left(\frac{3}{(B_2^2 + \sqrt{d^*})\varrho}\right)^{(d^*)^2})$$

$$\lesssim \frac{\mathcal{K}\sqrt{L}}{\sqrt{n_s}} + \sqrt{\frac{\log n_s}{n_s}}$$
 (Lemma B.6 and set $\varrho = \mathcal{O}(1/\sqrt{n_s})$)
$$\lesssim \frac{\mathcal{K}\sqrt{L}}{\sqrt{n_s}}$$
 (If $\mathcal{K} \gtrsim \sqrt{\log n_s}$)

Where (33) stems from the fact that $\xi_i(\ell(\tilde{f}(\tilde{x}'^{(i)}), G) - \ell(\tilde{f}(\tilde{x}^{(i)}), G))$ has identical distribution with $\ell(\tilde{f}(\tilde{x}'^{(i)}), G) - \ell(\tilde{f}(\tilde{x}^{(i)}), G)$. As we have shown that $\|\ell\|_{\text{Lip}} < \infty$, just apply Lemma B.5 to obtain (34). Regarding (35), as $\mathcal{N}_{\mathcal{G}_1}(\rho)$ is a ρ -covering, for any fixed $G \in \mathcal{G}_1$, we can find a $H_G \in \mathcal{N}_{\mathcal{G}_1}(\rho)$ satisfying $||G - H_G||_F \leq \rho$, therefore we have

1700
1701
1702
1703

$$\mathbb{E}_{\xi} \Big[\max_{G \in \mathcal{G}_{1}} \Big| \frac{1}{n_{s}} \sum_{i=1}^{n_{s}} \sum_{j=1}^{d^{*}} \sum_{k=1}^{d^{*}} \xi_{i,j,k} \Big((H_{G})_{jk} + G_{jk} - (H_{G})_{jk} \Big) \Big| \Big]$$
1703
1703
1704
1705
1705
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
1707
17

$$\sum_{\substack{i=1\\G\in\mathcal{G}_1}} \sum_{i=1}^{n_s} \sum_{i=1}^{a} \sum_{j=1}^{a} \sum_{k=1}^{a} \xi_{i,j,k}(H_G)_{jk} \Big| \Big| + \mathbb{E}_{\xi} \Big[\max_{\substack{G\in\mathcal{G}_1\\G\in\mathcal{G}_1}} \Big| \frac{1}{n_s} \sum_{i=1}^{n_s} \sum_{j=1}^{a} \sum_{k=1}^{a} \xi_{i,j,k} \Big(G_{jk} - (H_G)_{jk} \Big) \Big| \Big]$$

$$\leq \mathbb{E}_{\xi} \Big[\max_{G \in \mathcal{N}_{\mathcal{G}_{1}}(\rho)} \Big| \frac{1}{n_{s}} \sum_{i=1}^{n_{s}} \sum_{j=1}^{d^{*}} \sum_{k=1}^{d^{*}} \xi_{i,j,k} G_{jk} \Big| \Big] + \frac{1}{n_{s}} \sqrt{(d^{*})^{2} n_{s}} \sqrt{n_{s}} \sum_{j=1}^{d^{*}} \sum_{k=1}^{d^{*}} \left(G_{jk} - (H_{G})_{jk} \right)^{2}$$
(Cauchy-Schwarz inequality)

$$\leq \mathbb{E}_{\xi} \Big[\max_{G \in \mathcal{N}_{\mathcal{G}_{1}}(\rho)} \Big| \frac{1}{n_{s}} \sum_{i=1}^{n_{s}} \sum_{j=1}^{d^{*}} \sum_{k=1}^{d^{*}} \xi_{i,j,k} G_{jk} \Big| \Big] + d^{*} \rho$$

To turn out the last term of (36), notice that $||G||_F \leq B_2^2 + \sqrt{d^*}$ implies that $\sum_{i=1}^{d^*} \sum_{k=1}^{d^*} \xi_{i,j,k} G_{jk} \sim$ $\operatorname{subG}(B_2^2 + \sqrt{d^*})$, therefore $\frac{1}{n_s} \sum_{i=1}^{n_s} \sum_{j=1}^{d^*} \sum_{k=1}^{d^*} \xi_{i,j,k} G_{jk} \sim \operatorname{subG}(B_2^2 + \sqrt{d^*})$, just apply Lemma B.7 to finish the proof.

B.2.9 BOUND $\mathcal{E}_{\mathcal{F}}$

If we denote

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

where $C([0,1]^d)$ is the space of continuous functions on $[0,1]^d$ equipped with the sup-norm. Theo-rem 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

Lemma B.11 (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 $\mathcal{K} \ge 1$, any $W \ge c\mathcal{K}^{(2d+\alpha)/(2d+2)}$ and $L \ge 2\lceil \log_2(d+r) \rceil + 2,$

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

For utilizing this conclusion, first notice that

$$\inf_{f\in\mathcal{NN}_{d,d^*}(W,L,\mathcal{K})}\|f(oldsymbol{u})-f^*(oldsymbol{u})\|_2$$

1737
1738
1739
1740
$$= \inf_{f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K})} \sqrt{\sum_{i=1}^{d^*} (f_i(u) - f_i^*(u))^2}$$

$$\leq \inf_{f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K})} \sqrt{\sum_{i=1}^{d^*} \|f_i - f_i^*\|_{C([0,1]^d)}^2}$$

1744
1745
$$\leq \sup_{g \in \mathcal{H}^{\alpha}} \inf_{f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K})} \sqrt{\sum_{i=1}^{d^*} \|f_i - g\|_{C([0,1]^d)}^2}$$

1746

$$\begin{aligned} & 1747 \\ & 1748 \\ & 1749 \\ & 1750 \\ & 1751 \end{aligned} \leq \sup_{g \in \mathcal{H}^{\alpha}} \sqrt{\sum_{i=1}^{d^*} \inf_{f \in \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K})} \|f - g\|_{C([0,1]^d)}^2} \\ & \leq \sqrt{d^*} \mathcal{E}(\mathcal{H}^{\alpha}, \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K})) \end{aligned}$$

1752
$$\leq \mathcal{K}^{-\alpha/(d+1)},$$

where the third to last line inequality is from following reason: if $f_i \in \mathcal{NN}_{d,1}(|W/d^*|, L, \mathcal{K})$, 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 B.2, we have $f \in \mathcal{NN}_{d,d^*}(W, L, \mathcal{K})$.

Proposition B.2 ((iii) of Proposition 2.5 in Jiao et al. (2023)). Let $\phi_1 \in \mathcal{NN}_{d,d_1^*}(w_1, L_1, \mathcal{K}_1)$ and $\phi_2 \in \mathcal{NN}_{d,d_2^*}(W_2, L_2, \mathcal{K}_2)$, define $\phi(x) := (\phi_1(x), \phi_2(x))$, then $\phi \in \mathcal{NN}_{d,d_1^*+d_2^*}(W_1 + d_2^*)$ $W_2, \max\{L_1, L_2\}, \max\{\mathcal{K}_1, \mathcal{K}_2\}).$

Above conclusion implies optimal approximation element of f^* in $\mathcal{NN}_{d,d^*}(W,L,\mathcal{K})$ can be ar-bitrarily 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 $B_1 \le \|f^*\|_2 \le B_2.$

Therefore, if we denote

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

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

· . . .

$$\begin{aligned} \mathcal{E}_{\mathcal{F}} &= \inf_{f \in \mathcal{F}} \{ \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f, G) - \sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G) \} \\ &= \inf_{f \in \mathcal{F}} \{ \mathcal{T}(f) - \mathcal{T}(f^*) \} \\ &= \inf_{f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K})} \{ \mathcal{T}(f) - \mathcal{T}(f^*) \} \\ &\leq \| \ell \|_{\mathrm{Lip}} \inf_{f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K})} \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\tilde{\boldsymbol{x}}} \| \widetilde{f}(\tilde{\boldsymbol{x}}) - \widetilde{f}^*(\tilde{\boldsymbol{x}}) \|_2 \qquad (\text{Proposition B.1}) \\ &\leq \| \ell \|_{\mathrm{Lip}} \inf_{f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K})} \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}' \in \mathcal{A}(\boldsymbol{x})} \sqrt{2 \sum_{i=1}^{d^*} (f_i(\boldsymbol{x}') - f_i^*(\boldsymbol{x}'))^2} \end{aligned}$$

1779
1779
$$\leq \|\ell\|_{\operatorname{Lip}} \inf_{f \in \mathcal{NN}_{d,d^*}(W,L,\mathcal{K})} \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}' \in \mathcal{A}(\boldsymbol{x})} \bigvee 2 \sum_{i=1}^{\infty} (f_i(\boldsymbol{x}') - f_i^*(\boldsymbol{x}'))$$
1780

1781
$$\leq \sqrt{2d^*} \|\ell\|_{\operatorname{Lip}} \sup_{g \in \mathcal{H}^{\alpha}} \inf_{f \in \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K}/\sqrt{d^*})} \|f - g\|_{C([0,1]^d)}$$

1782
1783
1784
$$\leq \sqrt{2d^*} \|\ell\|_{\operatorname{Lip}} \mathcal{E}(\mathcal{H}^{\alpha}, \mathcal{NN}_{d,1}(\lfloor W/d^* \rfloor, L, \mathcal{K}/\sqrt{d^*}))$$

$$\lesssim \mathcal{K}^{-\alpha/(d+1)}.$$

B.2.10 BOUND $\mathcal{E}_{\widehat{G}}$

 $\operatorname{Recall} \mathcal{E}_{\widehat{\mathcal{G}}} = \sup_{f \in \mathcal{F}} \{ \sup_{G \in \mathcal{G}(f)} \mathcal{L}(f,G) - \mathcal{L}(f,\widehat{G}(f)) \} + 2(B_2^2 + \sqrt{d^*}) \sup_{f \in \mathcal{F}} \{ \mathbb{E}_{D_s} \left[\|\widehat{G}(f)\|_F \right] - C(f,\widehat{G}(f)) \} + 2(B_2^2 + \sqrt{d^*}) \sum_{f \in \mathcal{F}} \{ \mathbb{E}_{D_s} \left[\|\widehat{G}(f)\|_F \right] + C(f,\widehat{G}(f)) \} + 2(B_2^2 + \sqrt{d^*}) \sum_{f \in \mathcal{F}} \{ \mathbb{E}_{D_s} \left[\|\widehat{G}(f)\|_F \right] + C(f,\widehat{G}(f)) \} + C(f,\widehat{G}(f)) \}$ $||G^*(f)||_F$, then for the first item of $\mathcal{E}_{\widehat{G}}$, we have $\sup_{f\in\mathcal{F}}\{\sup_{G\in\mathcal{G}(f)}\mathcal{L}(f,G)-\mathcal{L}(f,\widehat{G}(f))\}$ $= \sup_{f \in \mathcal{F}} \{ \mathcal{L}(f, G^*(f)) - \mathcal{L}(f, \widehat{G}(f)) \}$

 $\leq \sqrt{2}B_2 \sup_{f \in \mathcal{T}} \|G^*(f) - \widehat{G}(f)\|_F \qquad (\text{Look up Table 2 to yield } \ell(\boldsymbol{u}, \cdot) \in \operatorname{Lip}(\sqrt{2}B_2))$

 $\leq \sqrt{2}B_2 \sup_{f \in \mathcal{F}} \left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - \frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top \right\|_F.$

And regrading to the second term, we can yield

 $\sup_{f \in \mathcal{F}} \{ \mathbb{E}_{D_s} \big[\| \widehat{G}(f) \|_F \big] - \| G^*(f) \|_F \}$ $= \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{D_s} \left[\left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - I_{d^*} \right\|_F - \left\| \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] - I_{d^*} \right\|_F \right] \right\}$ $\leq \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{D_s} \left[\left\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] \right\|_F \right] \right\}$ $\leq \mathbb{E}_{D_s} \bigg[\sup_{f \in \mathcal{F}} \Big\{ \Big\| \frac{1}{n_s} \sum_{i=1}^{n_s} f(\boldsymbol{x}_1^{(i)}) f(\boldsymbol{x}_2^{(i)})^\top - \mathbb{E}_{\boldsymbol{x}} \mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})} [f(\boldsymbol{x}_1) f(\boldsymbol{x}_2)^\top] \Big\|_F \Big\} \bigg]$

Combine above two inequalities to turn out

$$\begin{split} \mathbb{E}_{D_s}[\mathcal{E}_{\widehat{\mathcal{G}}}] \lesssim \mathbb{E}_{D_s}\Big[\sup_{f \in \mathcal{F}} \left\|\mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})}\Big[\frac{1}{n_s}\sum_{i=1}^{n_s}\left[\mathcal{M}(\widetilde{f}(\widetilde{\boldsymbol{x}})) - \mathcal{M}(\widetilde{f}(\widetilde{\boldsymbol{x}}^{(i)}))\right]\Big]\right\|_F \\ \leq \|\mathcal{M}\|_{\mathrm{Lip}}\mathbb{E}_{D_s}\Big[\left\|\mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})}[\widetilde{f}(\widetilde{\boldsymbol{x}})] - \frac{1}{n_s}\sum_{i=1}^{n_s}\widetilde{f}(\widetilde{\boldsymbol{x}}^{(i)})\right\|_2\Big] \end{split}$$

where $\mathcal{M}(\boldsymbol{u}) = \boldsymbol{u}_1 \boldsymbol{u}_2^{ op}$, where $\boldsymbol{u}_1, \boldsymbol{u}_2 \in \mathbb{R}^{d^*}$, we have shown it is a Lipchitz map on $\{\boldsymbol{u} \in$ \mathbb{R}^{2d^*} : $u \leq \sqrt{2}B_2$ } in Proposition B.1. By Multidimensional Chebyshev's inequality, we know that $P_s\left(\left\|\frac{1}{n_s}\sum_{i=1}^{n_s}\widetilde{f}(\widetilde{\boldsymbol{x}}^{(i)}) - \mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1,\boldsymbol{x}_2\in\mathcal{A}(\boldsymbol{x})}[\widetilde{f}(\widetilde{\boldsymbol{x}})]\right\|_2 \geq \frac{1}{n_s^{1/4}}\right) \leq \frac{\mathbb{E}\|\widetilde{f}(\widetilde{\boldsymbol{x}}) - \mathbb{E}[\widetilde{f}(\widetilde{\boldsymbol{x}})]\|_2^2}{\sqrt{n_s}} \leq \frac{8B_2^2}{\sqrt{n_s}}$ as $\|\tilde{f}(\tilde{\boldsymbol{x}})\|_2 \leq \sqrt{2}B_2$. Thus we have

$$\mathbb{E}_{D_s}[\mathcal{E}_{\widehat{\mathcal{G}}}] \lesssim \frac{1}{n_s^{1/4}} \cdot P_s\left(\left\|\frac{1}{n_s}\sum_{i=1}^{n_s} \widetilde{f}(\widetilde{\boldsymbol{x}}^{(i)}) - \mathbb{E}_{\boldsymbol{x}}\mathbb{E}_{\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{A}(\boldsymbol{x})}[\widetilde{f}(\widetilde{\boldsymbol{x}})]\right\|_2 \ge \frac{1}{n_s^{1/4}}\right) + 2\sqrt{2}B_2 \cdot \frac{8B_2^2}{\sqrt{n_s}}$$

$$(As \|\widetilde{f}(\widetilde{\boldsymbol{x}})\|_2 \le \sqrt{2}B_2)$$

$$\leq rac{1}{n_s^{1/4}} + 16\sqrt{2}B_2^3rac{1}{\sqrt{n_s}} \ \lesssim rac{1}{n_s^{1/4}}.$$

B.2.11 TRADE OFF BETWEEN STATISTICAL ERROR AND APPROXIMATION ERROR

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

$$\mathbb{E}_{D_s}[\mathcal{E}(\hat{f}_{n_s})] \le 2\mathbb{E}_{D_s}[\mathcal{E}_{\text{sta}}] + \mathcal{E}_{\mathcal{F}} + 2\mathbb{E}_{D_s}[\mathcal{E}_{\widehat{\mathcal{G}}}] \lesssim \frac{\mathcal{K}}{\sqrt{n_s}} + \mathcal{K}^{-\alpha/(d+1)}.$$

Taking $\mathcal{K} = n_s^{\frac{d+1}{2(\alpha+d+1)}}$ to yield $\mathbb{E}_{D_s}[\mathcal{E}(\hat{f}_{n_s})] \lesssim n_s^{-\frac{\alpha}{2(\alpha+d+1)}}$ As we have shown that $\sup_{G \in \mathcal{G}(f^*)} \mathcal{L}(f^*, G) = 0$, above inequality implies $\mathbb{E}_{D_s}[\sup_{G\in\mathcal{G}(\hat{f}_{n_s})}\mathcal{L}(\hat{f}_{n_s},G)] \lesssim n_s^{-\frac{\alpha}{2(\alpha+d+1)}}.$ To ensure above deduction holds, We need to set $W \ge cn_s^{\frac{2d+\alpha}{4(\alpha+d+1)}}$ and $L \ge 2\lceil \log_2(d+r) \rceil + 2$. B.2.12 THE PROOF OF MAIN THEOREM Next, we are going to prove our main theorem 4.2. We will state its formal version at first and then conclude Theorem 4.2 as a corollary. To notation conciseness, let $p = \frac{\sqrt{\frac{2}{\underset{i\neq j}{\min p_s(i)p_s(j)}} \left(\frac{C}{\lambda} n_s^{-\frac{\alpha}{2(\alpha+d+1)}} + \psi(n_s)\right)} + 2\sqrt{d^*} B_2 M n_s^{-\frac{\nu}{2(\alpha+d+1)}}}{B_2^2 \Theta(\sigma_s^{(n_s)}, \delta_s^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s})}, \text{ where } C \text{ is a constant, } 0 \le \psi(n_s) \lesssim (1 - \sigma_s^{(n_s)} + n_s^{-\frac{\min\{\alpha,\nu,\varsigma,\tau\}}{4(\alpha+d+1)}})^2 + (1 - \sigma_s^{(n_s)}) + n_s^{-\frac{\min\{\alpha,\nu,\varsigma,\tau\}}{8(\alpha+d+1)}}, \text{ then } 0 \le \psi(n_s) \lesssim (1 - \sigma_s^{(n_s)}) + n_s^{-\frac{\min\{\alpha,\nu,\varsigma,\tau\}}{4(\alpha+d+1)}} = 0$ the formal version of our main theoretical result can be stated as following the formal version of α and α and β are the formal version of α are the formal version of α and β are the formal version of α are the formal version of α and β are the formal version of α are the format version of α **Lemma B.12.** When Assumption 4.1-4.5 all hold, set $\varepsilon_{n_s} = m^2 n_s^{-\frac{\min\{\alpha,\nu,\varsigma,\tau\}}{8(\alpha+d+1)}}, W \ge c n_s^{\frac{2d+\alpha}{4(\alpha+d+1)}}.$ $L \ge 2\lceil \log_2(d+r) \rceil + 2, \mathcal{K} = n_s^{\frac{d+1}{2(\alpha+d+1)}}$ and $\mathcal{A} = \mathcal{A}_{n_s}$ in Assumption 4.3, then we have $\mathbb{E}_{D_s}[R_t^2(\varepsilon_{n_s}, \hat{f}_{n_s})] \lesssim n_s^{-\frac{\min\{\alpha, \nu, \varsigma\}}{4(\alpha+d+1)}}$ (37) and $\mathbb{E}_{D_s}[\max_{i \neq j} |\mu_t(i)^\top \mu_t(j)|] \lesssim 1 - \sigma_s^{(n_s)} + n_s^{-\frac{\min\{\alpha, 2\tau\}}{4(\alpha+d+1)}}.$ (38)Furthermore, If $\Theta(\sigma_s^{(n_s)}, \delta_s^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) > 0$, then with probability at least 1 - p, we have $\mathbb{E}_{D_s}[\operatorname{Err}(Q_{\hat{f}_{n_s}})] \le (1 - \sigma_t^{(n_s)}) + \mathcal{O}(n_s^{-\frac{\min\{\alpha,\nu,\varsigma\}}{8(\alpha+d+1)}}).$ *Proof.* First recall the conclusion we've got in Theorem B.1 $\mathbb{E}_{D_s}[R_t^2(\varepsilon, \hat{f}_{n_s})] \le \frac{m^4}{\varepsilon^2} \big(\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s}, G)}] + 8B_2 d^* M \mathcal{K}\rho + 4B_2^2 d^* K \eta \big),$ $\mathbb{E}_{D_s}[\max_{i\neq j} |\mu_t(i)^\top \mu_t(j)|] \le \sqrt{\frac{2}{\min_{i\neq j} p_s(i) p_s(j)} \left(\frac{1}{\lambda} \mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + \mathbb{E}_{D_s}[\psi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})]\right)}$ $+2\sqrt{d^*}B_2M\mathcal{K}\rho$ and with probability at least $1 - \frac{\sqrt{\frac{2}{\min_{i \neq j} p_s(i) p_s(j)} \left(\frac{1}{\lambda} \mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + \psi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s})\right)}{B_2^2 \Theta(\sigma_t, \delta_t, \varepsilon, \hat{f}_{n_s})} + 2\sqrt{d^*} B_2 M \mathcal{K} \rho$ we have

 $\mathbb{E}_{D_s}[\operatorname{Err}(Q_{\hat{f}_{n_s}})] \le (1 - \sigma_t) + \frac{m^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] + 8B_2 d^* M \mathcal{K} \rho + 4B_2^2 d^* K \eta},$

where

$$\begin{aligned} & \frac{1890}{1891} \qquad \psi(\sigma_s, \delta_s, \varepsilon, \hat{f}_{n_s}) = 4B_2^2 \Big[\Big(1 - \sigma_s + \frac{\kappa \delta_s + 2\varepsilon}{2B_2} \Big)^2 + (1 - \sigma_s) + \frac{Km^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s} [\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)]} \Big(3 - \frac{1}{2} \sum_{i=1}^{N} \frac{1}{\varepsilon} \Big(\frac{1}{\varepsilon} \Big)^2 \Big] \Big] \\ & \frac{1}{\varepsilon} \Big(\frac{1}{\varepsilon} \Big)^2 \Big[\left(\frac{1}{\varepsilon} - \frac{1}{\varepsilon} \sum_{i=1}^{N} \frac{1}{\varepsilon} \Big)^2 \Big] \Big] \Big(\frac{1}{\varepsilon} \Big)^2 \Big] \\ & \frac{1}{\varepsilon} \sum_{i=1}^{N} \frac{1}{\varepsilon} \sum_{i=1}^{$$

 $2\sigma_s + \frac{\kappa\delta_s + 2\varepsilon}{B_2} + \frac{m^4}{\varepsilon^2} \mathbb{E}_{D_s} [\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] \Big(\sum_{k=1}^K \frac{1}{p_s(k)}\Big) \Big] + B_2 \Big(\varepsilon^2 + \frac{4B_2^2m^2}{\varepsilon} \sqrt{\mathbb{E}_{D_s} [\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)]} \Big)^{\frac{1}{2}}.$

To obtain the conclusion shown in this theorem from above formulations, first notice $\rho = n_s^{-\frac{\nu+d+1}{2(\alpha+d+1)}}$ and $\eta = n_s^{-\frac{\varsigma}{2(\alpha+d+1)}}$ by comparing Assumption 4.4 and Assumption B.1, apart from that, we have shown $\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)] \lesssim n_s^{-\frac{\alpha}{2(\alpha+d+1)}}$ in B.2.11 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 $\mathcal{K}\rho/\varepsilon_{n_s}^2 \leq n_s^{-\frac{\tau}{2(\alpha+d+1)}}, \eta/\varepsilon_{n_s}^2 \leq n_s^{-\frac{\tau}{2(\alpha+d+1)}}, \mathcal{K}\delta_s^{(n_s)} \leq n_s^{-\frac{\tau}{2(\alpha+d+1)}}$ and $\mathbb{E}_{D_s}[\sup_{G \in \mathcal{G}(\hat{f}_{n_s})} \mathcal{L}(\hat{f}_{n_s}, G)]/\varepsilon_{n_s}^2 \leq n_s^{-\frac{\alpha}{4(\alpha+d+1)}}.$

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

Let us first state the formal version of Theorem 4.2 and then prove it.

Theorem B.3 (Formal version of Theorem 4.2). If Assumptions 4.1-4.5 all hold, set $W \ge cn_s^{\frac{2d+\alpha}{4(\alpha+d+1)}}$, $L \ge 2\lceil \log_2(d+r) \rceil + 2$, $\mathcal{K} = n_s^{\frac{d+1}{2(\alpha+d+1)}}$ and $\mathcal{A} = \mathcal{A}_{n_s}$ in Assumption 4.3, 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

$$\mathbb{E}_{D_s}[\operatorname{Err}(Q_{\hat{f}_{n_s}})] \le (1 - \sigma_t^{(n_s)}) + \mathcal{O}(n_s^{-\frac{\min\{\alpha,\nu,\varsigma\}}{8(\alpha+d+1)}})$$

1922 To show this, first recall

$$\begin{array}{l} \begin{array}{l} 1923\\ 1924\\ 1924\\ 1925\\ 1926\\ 1926\\ 1927 \end{array} & \Theta(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) = \Gamma_{\min}(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) - \sqrt{2 - 2\Gamma_{\min}(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s})} - \frac{\Delta_{\hat{\mu}_t}}{2} \\ - \frac{2\max_{k \in [K]} \|\hat{\mu}_t(k) - \mu_t(k)\|_2}{B_2}. \end{array}$$

1928 Note (31) and dominated convergence theorem imply $R_t(\varepsilon_{n_s}, \hat{f}_{n_s}) \to 0$ a.s., thus

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

1935 Combining with the fact that $\frac{\Delta_{\hat{\mu}_t}}{2} = \frac{1 - \min_{k \in [K]} \|\hat{\mu}_t(k)\|^2 / B_2^2}{2} < \frac{1}{2}$ can yield

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

1939 if we select proper B_1 and B_2 .

¹⁹⁴⁰ Besides that, by Multidimensional Chebyshev's inequality, we know that

$$P_t \left(\|\hat{\mu}_t(k) - \mu_t(k)\|_2 \ge \frac{B_2}{8} \right) \le \frac{64\sqrt{\mathbb{E}_{\boldsymbol{z} \in \tilde{C}_t(k)} \mathbb{E}_{\boldsymbol{z}' \in \mathcal{A}(\boldsymbol{z})} \|f(\boldsymbol{z}') - \mu_t(k)\|_2^2}}{B_2^2 \sqrt{2n_t(k)}} \le \frac{128}{B_2 \sqrt{n_t(k)}}$$

so that $\Theta(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) \geq \frac{1}{4}$ with probability at least $1 - \frac{128K}{B_2\sqrt{\min_k n_t(k)}}$ if n_s is large enough, of course the condition $\Theta(\sigma_t^{(n_s)}, \delta_t^{(n_s)}, \varepsilon_{n_s}, \hat{f}_{n_s}) > 0$ in Theorem B.12 can be satisfied.

1948 Therefore, with probability at least

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

¹⁹⁵⁴ we have the conclusions shown in Theorem 4.2, which completes the proof.

C EXPERIMENTAL DETAILS

1958 **Implementation details.** Except for tuning λ for different dataset, all other hyper parameters used 1959 in our experiments are align with Ermolov et al. (2021). To be specific, we train 1,000 epochs with 1960 learning rate 3×10^{-3} for CIFAR-10, CIFAR-100 and 2×10^{-3} for Tiny ImageNet. The learning rate 1961 warm-up is used for the first 500 iterations of the optimizer, in addition to a 0.2 learning rate drop 1962 50 and 25 epochs before the end. We adopt a mini-batch size of 256. Same as W-MSE 4 of Ermolov 1963 et al. (2021), we also set 4 as the number of positive samples per image. The dimension of the hidden layer of the projection head is set as 1024. The weight decay is 10^{-6} . We adopt an embedding size 1964 (d^*) of 64 for CIFAR10, CIFAR100 and 128 for Tiny ImageNet and employ the trick mentioned in 1965 Ermolov et al. (2021) during the pretraining process. The embedding size of BarlowTwins (Zbontar 1966 et al., 2021) is different from above as BarlowTwins need much larger representation size (1024) 1967 to guarantee its performance. As we see, the performance of our model can sufficiently outperform 1968 BarlowTwins, revealing the alignment term is pretty crucial for downstream performance practically. 1969 The backbone network used in our implementation is ResNet-18. 1970

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; 0.1) 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 pretraining, 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.



1984

1949

1955 1956

- 1985
- 1986
- 1987 1988
- 1989
- 1990
- 1991
- 1992
- 1993
- 1994
- 1995
- 1996
- 1997