
Understanding Transferable Representation Learning and Zero-shot Transfer in CLIP

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Multi-modal learning has become increasingly popular due to its ability to leverage
2 information from different data sources. Recently, CLIP has emerged as an effective
3 approach that employs vision-language contrastive pretraining to learn joint image
4 and text representations and exhibits remarkable performance in zero-shot learning
5 and text-guided natural image generation. Despite the huge practical success of
6 CLIP, its theoretical understanding remains elusive. In this paper, we formally
7 study transferrable representation learning underlying CLIP and demonstrate how
8 features from different modalities get aligned. We also analyze its zero-shot transfer
9 performance on the downstream tasks. Inspired by our analysis, we propose a
10 new CLIP-type approach, which achieves better performance than CLIP and other
11 state-of-the-art methods on benchmark datasets.

12 1 Introduction

13 Recently, CLIP (Radford et al., 2021) emerged as a milestone work that leverages vision-language con-
14 trastive pretraining to jointly learn image and text embeddings, using the vast amounts of image-text
15 data available on the web. This approach has achieved remarkable success in zero-shot transfer (Lei Ba
16 et al., 2015). Inspired by CLIP’s groundbreaking zero-shot capabilities, subsequent studies (Yao
17 et al., 2022; Li et al., 2022; Mu et al., 2022; Goel et al., 2022; Zhai et al., 2022; Alayrac et al., 2022)
18 emerged with the primary objective of further enhancing CLIP’s zero-shot performance. Despite
19 the empirical success of CLIP in zero-shot transfer, the theoretical understanding of how it works
20 remains elusive.

21 This paper delves into the mechanisms through which CLIP learns transferable representations and
22 demonstrates how such representations ensure successful zero-shot transfer for downstream tasks.
23 We present our theoretical result for transferable representation learning in CLIP and summarize our
24 contributions as follows.

- 25 • We theoretically examine transferable representation learning in CLIP. Our analysis shows that if a
26 near-optimal network is obtained on the training data, features from different modalities become
27 aligned, enabling zero-shot learning if appropriate prompts are issued.
- 28 • Building upon our general theoretical findings, we delve deeper into specific cases. We illustrate
29 how multi-modal learning aligns different features and reveal when the learned features obtained
30 by CLIP can outperform those obtained through naive square loss.
- 31 • We conduct experiments on real data to confirm our theoretical predictions. Furthermore, inspired
32 by our theoretical findings, we propose a new regularization technique for CLIP that effectively
33 leads to improved zero-shot performance.

34 2 Problem Setting and Preliminaries

35 2.1 Data Distribution

36 In our paper, we focus on the setting where the image x and the text y are conditionally independent
37 given the shared feature z .

38 **Assumption 2.1.** Let (\mathbf{x}, \mathbf{y}) be generated from the joint distribution $\mathcal{D}_{\mathbf{x} \times \mathbf{y}}$. We assume \mathbf{z} to be
 39 a shared feature of \mathbf{x}, \mathbf{y} satisfying $\mathbf{x} \perp \mathbf{y} | \mathbf{z}$, and further denote $(\mathbf{x}, \mathbf{y}, \mathbf{z})$ that follows the joint
 40 distribution $\mathcal{D}_{\mathbf{x} \times \mathbf{y} \times \mathbf{z}}$ with marginal distributions $\mathcal{D}_{\mathbf{x} \times \mathbf{z}}, \mathcal{D}_{\mathbf{y} \times \mathbf{z}}$. We further assume \mathbf{z} to be a discrete
 41 and sparse random variable $\mathbf{z} \in \mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_K\}$ with $p_k := \mathbb{P}(\mathbf{z} = \mathbf{v}_k)$.

42 2.2 Learning via Contrastive Loss

43 The CLIP architecture has three main components: (i) an image encoder network \mathbf{g} that can encode
 44 the image \mathbf{x} into the embedding $\mathbf{g}(\mathbf{x}) \in \mathbb{R}^d$; (ii) a text encoder network \mathbf{h} that can encode the
 45 text \mathbf{y} into an embedding vector $\mathbf{h}(\mathbf{y}) \in \mathbb{R}^d$; and (iii) a score function $f(\mathbf{x}, \mathbf{y}) = \mathbf{sim}(\mathbf{g}, \mathbf{h})$ that
 46 measures the similarity between the image \mathbf{x} and the text \mathbf{y} given their embeddings \mathbf{g}, \mathbf{h} (e.g.,
 47 $f(\mathbf{x}, \mathbf{y}) = \langle \mathbf{g}(\mathbf{x}), \mathbf{h}(\mathbf{y}) \rangle$). During the training, we will sample a batch of image-captions pairs
 48 $S' = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^B \subseteq S$. The contrastive loss in CLIP aims to align the image representation $\mathbf{g}(\mathbf{x})$ and
 49 text representations $\mathbf{h}(\mathbf{y})$ by minimizing the empirical version of the following population loss,

$$L_{\mathcal{D}^B}(f, \tau) = \mathbb{E} \left[\log \left(\sum_{t \in [B]} \exp([f(\mathbf{x}_t, \mathbf{y}_t) - f(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right] \\ + \mathbb{E} \left[\log \left(\sum_{t \in [B]} \exp([f(\mathbf{x}_t, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right], \quad (2.1)$$

50 where $\tau > 0$ is a temperature parameter and the expectation is taken with respect to all B random
 51 pairs $(\mathbf{x}_t, \mathbf{y}_t)$ i.i.d. sampled from $\mathcal{D}_{\mathbf{x} \times \mathbf{y}}$. Therefore, CLIP learns the score function f with the
 52 corresponding representations \mathbf{g} and \mathbf{h} by minimizing $L_{\mathcal{D}^B}(f, \tau)$. In fact, we can divide the training
 53 dataset S into n batches $\cup_{k \in [n]} \mathcal{S}_k$. Further discussion of problem setting is deferred to Appendix B.

54 3 Zero-shot Transfer

55 The key idea of CLIP is to pull the embeddings of positive image-text pairs together while pushing the
 56 embeddings of negative pairs apart. For the data pair $(\mathbf{x}, \mathbf{y}')$ generated with $\mathbf{x} \sim \mathcal{D}_{\mathbf{x} | \mathbf{z}}, \mathbf{y}' \sim \mathcal{D}_{\mathbf{x} | \mathbf{z}'}$,
 57 $(\mathbf{x}, \mathbf{y}')$ is a positive pair if $\mathbf{z} = \mathbf{z}'$ and a negative pair if $\mathbf{z} \neq \mathbf{z}'$.

58 **Assumption 3.1** ((α, β, γ) -Completeness). There exists a score function f^* bounded by 1 (i.e.,
 59 $|f^*| \leq 1$) with $f^* = \mathbf{sim}(\mathbf{g}^*, \mathbf{h}^*)$ satisfying the following properties,

- 60 • For any $\mathbf{z} \neq \mathbf{z}'$, let $\mathbf{x} \sim \mathcal{D}_{\mathbf{x} | \mathbf{z}}, \mathbf{y} \sim \mathcal{D}_{\mathbf{y} | \mathbf{z}}, \mathbf{x}' \sim \mathcal{D}_{\mathbf{x} | \mathbf{z}'}, \mathbf{y}' \sim \mathcal{D}_{\mathbf{y}' | \mathbf{z}'}$. With probability at least $1 - \alpha$,
 61 we have $f^*(\mathbf{x}', \mathbf{y}) \leq f^*(\mathbf{x}, \mathbf{y}) - \gamma$ and $f^*(\mathbf{x}, \mathbf{y}') \leq f^*(\mathbf{x}, \mathbf{y}) - \gamma$.
- 62 • Let $(\mathbf{x}, \mathbf{y}, \mathbf{z}) \sim \mathcal{D}_{\mathbf{x} \times \mathbf{y} \times \mathbf{z}}$, assume $\mathbb{E}_{(\mathbf{y}, \mathbf{z})} [\text{Var}_{\mathbf{x} | \mathbf{z}}(f^*(\mathbf{x}, \mathbf{y}))], \mathbb{E}_{(\mathbf{x}, \mathbf{z})} [\text{Var}_{\mathbf{y} | \mathbf{z}}(f^*(\mathbf{x}, \mathbf{y}))] \leq \beta$.

63 Further discussion on Assumption 3.1 can be found in Appendix G. In the zero-shot transfer task, we
 64 have K prompts $\{\mathbf{y}_k, k \in [K]\}$ where $\mathbf{y}_k \sim \mathcal{D}_{\mathbf{y} | \mathbf{v}_k}$. For a new image \mathbf{x} generated from $\mathcal{D}_{\mathbf{x}}$, we want
 65 to predict the label of the shared feature \mathbf{z} in \mathbf{x} . The following theorem provides the guarantee of
 66 zero-shot transfer learning for CLIP.

67 **Theorem 3.2** (Informal). Suppose Assumption 3.1 hold and we can find an ϵ approximate minimum
 68 $\hat{f} \in \mathcal{F}$ with respect to the temperature τ such that \hat{f} is bounded by M and

$$L_{\mathcal{D}^B}(\hat{f}, \tau) \leq L_{\mathcal{D}^B}(f^*, \tau) + \epsilon. \quad (3.1)$$

69 For the zero-shot downstream task, we calculate the similarity score $\hat{f}(\mathbf{x}, \mathbf{y}_k)$ for all $k \in [K]$ and pick
 70 the indices of the top- r scores within the set $\{\hat{f}(\mathbf{x}, \mathbf{y}_k)\}$ as the predictions of the image \mathbf{x} . The top- r
 71 error is bounded by $\epsilon' / \log(1 + r)$, where $\epsilon' = (C_B + 2) \cdot [\epsilon + C\tau^{-1}MB\alpha + C\tau^{-1}(\beta MB)^{1/3} +$
 72 $2B \exp(-\gamma/\tau)]$ and $C = \tilde{O}(1), C_B = \tilde{O}(\max_k p_k^{-1}/B)$.

73 A formal discussion is presented in Appendix G and H. Next, we will introduce a specific problem to
 74 illustrate how CLIP can learn transferable features with distinguishable margins, which is hard to
 75 achieve by simple square loss.

76 **Definition 3.3** (A Case Study). Let shared feature $\mathbf{z} \in \mathbb{R}^{K_1}$ be random variable uniformly drawn from
 77 the set $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_K\}$ where $\|\mathbf{v}_k\|_2 = 1, \max_{k \neq k'} \langle \mathbf{v}_k, \mathbf{v}_{k'} \rangle = 1 - \gamma$. Let $\boldsymbol{\xi} \in \mathbb{R}^{K_2}, \boldsymbol{\zeta} \in \mathbb{R}^{K_3}$
 78 be unique random features satisfying $\|\boldsymbol{\xi}\|_2, \|\boldsymbol{\zeta}\|_2 \leq R$ and are mutually independent given \mathbf{z} . The
 79 image-text pair is generated as

$$\mathbf{x} = \mathbf{G} \begin{bmatrix} \mathbf{z} \\ \boldsymbol{\xi} \end{bmatrix} = \mathbf{G}_1 \mathbf{z} + \mathbf{G}_2 \boldsymbol{\xi}, \quad \mathbf{y} = \mathbf{H} \begin{bmatrix} \mathbf{z} \\ \boldsymbol{\zeta} \end{bmatrix} = \mathbf{H}_1 \mathbf{z} + \mathbf{H}_2 \boldsymbol{\zeta},$$

80 where $\mathbf{G} \in \mathbb{R}^{d_1 \times (K_1 + K_2)}$ is the image dictionary with full rank $(K_1 + K_2)$, $\mathbf{H} \in \mathbb{R}^{d_2 \times (K_1 + K_3)}$ is
 81 the text dictionary with full rank $(K_1 + K_3)$.

82 We verify Assumptions 3.1 for the specified distribution in Appendix H. The following theorem gives
 83 convergence guarantees for CLIP and provides the upper bound of its zero-shot transfer error.

84 **Theorem 3.4.** For sufficiently large n , set the learning rate $\eta = O(\epsilon \tau^2 \|\mathbf{G}\|^{-2} \|\mathbf{H}\|_2^{-2} (1 + R)^{-4})$,
 85 gradient descent can find $\widehat{\mathbf{W}}$ within $4\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 / (\eta \epsilon)$ iterations such that $L_{\mathcal{D}^B}(\widehat{f}, \tau) \leq$
 86 $L_{\mathcal{D}^B}(f^*, \tau) + \epsilon$ where $\widehat{f} = \langle \widehat{\mathbf{W}}\mathbf{x}, \mathbf{y} \rangle$. In addition, the top- r zero-shot transfer error is bounded by
 87 $\epsilon' / \log(1 + r)$, where $\epsilon' = (C_B + 2) \cdot \left[\epsilon + 2B \exp(-\gamma/\tau) \right]$ and $C_B = \widetilde{O}(K/B)$.

88 **Square Loss Fails Zero-Shot Learning.** Suppose square loss $\mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbf{y}\|_2^2]$ is used to learn the
 89 embedding \mathbf{g} . We find that even if we can train with population risk and get the Bayesian optimal
 90 predictor, the learned representation \mathbf{g} is not suitable for the zero-shot transfer. We consider the data
 91 introduced in Definition 3.3 for the following.

92 **Theorem 3.5.** The Bayesian optimal representation \mathbf{g} is $\mathbf{g}(\mathbf{x}) = \mathbf{H} \begin{bmatrix} \mathbf{z} \\ \mathbb{E}[\boldsymbol{\zeta}|\mathbf{z}] \end{bmatrix}$.

93 The following corollary formally states the negative result.

94 **Corollary 3.6.** For the distribution in Definition 3.3 with $\mathbf{H} = \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix}$, margin $\gamma < 1/3$, text unique
 95 feature $\boldsymbol{\zeta} \in \mathbb{R}^{K_3}$ drawn from $\{\mathbf{e}_1, \mathbf{e}_2\}$ with probability $1/3, 2/3$ respectively. Then, the zero-shot
 96 top-1 error is at least $1/(3K)$ under various similarity scores, including cosine similarity.

97 **Remark 3.7.** By Theorem 3.4, we can achieve arbitrarily small top-1 error by CLIP as long as ϵ and
 98 τ are sufficiently small. However, for the representation learned from the square loss, the top-1 error
 99 is at least a constant even if we can achieve the Bayesian optimal predictor.

100 4 Learn Better Representation via Regularization

101 In Corollary 3.2, we know that CLIP can achieve a small error for zero-shot transfer tasks. In this
 102 section, we investigate how large the margin can be achieved between different features \mathbf{z} 's. Under
 103 the same condition of Corollary 3.2, we present the following corollary.

104 **Corollary 4.1.** Suppose the result of Theorem G.1 holds for the learned similarity function \widehat{f} . We
 105 calculate the similarity score $\widehat{f}(\mathbf{x}, \mathbf{y}_k)$ for all $k \in [K]$. Then with probability at least $1 - 4\epsilon'$, the
 106 top-1 result gives the correct answer with a margin τ .

107 Here, the margin depends on the temperature parameter τ . Note that we only achieve the margin
 108 with τ instead of γ guaranteed in the Assumption 3.1. This indicates a theoretical gap in the learned
 109 margin.

110 **Theorem 4.2.** Under the same condition as Theorem 3.4, there exists a special case with initialization
 111 $\mathbf{W}^{(0)}$, such that when we train the model with polynomial iterations $T = \text{poly}(\eta^{-1}, \epsilon, d_1, d_2)$, with
 112 probability at least 0.99, the top-1 result can only give the correct answer with a margin $\widetilde{O}(\tau)$.

113 Such a phenomenon also exists in real data: the margin will decrease when temperature τ decreases
 114 (see Figure 1). To obtain a larger margin, we propose to use the following regularization,

$$R(f) = -\frac{1}{|S|} \sum_{(\mathbf{x}, \mathbf{y}) \in S} f(\mathbf{x}, \mathbf{y}).$$

115 The following theorem shows that the regularization can improve the margin.

116 **Theorem 4.3.** Under the same condition as Theorem 4.2, with sufficiently small τ and appropriately
 117 chosen λ , within polynomial iterations $T = \text{poly}(\eta^{-1}, \epsilon, d_1, d_2)$, we can find a score function \widehat{f} with
 118 large margin. In particular, with a probability of at least 0.99, the top-1 result gives the correct label
 119 with a margin $\widetilde{\Omega}(\gamma)$.

120 5 Experiments

121 **Datasets.** For performance evaluation, we consider Conceptual Captions 3M (CC3M) (Sharma
 122 et al., 2018) and MSCOCO (Chen et al., 2015) as the pretraining datasets, in alignment with prior
 123 literature (Li et al., 2022; Goel et al., 2022).

124 **Architectures.** We consider the same setting for experiments on all baseline CLIP-objectives.
 125 Following the original CLIP paper, we employ ResNet (He et al., 2016) as the image encoder and the
 126 Transformer architecture (Vaswani et al., 2017) as the text encoder. We use pre-trained weights for
 127 both encoders. Detailed hyperparameters and additional experiments are presented in Appendix E.

128 5.1 Effect of Temperature on Margin

129 In support of our theoretical discussions in Corollary 4.1 and Theorem 4.2 that find the positive
 130 correlation between the margin and the temperature parameter, we conduct real data experiments. In
 131 Figure 1, we examine the margin distribution of CLIP models trained at varying temperatures. The
 132 margin is considered as the difference between a diagonal value and an off-diagonal value within
 133 a batch: $f(\mathbf{x}_i, \mathbf{y}_i) - f(\mathbf{x}_j, \mathbf{y}_i)$ and $f(\mathbf{x}_i, \mathbf{y}_i) - f(\mathbf{x}_i, \mathbf{y}_j)$. We collect the margins of untrained and
 134 trained CLIP models on all batches within the MSCOCO training dataset.

135 As depicted in Figure 1, a CLIP model with random initialization at the projection layers has margins
 136 normally distributed near zero, whereas trained models exhibit positive margins, signifying successful
 137 training. Furthermore, we consider CLIP models trained at fixed temperature values of 0.07 and
 138 0.01. As observed in the figure, the margin distribution shifts to the left as temperature τ decreases,
 139 suggesting that an extremely small τ leads to small margins, aligning with the results in Corollary 4.1.

140 5.2 Zero-shot Transfer

141 To confirm Theorem 4.3, we investigate the advantages of
 142 incorporating our regularization term during training by
 143 evaluating zero-shot transfer accuracy and linear probing
 144 on various datasets. We consider the following training
 145 objectives when adding our regularization: (1) the original
 146 CLIP (Radford et al., 2021), and (2) CyCLIP (Goel et al.,
 147 2022) with cross-modal and in-modal consistency regular-
 148 izations, adopting the same hyperparameters for the reg-
 149 ularizations as outlined in Goel et al. (2022). All models
 150 are trained on CC3M using the same model architecture,
 151 batch size, and optimizer settings. Further experimental
 152 details are provided in Appendix E.

153 In Table 1, we present the zero-shot test accuracy of CLIP
 154 models trained with the original CLIP objective and the
 155 CyCLIP objective. Firstly, we demonstrate the model’s
 156 performance when training solely on the regularization
 157 term (L2) and compare to that of the CLIP objective. In alignment with our Corollary 3.6, we can
 158 observe on real data that training exclusively on the L2 objective leads to a large error and even
 159 random guessing on the zero-shot datasets. Combining with our theoretical analysis, we show that a
 160 naive square loss fails to learn transferable representations. In the context of multi-modal learning,
 161 contrastive loss is important. Moreover, confirming our result from Theorem 4.3, incorporating
 162 the regularization term into the contrastive objective effectively enhances performance across the
 163 majority of zero-shot transfer tasks. It improves over the baseline on 5 out of 6 datasets by a good
 164 margin. The best performance achieved by adding regularization to the CLIP objective outperforms
 165 its original objective by 3.62% on CIFAR10 and by 2.06% on average of all datasets.

Table 1: Zero-shot top-1 accuracy (%). Notably, adding the regularization term successfully improves the baselines on 5 out of the 6 datasets.

	CIFAR10	CIFAR100	STL10	Food101	ImageNetV2	DTD	Average
Reg	10.04	1.05	9.95	1.08	0.11	2.07	3.47
CLIP	63.85	31.17	90.35	8.39	20.24	21.22	39.20
CyCLIP	60.71	28.87	89.98	9.72	19.66	20.21	38.19
CLIP+Reg	67.47	33.33	92.64	12.14	22.36	19.63	41.26

166 6 Conclusion

167 In this paper, we rigorously investigated the theoretical underpinnings of transferable representation
 168 learning in CLIP. We provided insights through specific cases and corroborated our theory with
 169 empirical evidence. Lastly, we proposed a regularization term grounded in our theoretical findings to
 170 enhance CLIP’s performance in zero-shot transfer.

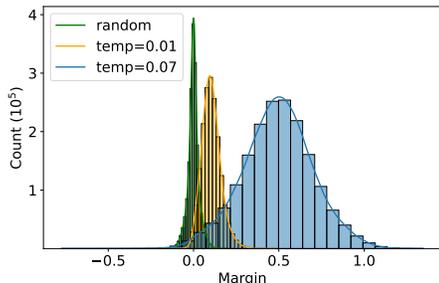


Figure 1: The distribution of the margins with regard to CLIP models trained at different temperature values. Margin is computed within each batch of the data.

References

- 171
172 ALAYRAC, J.-B., DONAHUE, J., LUC, P., MIECH, A., BARR, I., HASSON, Y., LENC, K., MENSCH,
173 A., MILLICAN, K., REYNOLDS, M. ET AL. (2022). Flamingo: a visual language model for
174 few-shot learning. *Advances in Neural Information Processing Systems* **35** 23716–23736.
- 175 BARTLETT, P. L. and MENDELSON, S. (2002). Rademacher and Gaussian complexities: Risk
176 bounds and structural results. *Journal of Machine Learning Research* **3** 463–482.
- 177 BOSSARD, L., GUILLAUMIN, M. and VAN GOOL, L. (2014). Food-101—mining discriminative
178 components with random forests. In *Computer Vision—ECCV 2014: 13th European Conference,*
179 *Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VI 13*. Springer.
- 180 CHEN, X., FANG, H., LIN, T.-Y., VEDANTAM, R., GUPTA, S., DOLLÁR, P. and ZITNICK,
181 C. L. (2015). Microsoft coco captions: Data collection and evaluation server. *arXiv preprint*
182 *arXiv:1504.00325* .
- 183 CIMPOI, M., MAJI, S., KOKKINOS, I., MOHAMED, S. and VEDALDI, A. (2014). Describing
184 textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern*
185 *recognition*.
- 186 COATES, A., NG, A. and LEE, H. (2011). An analysis of single-layer networks in unsupervised
187 feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence*
188 *and statistics*. JMLR Workshop and Conference Proceedings.
- 189 DESAI, K. and JOHNSON, J. (2021). Virtex: Learning visual representations from textual annotations.
190 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*.
- 191 FUKUMIZU, K., BACH, F. R. and JORDAN, M. I. (2004). Dimensionality reduction for supervised
192 learning with reproducing kernel hilbert spaces. *Journal of Machine Learning Research* **5** 73–99.
- 193 FUKUMIZU, K., BACH, F. R. and JORDAN, M. I. (2009). Kernel dimension reduction in regression .
- 194 GAO, Y., LIU, J., XU, Z., ZHANG, J., LI, K., JI, R. and SHEN, C. (2022). Pyramidclip: Hierar-
195 chical feature alignment for vision-language model pretraining. *Advances in Neural Information*
196 *Processing Systems* **35** 35959–35970.
- 197 GOEL, S., BANSAL, H., BHATIA, S., ROSSI, R., VINAY, V. and GROVER, A. (2022). Cyclip: Cyclic
198 contrastive language-image pretraining. *Advances in Neural Information Processing Systems* **35**
199 6704–6719.
- 200 GOMEZ, L., PATEL, Y., RUSINOL, M., KARATZAS, D. and JAWAHAR, C. (2017). Self-supervised
201 learning of visual features through embedding images into text topic spaces. In *Proceedings of the*
202 *IEEE conference on computer vision and pattern recognition*.
- 203 HAOCHEN, J. Z., WEI, C., GAIDON, A. and MA, T. (2021). Provable guarantees for self-supervised
204 deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*
205 **34**.
- 206 HE, K., ZHANG, X., REN, S. and SUN, J. (2016). Deep residual learning for image recognition. In
207 *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- 208 HUANG, Y., DU, C., XUE, Z., CHEN, X., ZHAO, H. and HUANG, L. (2021). What makes multi-
209 modal learning better than single (provably). *Advances in Neural Information Processing Systems*
210 **34** 10944–10956.
- 211 JIA, C., YANG, Y., XIA, Y., CHEN, Y.-T., PAREKH, Z., PHAM, H., LE, Q., SUNG, Y.-H., LI, Z.
212 and DUERIG, T. (2021). Scaling up visual and vision-language representation learning with noisy
213 text supervision. In *International Conference on Machine Learning*. PMLR.
- 214 KRIZHEVSKY, A. (2009). Learning multiple layers of features from tiny images. Tech. rep.
- 215 LEE, J. D., LEI, Q., SAUNSHI, N. and ZHUO, J. (2020). Predicting what you already know helps:
216 Provable self-supervised learning. *arXiv preprint arXiv:2008.01064* .

- 217 LEE, J. D., LEI, Q., SAUNSHI, N. and ZHUO, J. (2021). Predicting what you already know
218 helps: Provable self-supervised learning. *Advances in Neural Information Processing Systems* **34**
219 309–323.
- 220 LEI BA, J., SWERSKY, K., FIDLER, S. ET AL. (2015). Predicting deep zero-shot convolutional
221 neural networks using textual descriptions. In *Proceedings of the IEEE international conference*
222 *on computer vision*.
- 223 LI, Y., LIANG, F., ZHAO, L., CUI, Y., OUYANG, W., SHAO, J., YU, F. and YAN, J. (2022).
224 Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm.
225 In *International Conference on Learning Representations*.
- 226 LIANG, P. P., DENG, Z., MA, M., ZOU, J., MORENCY, L.-P. and SALAKHUTDINOV, R.
227 (2023). Factorized contrastive learning: Going beyond multi-view redundancy. *arXiv preprint*
228 *arXiv:2306.05268*.
- 229 MITROVIC, J., MCWILLIAMS, B., WALKER, J., BUESING, L. and BLUNDELL, C. (2020). Repre-
230 sentation learning via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922*.
- 231 MU, N., KIRILLOV, A., WAGNER, D. and XIE, S. (2022). Slip: Self-supervision meets language-
232 image pre-training. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel,*
233 *October 23–27, 2022, Proceedings, Part XXVI*. Springer.
- 234 NAKADA, R., GULLUK, H. I., DENG, Z., JI, W., ZOU, J. and ZHANG, L. (2023). Understanding
235 multimodal contrastive learning and incorporating unpaired data. In *International Conference on*
236 *Artificial Intelligence and Statistics*. PMLR.
- 237 NILSBACK, M.-E. and ZISSERMAN, A. (2008). Automated flower classification over a large number
238 of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*.
239 IEEE.
- 240 PARKHI, O. M., VEDALDI, A., ZISSERMAN, A. and JAWAHAR, C. (2012). Cats and dogs. In *2012*
241 *IEEE conference on computer vision and pattern recognition*. IEEE.
- 242 PHAM, H., DAI, Z., GHIASI, G., KAWAGUCHI, K., LIU, H., YU, A. W., YU, J., CHEN, Y.-T.,
243 LUONG, M.-T., WU, Y. ET AL. (2021). Combined scaling for zero-shot transfer learning. *arXiv*
244 *preprint arXiv:2111.10050*.
- 245 RADFORD, A., KIM, J. W., HALLACY, C., RAMESH, A., GOH, G., AGARWAL, S., SASTRY, G.,
246 ASKELL, A., MISHKIN, P., CLARK, J. ET AL. (2021). Learning transferable visual models from
247 natural language supervision. In *International Conference on Machine Learning*. PMLR.
- 248 RECHT, B., ROELOFS, R., SCHMIDT, L. and SHANKAR, V. (2019). Do imagenet classifiers
249 generalize to imagenet? In *International conference on machine learning*. PMLR.
- 250 SAITO, K., SOHN, K., ZHANG, X., LI, C.-L., LEE, C.-Y., SAENKO, K. and PFISTER, T. (2022).
251 Prefix conditioning unifies language and label supervision. *arXiv preprint arXiv:2206.01125*.
- 252 SARIYILDIZ, M. B., PEREZ, J. and LARLUS, D. (2020). Learning visual representations with
253 caption annotations. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK,*
254 *August 23–28, 2020, Proceedings, Part VIII 16*. Springer.
- 255 SAUNSHI, N., ASH, J., GOEL, S., MISRA, D., ZHANG, C., ARORA, S., KAKADE, S. and
256 KRISHNAMURTHY, A. (2022). Understanding contrastive learning requires incorporating inductive
257 biases. *arXiv preprint arXiv:2202.14037*.
- 258 SAUNSHI, N., PLEVRAKIS, O., ARORA, S., KHODAK, M. and KHANDEPARKAR, H. (2019). A the-
259 oretical analysis of contrastive unsupervised representation learning. In *International Conference*
260 *on Machine Learning*. PMLR.
- 261 SHARIATNIA, M. M. (2021). Simple CLIP.
- 262 SHARMA, P., DING, N., GOODMAN, S. and SORICUT, R. (2018). Conceptual captions: A cleaned,
263 hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th*
264 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

- 265 THOMEE, B., SHAMMA, D. A., FRIEDLAND, G., ELIZALDE, B., NI, K., POLAND, D., BORTH, D.
266 and LI, L.-J. (2016). Yfcc100m: The new data in multimedia research. *Communications of the*
267 *ACM* **59** 64–73.
- 268 TIAN, Y., YU, L., CHEN, X. and GANGULI, S. (2020). Understanding self-supervised learning with
269 dual deep networks. *arXiv preprint arXiv:2010.00578* .
- 270 TOSH, C., KRISHNAMURTHY, A. and HSU, D. (2021a). Contrastive estimation reveals topic
271 posterior information to linear models. *Journal of Machine Learning Research* **22** 1–31.
- 272 TOSH, C., KRISHNAMURTHY, A. and HSU, D. (2021b). Contrastive estimation reveals topic
273 posterior information to linear models. *Journal of Machine Learning Research* **22** 1–31.
- 274 TSAI, Y.-H. H., WU, Y., SALAKHUTDINOV, R. and MORENCY, L.-P. (2020). Demystifying
275 self-supervised learning: An information-theoretical framework. *arXiv preprint arXiv:2006.05576*
276 .
- 277 VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER,
278 Ł. and POLOSUKHIN, I. (2017). Attention is all you need. In *Advances in neural information*
279 *processing systems*.
- 280 WANG, T. and ISOLA, P. (2020). Understanding contrastive representation learning through align-
281 ment and uniformity on the hypersphere. In *International Conference on Machine Learning*.
282 PMLR.
- 283 WEN, Z. and LI, Y. (2021). Toward understanding the feature learning process of self-supervised
284 contrastive learning. In *International Conference on Machine Learning*. PMLR.
- 285 WIGHTMAN, R. (2019). Pytorch image models. [https://github.com/rwightman/
286 pytorch-image-models](https://github.com/rwightman/pytorch-image-models).
- 287 WOLF, T., DEBUT, L., SANH, V., CHAUMOND, J., DELANGUE, C., MOI, A., CISTAC, P., RAULT,
288 T., LOUF, R., FUNTOWICZ, M., DAVISON, J., SHLEIFER, S., VON PLATEN, P., MA, C., JERNITE,
289 Y., PLU, J., XU, C., SCAO, T. L., GUGGER, S., DRAME, M., LHOEST, Q. and RUSH, A. M.
290 (2020). Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020*
291 *Conference on Empirical Methods in Natural Language Processing: System Demonstrations*.
292 Association for Computational Linguistics, Online.
- 293 WU, K., ZHANG, J., PENG, H., LIU, M., XIAO, B., FU, J. and YUAN, L. (2022). Tinyvit: Fast
294 pretraining distillation for small vision transformers. In *European Conference on Computer Vision*.
295 Springer.
- 296 YAO, L., HUANG, R., HOU, L., LU, G., NIU, M., XU, H., LIANG, X., LI, Z., JIANG, X. and
297 XU, C. (2022). FILIP: Fine-grained interactive language-image pre-training. In *International*
298 *Conference on Learning Representations*.
- 299 ZADEH, A., LIANG, P. P. and MORENCY, L.-P. (2020). Foundations of multimodal co-learning.
300 *Information Fusion* **64** 188–193.
- 301 ZHAI, X., WANG, X., MUSTAFA, B., STEINER, A., KEYSERS, D., KOLESNIKOV, A. and BEYER,
302 L. (2022). Lit: Zero-shot transfer with locked-image text tuning. In *Proceedings of the IEEE/CVF*
303 *Conference on Computer Vision and Pattern Recognition*.
- 304 ZHANG, T. (2002). Covering number bounds of certain regularized linear function classes. *Journal*
305 *of Machine Learning Research* **2** 527–550.
- 306 ZHANG, Y., JIANG, H., MIURA, Y., MANNING, C. D. and LANGLOTZ, C. P. (2022). Contrastive
307 learning of medical visual representations from paired images and text. In *Machine Learning for*
308 *Healthcare Conference*. PMLR.

309 A Related Work

310 **Vision-Language Pre-Training.** While labeled data are expensive and relatively scarce, images
 311 paired with text descriptions are available in much larger volumes (Thomee et al., 2016). Conse-
 312 quently, numerous studies (Gomez et al., 2017; Sariyildiz et al., 2020; Desai and Johnson, 2021;
 313 Zhang et al., 2022; Liang et al., 2023) have focused on leveraging free-form natural language su-
 314 pervision to learn visual representations. Recently, CLIP (Radford et al., 2021) and ALIGN (Jia
 315 et al., 2021) have emerged as prominent works extending contrastive learning to the vision-language
 316 pre-training framework. Built upon CLIP’s success, several studies (Pham et al., 2021; Gao et al.,
 317 2022; Saito et al., 2022) have refined CLIP’s contrastive methodology to better learn from web-scale
 318 image-text data. Notably, FILIP (Yao et al., 2022) introduces a fine-grained contrastive loss tailored
 319 for transformer architectures. DeCLIP (Li et al., 2022) and SLIP (Mu et al., 2022) additionally
 320 incorporate single-modality self-supervised learning. CyCLIP (Goel et al., 2022) introduces two
 321 regularizing terms enforcing cross-modal and in-modal consistency. LiT (Zhai et al., 2022) and
 322 Flamingo (Alayrac et al., 2022) consider training from pre-trained single-modality models. In our
 323 empirical validation of theoretical findings, we employ the same setting and train from pre-trained
 324 image and text encoders.

325 **Theory of self-supervised learning.** To understand self-supervised learning, numerous studies have
 326 been conducted, particularly focusing on *unimodal* contrastive learning, a widely used self-supervised
 327 learning approach rooted in data augmentation (Saunshi et al., 2019; Tsai et al., 2020; Mitrovic
 328 et al., 2020; Tian et al., 2020; Wang and Isola, 2020; Tosh et al., 2021a,b; HaoChen et al., 2021;
 329 Wen and Li, 2021; Saunshi et al., 2022). In multimodal learning, theoretical explanation has been
 330 explored in several studies (Zadeh et al., 2020; Huang et al., 2021; Lee et al., 2020; Nakada et al.,
 331 2023). These works have established that multimodal learning can surpass unimodal learning in
 332 terms of performance. For instance, Lee et al. (2020) employed square loss prediction to learn
 333 image representations under certain conditional independence assumptions, offering generalization
 334 performance guarantees. Meanwhile, Nakada et al. (2023) examined CLIP within specific linear
 335 representation settings and emphasized its correlation with singular value decomposition (SVD). We
 336 note that, these related works have not considered the zero-shot transfer mechanism and thus can’t
 337 adequately explain the zero-shot transfer capability of CLIP.

338 B Preliminaries

339 **Notation.** We use lowercase letters, lowercase boldface letters, and uppercase boldface letters
 340 to denote scalars, vectors, and matrices, respectively. For a vector \mathbf{x} , we use $\|\mathbf{x}\|_2$ to denote its
 341 Euclidean norm. For a matrix \mathbf{W} , we use $\|\mathbf{W}\|_F$ to denote its Frobenius norm. Given two sequences
 342 $\{x_n\}$ and $\{y_n\}$, we denote $x_n = \mathcal{O}(y_n)$ if $|x_n| \leq C_1|y_n|$ for some absolute positive constant
 343 C_1 , $x_n = \Omega(y_n)$ if $|x_n| \geq C_2|y_n|$ for some absolute positive constant C_2 , and $x_n = \Theta(y_n)$ if
 344 $C_3|y_n| \leq |x_n| \leq C_4|y_n|$ for some absolute constants $C_3, C_4 > 0$. We also use $\tilde{\mathcal{O}}(\cdot)$ to hide
 345 logarithmic factors of d in $\mathcal{O}(\cdot)$. Additionally, we denote $x_n = \text{poly}(y_n)$ if $x_n = \mathcal{O}(y_n^D)$ for some
 346 positive constant D , and $x_n = \text{polylog}(y_n)$ if $x_n = \text{poly}(\log(y_n))$. We also denote by $x_n = o(y_n)$
 347 if $\lim_{n \rightarrow \infty} x_n/y_n = 0$. Finally we use $[N]$ to denote the index set $\{1, \dots, N\}$. In the function space,
 348 let $B_r(f)$ denote the ball of radius r centered at f , with the metrics $\|\cdot\|_\infty$. A set C is the covering
 349 of function class \mathcal{F} with radius r , if and only if $\mathcal{F} \subseteq \cup_{f \in C} B_r(f)$. The covering number of \mathcal{F} with
 350 radius r is the minimum cardinality of any covering of \mathcal{F} , denoted as $\mathcal{N}(\mathcal{F}, r)$.

351 **Loss function of CLIP.**

$$\begin{aligned}
 L_{S'}(f, \tau) &= \frac{1}{B} \sum_{i \in S'} -\log \left(\frac{\exp(f(\mathbf{x}_i, \mathbf{y}_i)/\tau)}{\sum_{j \in S'} \exp(f(\mathbf{x}_j, \mathbf{y}_i)/\tau)} \right) + \frac{1}{B} \sum_{i \in S'} -\log \left(\frac{\exp(f(\mathbf{x}_i, \mathbf{y}_i)/\tau)}{\sum_{j \in S'} \exp(f(\mathbf{x}_i, \mathbf{y}_j)/\tau)} \right) \\
 &= \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f(\mathbf{x}_j, \mathbf{y}_i) - f(\mathbf{x}_i, \mathbf{y}_i)]/\tau) \right) \\
 &\quad + \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f(\mathbf{x}_i, \mathbf{y}_j) - f(\mathbf{x}_i, \mathbf{y}_i)]/\tau) \right), \tag{B.1}
 \end{aligned}$$

352 where $\tau > 0$ is a temperature parameter. The training loss $L_{S'}$ over a single epoch can be viewed as
 353 the empirical version of (2.1).

354 **Remark B.1.** In Assumption 2.1, the assumption of conditional independence is frequently made
 355 in the analysis of self-supervised learning (Saunshi et al., 2019; Lee et al., 2021) and dimension

356 reduction algorithms (Fukumizu et al., 2004, 2009). Under the premise that \mathbf{x}, \mathbf{y} are conditionally
 357 independent (CI) given \mathbf{z} , it can be posited that any additional patterns found within $\mathbf{x}|\mathbf{z}$ and $\mathbf{y}|\mathbf{z}$
 358 should be interpreted as unique features. Notably, in the absence of discrete and sparse constraints,
 359 a suitable \mathbf{z} can always be found, given that one could simply assign $\mathbf{z} = \mathbf{x}$ or $\mathbf{z} = \mathbf{y}$. From
 360 the generative model’s point of view, Assumption 2.1 naively holds when the data are from some
 361 generator with $\mathbf{x} = T_1(\mathbf{z}, \xi)$ and $\mathbf{y} = T_2(\mathbf{z}, \zeta)$ where $\xi \perp \zeta | \mathbf{z}$.

362 Given the population loss in (2.1), the following theorem shows that the empirical loss $\widehat{\mathbb{E}}_S(f, \tau) :=$
 363 $(1/n) \sum_{k \in [n]} L_{S_k}(f, \tau)$ concentrates on the population loss when n is large enough.

364 **Theorem B.2.** Suppose $\delta \in (0, 1)$ and $n \geq (8\tau^{-1}\epsilon^{-2}M \log B) \log(2\mathcal{N}(\mathcal{F}, \epsilon/8M)/\delta)$, then with
 365 probability at least $1 - \delta$, we have

$$|\widehat{L}_S(f, \tau) - L_{\mathcal{D}^B}(f, \tau)| \leq \epsilon$$

366 for all function $f \in \mathcal{F}$ and $|f| \leq M$, where $\mathcal{N}(\mathcal{F}, \epsilon)$ is the covering number of \mathcal{F} .

367 Theorem B.2 shows that the generalization gap $|\widehat{L}_S(f, \tau) - L_{\mathcal{D}^B}(f, \tau)|$ approaches zero as the number
 368 of batches n increase. In practice, the batch size is limited by the GPU’s memory and is smaller than
 369 the number of batches (or the number of training examples). Therefore, instead of letting the batch
 370 size B go to infinity like in prior studies (Wang and Isola, 2020; Pham et al., 2021), we keep the batch
 371 size B as a constant in (2.1) and Theorem B.2 to enable the analysis of CLIP even for small batches.
 372 Pham et al. (2021) also provided the generalization gap for CLIP. However, their result is for $B \rightarrow \infty$
 373 and a loss function without the log term, i.e., $\exp(f(\mathbf{x}_i, \mathbf{y}_i)/\tau) / \left(\sum_{j \in S'} \exp(f(\mathbf{x}_j, \mathbf{y}_j)/\tau) \right)$.

374 C Discussion on the Margin

375 In Assumption 3.1, we introduce the (α, β, γ) completeness. In this section, we will discuss how to
 376 verify the assumption and formally measure the quality of the learned function.

377 Sample two independent tuple $(\mathbf{x}, \mathbf{y}, \mathbf{z}) \sim \mathcal{D}_{\mathbf{x} \times \mathbf{y} \times \mathbf{z}}$ and $(\mathbf{x}', \mathbf{y}', \mathbf{z}') \sim \mathcal{D}_{\mathbf{x} \times \mathbf{y} \times \mathbf{z}}$, we introduce a
 378 measure as follows.

$$\alpha_\gamma = \mathbb{P}(\mathbf{z} \neq \mathbf{z}', f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}') \leq \gamma) + \mathbb{P}(\mathbf{z} \neq \mathbf{z}', f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}', \mathbf{y}) \leq \gamma)$$

379 By Assumption 3.1, we know that we want to find a large γ with small α_γ .

380 However, in real applications, we can access $\mathcal{D}_{\mathbf{x} \times \mathbf{y}}$ but have little knowledge of the model \mathcal{V} and the
 381 latent variable \mathbf{z} . Thus, we introduce another measure $\widehat{\alpha}_\gamma$ as follows,

$$\widehat{\alpha}_\gamma = \mathbb{P}(f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}') \leq \gamma) + \mathbb{P}(f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}', \mathbf{y}) \leq \gamma)$$

382 $\widehat{\alpha}_\gamma$ differs from the α_γ since we didn’t extinguish different classes in the probability. Therefore we
 383 can easily calculate $\widehat{\alpha}_\gamma$ without observe \mathbf{z} . Besides, we have the following upper and low bounds,
 384 which show that $\widehat{\alpha}_\gamma$ can approximate α_γ .

385 **Theorem C.1.** Let $\gamma \geq 0$, then we have that

$$\widehat{\alpha}_\gamma \geq \alpha_\gamma \geq \widehat{\alpha}_\gamma - \sum_{k \in [K]} p_k^2.$$

386 where p_k is the probability of the classes in Assumption 2.1. Besides the second inequality become
 387 exact equality for $\gamma = 0$.

Proof.

$$\begin{aligned} \widehat{\alpha}_\gamma &= \mathbb{P}(f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}') \leq \gamma) + \mathbb{P}(f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}', \mathbf{y}) \leq \gamma) \\ &= \underbrace{\mathbb{P}(\mathbf{z} \neq \mathbf{z}', f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}') \leq \gamma) + \mathbb{P}(\mathbf{z} \neq \mathbf{z}', f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}', \mathbf{y}) \leq \gamma)}_{=\alpha_\gamma} \\ &\quad + \underbrace{\mathbb{P}(\mathbf{z} = \mathbf{z}', f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}') \leq \gamma) + \mathbb{P}(\mathbf{z} = \mathbf{z}', f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}', \mathbf{y}) \leq \gamma)}_{\text{ApproximateError}}. \end{aligned}$$

388 The Approximate Error has a naive lower bound of 0 and we can upper bound it as follows

$$\begin{aligned} \mathbb{P}(\mathbf{z} = \mathbf{z}', f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}') \leq \gamma) &= \mathbb{P}(f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}') \leq \gamma | \mathbf{z} = \mathbf{z}') \cdot \mathbb{P}(\mathbf{z} = \mathbf{z}') \\ &\leq \mathbb{P}(f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}') \leq 0 | \mathbf{z} = \mathbf{z}') \cdot \mathbb{P}(\mathbf{z} = \mathbf{z}') \\ &= 1/2 \sum_{k \in [K]} p_k^2. \end{aligned}$$

389 were the the inequality is due to fact that $\gamma \geq 0$ and the last equality is because \mathbf{y}' and \mathbf{y} are symmetric
390 give $\mathbf{z} = \mathbf{z}'$. Finally, the inequality is an exact equality for $\gamma = 0$. \square

391 By Theorem C.1, α_γ and $\hat{\alpha}_\gamma$ are close to each other if $\max_{k \in [K]} p_k$ is small, since

$$\sum_{k \in [K]} p_k^2 \leq \sum_{k \in [K]} p_k \cdot \max_{k \in [K]} p_k = \max_{k \in [K]} p_k \cdot \left(\sum_{k \in [K]} p_k \right) = \max_{k \in [K]} p_k.$$

392 **Relation with the Figure 2:** $\hat{\alpha}_\gamma$ has a strong relationship with Figure 2, where we have plot the
393 distribution of $f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}')$ and $f(\mathbf{x}, \mathbf{y}) - f(\mathbf{x}', \mathbf{y})$. The figure can be viewed as the figure of
394 the probability density function and $\hat{\alpha}_\gamma$ can be viewed as the cumulative probability function which is
395 the integral of probability mass smaller than γ . From Figure 2, we can deduce that the CLIP learned
396 with regularization has consistently smaller $\hat{\alpha}_\gamma$ for all $\gamma \geq 0$.

397 D Discussion on the Trainable Temperature Parameter τ

398 This section considers the setting where the temperature τ is also trainable with the following loss.

$$\begin{aligned} L_{\mathcal{D}^B}(f, \tau) &= \mathbb{E} \left[\log \left(\sum_{t \in [B]} \exp([f(\mathbf{x}_1, \mathbf{y}_t) - f(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right] \\ &\quad + \mathbb{E} \left[\log \left(\sum_{t \in [B]} \exp([f(\mathbf{x}_t, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right]. \end{aligned}$$

399 Suppose τ is clipped to be within the range $[\tau_{\min}, \tau_{\max}]$, it is natural to assume that we can obtain
400 function \hat{f} with temperature $\hat{\tau} \in [\tau_{\min}, \tau_{\max}]$ such that

$$L_{\mathcal{D}^B}(\hat{f}, \hat{\tau}) \leq \min_{\tau \in [\tau_{\min}, \tau_{\max}]} L_{\mathcal{D}^B}(f^*, \tau) + \epsilon \quad (\text{D.1})$$

$$= L_{\mathcal{D}^B}(f^*, \hat{\tau}) + \epsilon - \left(L_{\mathcal{D}^B}(f^*, \hat{\tau}) - \min_{\tau \in [\tau_{\min}, \tau_{\max}]} L_{\mathcal{D}^B}(f^*, \tau) \right) \quad (\text{D.2})$$

$$= L_{\mathcal{D}^B}(f^*, \hat{\tau}) + \tilde{\epsilon} \quad (\text{D.3})$$

401 where $\tilde{\epsilon} = \epsilon - \left(L_{\mathcal{D}^B}(f^*, \hat{\tau}) - \min_{\tau \in [\tau_{\min}, \tau_{\max}]} L_{\mathcal{D}^B}(f^*, \tau) \right) \leq \epsilon$. Since $\tilde{\epsilon}$ is smaller than ϵ , we
402 can get smaller ϵ' in Theorem G.1, and thus get smaller top-r error in zero-shot transfer task by
403 Corollary 3.2. This observation implies that the representation $(\hat{f}, \hat{\tau})$ found by trainable temperature
404 can be better than the representation $(\hat{f}', \hat{\tau})$ found with fixed temperature $\hat{\tau}$.

405 E Additional Experiment Results

406 We consider the same model architecture as CLIP (Radford et al., 2021) and consider ResNet-
407 50 (He et al., 2016) as the image encoder and transformer (Vaswani et al., 2017) achitecture as the
408 text encoder. Specifically, we use pre-trained weights for the encoders for faster convergence in
409 training. We follow the code framework in Shariatnia (2021) and use pre-trained ResNet-50 from the
410 PyTorch Image Models library (Wightman, 2019) and pre-trained DistilBERT from the Huggingface
411 Transformers library (Wolf et al., 2020). We further have linear projection layers on both image and
412 text encoder as the same as in CLIP and consider embedding dimension of 512. As we are training at
413 small-scale data with pre-trained encoders, we follow Shariatnia (2021) and use AdamW optimizer
414 with learning rate 1e-4 on the image encoder, 1e-5 on the text encoder, and 1e-3 on the projection
415 layers, with weight decay coefficient 1e-3. Our code is provided anonymously on Github¹.

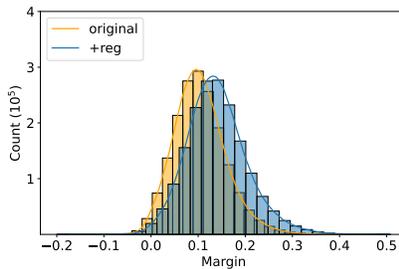
¹https://anonymous.4open.science/r/CLIP_theory-BC8F/README.md

416 **E.1 Effect of Temperature on Margin**

417 **Setup.** For lightweight exploration in section 5.1, we use the training dataset from MSCOCO (Chen
 418 et al., 2015) Image Captioning Task as the data for vision-language contrastive pre-training. Specifically,
 419 the dataset contains 82, 783 images where each image is coupled with 5 captions. We consider
 420 each image-caption pair as a data example in pre-training and therefore arrive at 413, 915 pre-training
 421 data pairs. We further randomly split the data to keep 20% of the data as validation set and stops
 422 training as the contrastive loss on validation data no longer decreases to avoid overfitting on the small
 423 dataset.

424 **Margin.** Given a training data batch, the margin is consider as the difference between a diagonal
 425 value and an off-diagonal value: $f(\mathbf{x}_i, \mathbf{y}_i) - f(\mathbf{x}_j, \mathbf{y}_i)$ and $f(\mathbf{x}_i, \mathbf{y}_i) - f(\mathbf{x}_i, \mathbf{y}_j)$. We consider CLIP
 426 models trained at fixed temperature $\tau = 0.07$ and $\tau = 0.01$. We note that 0.07 is the default value for
 427 τ to start training in CLIP and 0.01 is the clamping value (equivalently as the maximum logit scale of
 428 4.6052.) In Figure 1, we collected the margins from all batches of size 64 in the MSCOCO training
 429 data, where the data is randomly shuffled.

430 **Additional Experiments.** Here, we additionally compare the margin distribution of CLIP trained at
 431 temperature $\tau = 0.01$, without or with our regularization term. We could observe that the margin
 432 distribution shifts to the right with the regularization term, which alleviates the negative influence of
 an extremely small temperature value.



433 Figure 2: The distribution of the margins with regard to CLIP models trained $\tau = 0.01$ with or
 without regularization. Margin is computed within each batch of the data.

434 **E.2 Zero-shot Transfer and Linear Probing**

435 **Setup.** In the evaluation of zero-shot transfer and linear probing, we use CC3M (Sharma et al., 2018)
 436 as the pre-training dataset, which contains around 3, 318, 332 image-caption pairs gathered from the
 437 web. While some URLs are broken so that we cannot download the images, we eventually reached
 438 a pre-training dataset of 2, 786, 288 data pairs. When training CLIP models, we use the default
 439 coefficients of CyCLIP regularization terms of $\lambda_1 = 0.25$ and $\lambda_2 = 0.25$. For our regularization
 440 term, we use a coefficient of $\lambda = 0.1$. As in CLIP, we set the temperature τ from 0.07, equivalently
 441 having maximum logit scale at 2.6593. Lastly, we use a training batch size of 32 and trained for 8
 442 epochs in the results reported in section 5.2.

Table 2: Summary of datasets used for zero-shot transfer and linear probing.

Dataset	Classes	Class Description
CIFAR10	10	Categories of animals and vehicles
CIFAR100	100	Categories of objects including animals, foods, vehicles and people
STL10	10	Categories of animals and vehicles
Food101	101	Categories of foods/dishes
ImageNetV2	1000	Categories of objects including animals, foods, vehicles and people
DTD	47	Categories of textures
Flowers102	102	Categories of flower species
Oxford-IIIT Pet	37	Categories of cats and dogs

443 **Evaluations.** As similar in previous works (Radford et al., 2021; Yao et al., 2022; Mu et al., 2022;
 444 Goel et al., 2022), we consider the following image classification tasks for zero-shot transfer and
 445 linear probing: CIFAR10/100 (Krizhevsky, 2009), STL10 (Coates et al., 2011), Food101 (Bossard
 446 et al., 2014), ImageNetV2 (Recht et al., 2019), DTD (Describable Textures, Cimpoi et al. (2014)),
 447 Flowers102 (Nilsback and Zisserman, 2008) and Oxford-IIIT Pet (Parkhi et al., 2012). The dataset

448 statistics are reported in Table 2. For zero-shot transfer, we use the same prompt engineering and
 449 ensembling as the original CLIP and report the top-1 accuracy. For linear probing, as the same
 450 in CLIP, we train a logistic regression classifier on the image embeddings generated by the image
 451 encoder of pre-trained CLIP models on the training data from the considered datasets. The classifiers
 452 are all trained to convergence and we report the test accuracy on each of the test dataset of the tasks.
 453 We note that, due to the limitation of the training data CC3M, the zero-shot test accuracy of all
 454 CLIP-objectives on Flowers102 and Oxford-IIIT Pet are near random guesses. Therefore, we omit
 455 these datasets for zero-shot transfer.

456 **Additional Experiments.** In Table 3, we report the results of linear probing, where logistic regression
 457 classifiers are fitted to the embeddings learned by the image encoders of our compared models.
 458 This table offers an assessment of the visual representation learning for each training objective.
 459 Similarly supporting Corollary 3.6, training on the regularization term only results in learning bad
 460 representations that yield unsatisfactory performances on linear probing. Moreover, in alignment
 461 with Theorem 4.3, we observe that adding the regularization term consistently improves CLIP’s
 performance across various datasets by an average of 1.54%.

Table 3: Linear probing accuracy (%). All logistic regression models are trained till convergence. Adding our regularization term to CLIP provides decent improvements across all datasets. On CyCLIP, we also makes improvements on the majority of datasets.

	CIFAR10	CIFAR100	STL10	Food101	DTD	Flowers	OxfordPets	Average
Reg	14.09	2.17	17.86	1.73	3.40	2.18	4.12	6.51
CLIP	87.30	66.03	93.26	62.8	56.70	70.24	72.91	72.75
CyCLIP	86.31	63.93	93.69	61.57	56.86	70.56	70.46	71.91
CLIP+Reg	88.49	66.16	94.98	63.39	57.66	72.21	77.13	74.29

462 We additionally report the zero-shot transfer results of the original CLIP objective and adding our
 463 regularization term, on a different visual encoder architecture of TinyViT (Wu et al., 2022) with
 464 pre-trained weights from Huggingface.

Table 4: Zero-shot top-1 accuracy (%). Notably, adding the regularization term successfully improves the baselines on 5 out of the 6 datasets.

	CIFAR10	CIFAR100	STL10	Food101	ImageNetV2	DTD	Average
CLIP	52.02	15.57	81.89	7.92	16.91	11.80	31.02
CLIP+Reg	53.30	19.67	83.76	7.99	16.06	11.53	32.05

465

466 F Proof of Results in Section 2

467 *Proof of Theorem B.2.* We first prove that $L_{S'}(f, \tau)$ is upper bounded by $4M \log B/\tau$.

$$\begin{aligned}
 L_{S'}(f, \tau) &= \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp \left([f(\mathbf{x}_j, \mathbf{y}_i) - f(\mathbf{x}_i, \mathbf{y}_i)]/\tau \right) \right) \\
 &\quad + \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp \left([f(\mathbf{x}_i, \mathbf{y}_j) - f(\mathbf{x}_i, \mathbf{y}_i)]/\tau \right) \right) \\
 &\leq \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp(2M/\tau) \right) + \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp(2M/\tau) \right) \\
 &= 4M \log B/\tau. \tag{F.1}
 \end{aligned}$$

468 where the inequality is by the fact the $|f| \leq M$. On the other hand, we have that

$$\begin{aligned}
 L_{S'}(f, \tau) &= \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp \left([f(\mathbf{x}_j, \mathbf{y}_i) - f(\mathbf{x}_i, \mathbf{y}_i)]/\tau \right) \right) \\
 &\quad + \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp \left([f(\mathbf{x}_i, \mathbf{y}_j) - f(\mathbf{x}_i, \mathbf{y}_i)]/\tau \right) \right) \\
 &\geq \frac{2}{B} \sum_{i \in S'} \log \left(\exp \left([f(\mathbf{x}_i, \mathbf{y}_i) - f(\mathbf{x}_i, \mathbf{y}_i)]/\tau \right) \right) \\
 &\geq 0.
 \end{aligned}$$

469 where the inequality is because Exp function is greater than 0. Therefore we have proved that
 470 $L_{S'}(f, \tau) \in (0, 4M \log(B)/\tau]$. For all $f_1, f_2 \in \mathcal{F}$ and any batch S' with size B , we have that

$$\begin{aligned}
 L_{S'}(f_1, \tau) - L_{S'}(f_2, \tau) &= \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f_1(\mathbf{x}_j, \mathbf{y}_i) - f_1(\mathbf{x}_i, \mathbf{y}_i)]/\tau) \right) \\
 &\quad - \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f_2(\mathbf{x}_j, \mathbf{y}_i) - f_2(\mathbf{x}_i, \mathbf{y}_i)]/\tau) \right) \\
 &\quad + \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f_1(\mathbf{x}_i, \mathbf{y}_j) - f_1(\mathbf{x}_i, \mathbf{y}_i)]/\tau) \right) \\
 &\quad - \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f_2(\mathbf{x}_i, \mathbf{y}_j) - f_2(\mathbf{x}_i, \mathbf{y}_i)]/\tau) \right) \\
 &\leq \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f_1(\mathbf{x}_j, \mathbf{y}_i) - f_1(\mathbf{x}_i, \mathbf{y}_i)]/\tau) \right) \\
 &\quad - \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f_1(\mathbf{x}_j, \mathbf{y}_i) - f_1(\mathbf{x}_i, \mathbf{y}_i) - 2\|f_1 - f_2\|_\infty]/\tau) \right) \\
 &\quad + \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f_1(\mathbf{x}_i, \mathbf{y}_j) - f_1(\mathbf{x}_i, \mathbf{y}_i)]/\tau) \right) \\
 &\quad - \frac{1}{B} \sum_{i \in S'} \log \left(\sum_{j \in S'} \exp([f_1(\mathbf{x}_i, \mathbf{y}_j) - f_1(\mathbf{x}_i, \mathbf{y}_i) - 2\|f_1 - f_2\|_\infty]/\tau) \right) \\
 &= 4\|f_1 - f_2\|_\infty/\tau.
 \end{aligned}$$

471 Similarly, we can get another direction $L_{S'}(f_2, \tau) - L_{S'}(f_1, \tau) \leq 4\|f_1 - f_2\|_\infty/\tau$, which yields
 472 to $|L_{S'}(f_2, \tau) - L_{S'}(f_1, \tau)| \leq 4\|f_1 - f_2\|_\infty/\tau$. Taking the expectation gives that $|L_{\mathcal{D}^B}(f_2, \tau) -$
 473 $L_{\mathcal{D}^B}(f_1, \tau)| \leq 4\|f_1 - f_2\|_\infty/\tau$. By the definition of the covering set, the function class \mathcal{F} can
 474 be covered by K subsets $\mathcal{B}_1, \dots, \mathcal{B}_K$, that is $\mathcal{F} = \mathcal{B}_1 \cup \dots \cup \mathcal{B}_K$, where $K = \mathcal{N}(\mathcal{F}, \tau\epsilon/16)$ and
 475 $\mathcal{B}_1, \dots, \mathcal{B}_K$ are the balls of the radius $\tau \cdot \epsilon/16$ centered at f_1, \dots, f_K . Then we have that

$$\begin{aligned}
 &\mathbb{P}_{S \sim \mathcal{D}^n} \left[\sup_{f \in \mathcal{F}} |L_{\mathcal{D}^B}(f, \tau) - \widehat{L}_S(f, \tau)| \geq \epsilon \right] \\
 &\leq \sum_{k \in [K]} \mathbb{P}_{S \sim \mathcal{D}^n} \left[\sup_{f \in \mathcal{B}_k} |L_{\mathcal{D}^B}(f, \tau) - \widehat{L}_S(f, \tau)| \geq \epsilon \right] \\
 &\leq \sum_{k \in [K]} \mathbb{P}_{S \sim \mathcal{D}^n} \left[|L_{\mathcal{D}^B}(f_k, \tau) - \widehat{L}_S(f_k, \tau)| \geq \epsilon/2 \right] \\
 &= \sum_{k \in [K]} \mathbb{P}_{S \sim \mathcal{D}^n} \left[|L_{\mathcal{D}^B}(f_k, \tau) - (1/n) \sum_{i \in [n]} L_{S_i}(f_k, \tau)| \geq \epsilon/2 \right] \\
 &\leq 2K \exp \left(-\frac{n\epsilon^2\tau}{8M \log B} \right) \\
 &= 2\mathcal{N}(\mathcal{F}, \tau\epsilon/16) \exp \left(-\frac{n\epsilon^2\tau}{8M \log B} \right), \tag{F.2}
 \end{aligned}$$

476 the first inequality is by union bound, the second is by triangle inequality, and the
 477 third is by Hoeffding's inequality and (F.1). Finally, plugging the condition $n \geq$
 478 $(8\tau^{-1}\epsilon^{-2}M \log B) \log(2\mathcal{N}(\mathcal{F}, \epsilon/8M)/\delta)$ into (F.2) we have that

$$\mathbb{P}_{S \sim \mathcal{D}^n} \left[\sup_{f \in \mathcal{F}} |L_{\mathcal{D}^B}(f, \tau) - \widehat{L}_S(f, \tau)| \geq \epsilon \right] \leq \delta,$$

479 which completes the proof. \square

480 G Transferrable Representation Learning

481 **Discussion on Assumption 3.1.** In simple terms, Assumption 3.1 is made on the data distribution to
 482 allow the *existence* of good encoding functions \mathbf{g}^* and \mathbf{h}^* . Specifically, the first bullet guarantees
 483 that the data with different \mathbf{z} , the underlying shared feature, is well distinguishable with margin γ . If
 484 the data from different \mathbf{z} does not satisfy this condition, the majority of the diagonal term $f(\mathbf{x}_i, \mathbf{y}_i)$ in
 485 (B.1) can be smaller than the off-diagonal term $f(\mathbf{x}_j, \mathbf{y}_i)$, which is not favored by the mechanism of
 486 CLIP. In other words, all encoding functions may yield higher similarity score for negative pairs than
 487 positive pairs. The second bullet requires the similarity score within each underlying shared feature
 488 not vary too much, which is naturally satisfied if the learned embeddings $\mathbf{g}(\mathbf{x})$, $\mathbf{h}(\mathbf{y})$ are consistent
 489 and do not vary too much given the same \mathbf{z} .

490 **Theorem G.1 (Formal).** Suppose Assumption 3.1 hold and we can find an ϵ approximate minimum
 491 $\hat{f} \in \mathcal{F}$ with respect to the temperature τ such that \hat{f} is bounded by M and

$$L_{\mathcal{D}^B}(\hat{f}, \tau) \leq L_{\mathcal{D}^B}(f^*, \tau) + \epsilon. \quad (\text{G.1})$$

492 Then the following results hold:

493 1. For $(\mathbf{x}, \mathbf{z}) \sim \mathcal{D}_{\mathbf{x} \times \mathbf{z}}$, $\{\mathbf{y}_k \sim \mathcal{D}_{\mathbf{y} | \mathbf{v}_k}, k \in [K]\}$, let $\mathbf{y}^* = \sum_{k \in [K]} \mathbf{1}(\mathbf{z} = \mathbf{v}_k) \mathbf{y}_k$, we have

$$\mathbb{E} \left[\log \left(\sum_{k \in [K]} \exp \left([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)] / \tau \right) \right) \right] \leq \epsilon'. \quad (\text{G.2})$$

494 2. For $(\mathbf{y}, \mathbf{z}) \sim \mathcal{D}_{\mathbf{y} \times \mathbf{z}}$, $\{\mathbf{x}_k \sim \mathcal{D}_{\mathbf{x} | \mathbf{v}_k}, k \in [K]\}$, let $\mathbf{x}^* = \sum_{k \in [K]} \mathbf{1}(\mathbf{z} = \mathbf{v}_k) \mathbf{x}_k$, we have

$$\mathbb{E} \left[\log \left(\sum_{k \in [K]} \exp \left([\hat{f}(\mathbf{x}_k, \mathbf{y}) - \hat{f}(\mathbf{x}^*, \mathbf{y})] / \tau \right) \right) \right] \leq \epsilon'. \quad (\text{G.3})$$

495 3. For $(\mathbf{x}, \mathbf{y}, \mathbf{z}) \sim \mathcal{D}_{\mathbf{x} \times \mathbf{y} \times \mathbf{z}}$, variance $\mathbb{E}_{(\mathbf{y}, \mathbf{z})} [\text{Var}_{\mathbf{x} | \mathbf{z}}(\hat{f}(\mathbf{x}, \mathbf{y}))] + \mathbb{E}_{(\mathbf{x}, \mathbf{z})} [\text{Var}_{\mathbf{y} | \mathbf{z}}(\hat{f}(\mathbf{x}, \mathbf{y}))] \leq$
 496 $16M^2\epsilon'$.

497 where $\epsilon' = (C_B + 2) \cdot [\epsilon + C\tau^{-1}MB\alpha + C\tau^{-1}(\beta MB)^{1/3} + 2B \exp(-\gamma/\tau)]$ and $C = \tilde{O}(1)$, $C_B =$
 498 $\tilde{O}(\max_k p_k^{-1}/B)$.

499 **Remark G.2.** Theorem G.1 establishes a soft margin between CLIP's learned embeddings on data of
 500 different \mathbf{z} 's. For instance, if an image \mathbf{x} has a shared feature $\mathbf{z} = \mathbf{v}_1$, we have its accurate description
 501 $\mathbf{y}^* = \sum_{k \in [K]} \mathbf{1}(\mathbf{z} = \mathbf{v}_k) \mathbf{y}_k = \mathbf{y}_1$. From (G.2), it follows that $\log \left(\sum_{k \in [K]} \exp \left([\hat{f}(\mathbf{x}, \mathbf{y}_k) -$
 502 $\hat{f}(\mathbf{x}, \mathbf{y}_1)] / \tau \right) \right)$ is small. This can only occur when $\hat{f}(\mathbf{x}, \mathbf{y}_k) < \hat{f}(\mathbf{x}, \mathbf{y}_1)$ for all $k \geq 2$, i.e., the
 503 trained model always yield higher similarity score for this image-text pair as compared to all other
 504 texts generated on different topics. This outcome aligns with the expectation that image-text pairs
 505 with the same shared feature will yield the highest similarity score.

506 **Remark G.3 (Choice of temperature parameter).** In Theorem G.1, when the data is well separated
 507 (i.e., $\alpha, \beta = 0$), a smaller temperature will invariably lead to a smaller ϵ' and, consequently, better
 508 performance. In practice, τ is typically set to be 0.01, a sufficiently small value that ensures the term
 509 $\exp(-\gamma/\tau)$ is less than 0.000454 for $\gamma = 0.1$. However, when the data is nonseparable (i.e., α and
 510 β exceed 0), a balance must be struck between the terms related to τ . As a consequence, τ should not
 511 be too small. A reasonable choice would be $\tau = O(\gamma / \log(B/\epsilon))$.

512 **Remark G.4 (Batch size).** In Theorem G.1, while we do not demand an increasing batch size B ,
 513 our analysis does suggest a preference for larger batch sizes, as they can reduce the constant C_B and
 514 consequently ϵ' .

515 **Lemma G.5.** For $b_j \geq 0, j \in [m]$, we have that

$$\log \left(1 + \sum_{j \in [m]} b_j \right) \leq \sum_{j \in [m]} \log(1 + b_j).$$

516 *Proof.* Notice that

$$\prod_{j \in [J]} (1 + b_j) \geq 1 + \sum_{j \in [J]} b_j.$$

517 Taking the logarithm over both sides completes the proof. \square

518 **Lemma G.6.** Suppose that a_1, \dots, a_m are i.i.d random variable sample lies in $[-R, R]$ where $R \geq 1$,
519 with mean $\mu := \mathbb{E}[a_1]$ and variance $\sigma^2 := \mathbb{E}[(a_1 - \mathbb{E}[a_1])^2]$. Then we have that

$$\mathbb{E}[\log \left(\sum_{i=1}^m \exp(a_i) \right)] \geq \log(m) + \mu + \frac{m-1}{4mR^2} \sigma^2.$$

520 *Proof.* Let $\bar{a} = [\sum_{i=1}^m a_i]/m$

$$\begin{aligned} \log \left(\sum_{i=1}^m \exp(a_i) \right) &= \log(m) + \frac{1}{m} \sum_{i=1}^m a_i + \log \left(\frac{1}{m} \sum_{i=1}^m \exp(a_i - \bar{a}) \right) \\ &\geq \log(m) + \frac{1}{m} \sum_{i=1}^m a_i + \log \left(1 + \frac{1}{3mR^2} \sum_{i=1}^m [a_i - \bar{a}]^2 \right) \\ &\geq \log(m) + \frac{1}{m} \sum_{i=1}^m a_i + \frac{1}{4mR^2} \sum_{i=1}^m [a_i - \bar{a}]^2. \end{aligned}$$

521 where the first inequality is by $\exp(t) \geq 1 + t + t^2/(3R^2), \forall t \in [-R, R]$, the second inequality is
522 due to $\log(1+t) \geq 3t/4, \forall t \in [0, 1/3]$.

523 □

524 **Lemma G.7.** Suppose f^* is the function that satisfies Assumption 3.1, then we have that

$$L_{\mathcal{D}^B}(f^*, \tau) \leq 2\mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbb{1}(\mathbf{z}_t = \mathbf{z}_1) \right) \right] + 6MB\alpha/\tau + 3\sqrt[3]{6MB\beta}/\tau + 2B \exp(-\gamma/\tau)$$

525 *Proof.* Let the event \mathcal{E}_t be the case that either i) $\mathbf{z}_t = \mathbf{z}_1$ and $|f^*(\mathbf{x}_t, \mathbf{y}_1) - f^*(\mathbf{x}_1, \mathbf{y}_1)| \leq \rho$ or ii)
526 $\mathbf{z}_t \neq \mathbf{z}_1$ and $f^*(\mathbf{x}_t, \mathbf{y}_1) - f^*(\mathbf{x}_1, \mathbf{y}_1) \leq -\gamma$. We also denote the complementary set of \mathcal{E}_t to be \mathcal{E}_t^c .
527 By Assumption 3.1, we have that

$$\begin{aligned} \mathbb{P}(\mathcal{E}_t, \mathbf{z}_t = \mathbf{z}_1) &\leq \beta/\rho^2 \\ \mathbb{P}(\mathcal{E}_t, \mathbf{z}_t \neq \mathbf{z}_1) &\leq \alpha. \end{aligned}$$

528 the first inequality is by Chebyshev's inequality, and the second is by margin assumption. Therefore,
529 we have that $\mathbb{P}(\mathcal{E}_t^c) \leq \alpha + \beta/\rho^2$. Next, let us decompose $L_{\mathcal{D}^B}(f^*, \tau)$ into three parts,

$$\begin{aligned} L_{\mathcal{D}^B}(f^*, \tau) &= \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbb{1}(\mathbf{z}_t \neq \mathbf{z}_1) \mathbb{1}(\mathcal{E}_t) \exp([f^*(\mathbf{x}_t, \mathbf{y}_t) - f^*(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right. \right. \\ &\quad \left. \left. + \sum_{t \in [B]} \mathbb{1}(\mathcal{E}_t^c) \exp([f^*(\mathbf{x}_1, \mathbf{y}_t) - f^*(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right. \right. \\ &\quad \left. \left. + \sum_{t \in [B]} \mathbb{1}(\mathbf{z}_t = \mathbf{z}_1) \mathbb{1}(\mathcal{E}_t) \exp([f^*(\mathbf{x}_1, \mathbf{y}_t) - f^*(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right] \\ &\quad + \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbb{1}(\mathbf{z}_t \neq \mathbf{z}_1) \mathbb{1}(\mathcal{E}_t) \exp([f^*(\mathbf{x}_t, \mathbf{y}_1) - f^*(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right. \right. \\ &\quad \left. \left. + \sum_{t \in [B]} \mathbb{1}(\mathcal{E}_t^c) \exp([f^*(\mathbf{x}_t, \mathbf{y}_1) - f^*(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right. \right. \\ &\quad \left. \left. + \sum_{t \in [B]} \mathbb{1}(\mathbf{z}_t = \mathbf{z}_1) \mathbb{1}(\mathcal{E}_t) \exp([f^*(\mathbf{x}_t, \mathbf{y}_1) - f^*(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right] \\ &\leq 2\mathbb{E} \left[\log \left(1 + B \exp(-\gamma/\tau) + \sum_{t \geq 2} \mathbb{1}(\mathcal{E}_t^c) \exp(2M/\tau) + \sum_{t \geq 2} \mathbb{1}(\mathbf{z}_t = \mathbf{z}_1) \exp(\rho/\tau) \right) \right] \\ &\leq 2\mathbb{E} \left[\underbrace{\log \left(1 + B \exp(-\gamma/\tau) \right)}_{I_1} \right] + \sum_{t \geq 2} 2\mathbb{E} \left[\underbrace{\log \left(1 + \mathbb{1}(\mathcal{E}_t^c) \exp(2M/\tau) \right)}_{I_2} \right] \end{aligned}$$

$$+ 2 \underbrace{\mathbb{E} \left[\log \left(1 + \sum_{t \geq 2} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp(\rho/\tau) \right) \right]}_{I_3} \quad (\text{G.4})$$

530 where the first inequality is by Assumption 3.1, the second inequality is due to Lemma G.5. Next, we
531 will bound I_1, I_2, I_3 separately.

$$I_1 \leq B \exp(-\gamma/\tau), \quad (\text{G.5})$$

532 where the inequality is due to the fact that $\log(1+x) \leq x$.

$$I_2 = \mathbb{E} \left[\mathbf{1}(\mathcal{E}_t^c) \log \left(1 + \exp(2M/\tau) \right) \right] \leq \mathbb{P}(\mathcal{E}_t^c) \frac{3M}{\tau} = (\alpha + \beta/\rho^2) \cdot \frac{3M}{\tau}. \quad (\text{G.6})$$

533 where the first equality is due to $\log(1 + \mathbf{1}(\mathcal{E}_t^c) \exp(2M/\tau)) = 0$ when $\mathbf{1}(\mathcal{E}_t^c) = 0$, the first
534 inequality is due to $\log(1 + \exp(2M/\tau)) \leq 3M/\tau$. The last inequality is due to $\mathbb{P}(\mathcal{E}_t^c) \leq \alpha + \beta/\rho^2$.

$$\begin{aligned} I_3 &\leq \mathbb{E} \left[\log \left(\exp(\rho/\tau) + \sum_{t \geq 2} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp(\rho/\tau) \right) \right] \\ &= \rho/\tau + \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \right) \right]. \end{aligned} \quad (\text{G.7})$$

535 where the inequality is because $1 \leq \exp(\rho/\tau)$.

536 Plugging (G.5), (G.6) and (G.7) into (G.4) gives that,

$$\begin{aligned} L_{\mathcal{D}^B}(f^*, \tau) &\leq 2B \exp(-\gamma/\tau) + 6MB\alpha/\tau + 6MB\beta/(\tau\rho^2) + 2\rho/\tau + 2\mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \right) \right] \\ &\leq 2\mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \right) \right] + 6MB\alpha/\tau + 3\sqrt[3]{6MB\beta}/\tau + 2B \exp(-\gamma/\tau), \end{aligned}$$

537 where the second inequality is by choosing $\rho = \sqrt[3]{6MB\beta}$. □

538 *Proof of Theorem G.1.* First by Lemma G.7, we have that

$$L_{\mathcal{D}^B}(\hat{f}, \tau) \leq L_{\mathcal{D}^B}(f^*, \tau) + \epsilon \leq 2\mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \right) \right] + \epsilon' \quad (\text{G.8})$$

539 where $\epsilon' = \epsilon + 6MB\alpha/\tau + 3\sqrt[3]{6MB\beta}/\tau + 2B \exp(-\gamma/\tau)$. Notice that

$$\begin{aligned} L_{\mathcal{D}^B}(\hat{f}, \tau) &= \mathbb{E} \left[\log \left(\sum_{t \in [B]} \exp \left([\hat{f}(\mathbf{x}_1, \mathbf{y}_t) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau \right) \right) \right] \\ &\quad \underbrace{\hspace{10em}}_{I_1} \\ &\quad + \mathbb{E} \left[\log \left(\sum_{t \in [B]} \exp \left([\hat{f}(\mathbf{x}_t, \mathbf{y}_1) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau \right) \right) \right] \quad (\text{G.9}) \\ &\quad \underbrace{\hspace{10em}}_{I_2} \end{aligned}$$

540 Next, we prove the bullets in Theorem G.1 one by one.

541 **First and Second Bullet in Theorem G.1:** Denote the event \mathcal{E} as the case that for all $t \geq 1$, $\mathbf{z}_t \neq \mathbf{z}_1$,
542 which is the event that CLIP favored. We first lower bound I_1 .

$$\begin{aligned} I_1 &= \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp \left([\hat{f}(\mathbf{x}_t, \mathbf{y}_1) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau \right) \right) \right] \\ &\quad + \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp \left([\hat{f}(\mathbf{x}_t, \mathbf{y}_1) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau \right) \right) \right] \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right. \right. \\
&\quad \left. \left. + \sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right] \\
&\geq \mathbb{E} \left[\mathbf{1}(\mathcal{E}) \log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \right] \\
&\quad + \mathbb{E} \left[\mathbf{1}(\mathcal{E}^c) \log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right] \\
&= \mathbb{E} \left[\mathbf{1}(\mathcal{E}) \log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \right] \\
&\quad + \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right] \\
&\geq \mathbb{E} \left[\mathbf{1}(\mathcal{E}) \log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \right] \\
&\quad + \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp(\mathbb{E}[\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1) | \mathbf{z}_t, \mathbf{z}_1]/\tau) \right) \right] \\
&= \mathbb{E} \left[\mathbf{1}(\mathcal{E}) \log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \right] \\
&\quad + \mathbb{E} \left[\log \left(|\{t \in [B] | \mathbf{z}_t = \mathbf{z}_1\}| \right) \right]. \tag{G.10}
\end{aligned}$$

543 where the first inequality is because when \mathcal{E} holds $\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) -$
544 $\widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) = 1$ when \mathcal{E}^c holds $\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \geq 0$, the last
545 second equality is because when \mathcal{E} holds $\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) = 1$,
546 the second inequality is because LogSumExp function is convex, and the last equality is due to
547 $\mathbb{E}[\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1) | \mathbf{z}_t, \mathbf{z}_1] = 0$ when $\mathbf{z}_t = \mathbf{z}_1$. Similarly, we can prove

$$\begin{aligned}
I_2 &\geq \mathbb{E} \left[\mathbf{1}(\mathcal{E}) \log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp([\widehat{f}(\mathbf{x}_1, \mathbf{y}_t) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \right] \\
&\quad + \mathbb{E} \left[\log \left(|\{t \in [B] | \mathbf{z}_t = \mathbf{z}_1\}| \right) \right]. \tag{G.11}
\end{aligned}$$

548 Notice that when event \mathcal{E} holds, $\mathbf{z}_t \neq \mathbf{z}_1$ holds for all $t \geq 2$. Therefore, plugging the (G.10) and
549 (G.11) into (G.9) gives,

$$\mathbb{E} \left[\mathbf{1}(\mathcal{E}) \log \left(\sum_{t \geq 2} \exp([\widehat{f}(\mathbf{x}_t, \mathbf{y}_1) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \right] \leq \epsilon' \tag{G.12}$$

$$\mathbb{E} \left[\mathbf{1}(\mathcal{E}) \log \left(\sum_{t \geq 2} \exp([\widehat{f}(\mathbf{x}_1, \mathbf{y}_t) - \widehat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \right] \leq \epsilon'. \tag{G.13}$$

$$\tag{G.14}$$

550 Let us compute the probability of \mathcal{E} given \mathbf{z}_1 . Let $\mathbf{z}_1 = \mathbf{v}_1$ without loss of generality, we have that

$$\mathbb{P}(\mathcal{E} | \mathbf{z} = \mathbf{v}_1) = (1 - p_1)^{B-1}.$$

551 Therefore $\mathbb{P}(\mathcal{E} | \mathbf{z} = \mathbf{v}_1)$ is always positive and is greater than $1/2$ as long as $B \leq 1/p_1$.

552 Next, consider the following situation. Given $\mathbf{z}_1 = \mathbf{v}_1$, we generate sequence $\mathbf{z}'_1, \dots, \mathbf{z}'_L$ with length
553 $L = \lceil \log(2K)/(B-1) \min p_k \rceil (B-1)$, such that each $\mathbf{z}'_1, \dots, \mathbf{z}'_L$ are generated from $\mathcal{D}_{\mathbf{z} | \mathbf{z} \neq \mathbf{v}_1}$.

554 The probability that the sequence includes \mathbf{v}_k is

$$1 - (1 - p_k/(1 - p_k))^L \geq 1 - (1 - p_k)^L \geq 1 - \exp(-Lp_k) \geq 1 - \exp(-L \min p_k).$$

555 Therefore the probability that the sequence can cover all the other $K - 1$ classes is at least

$$1 - K \exp(-L \min p_k) \geq 1/2.$$

556 Then we look deeper into

$$\mathbb{E} \left[\log \left(\sum_{t \geq 2} \exp([\hat{f}(\mathbf{x}_t, \mathbf{y}_1) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \middle| \mathbf{z}_1 = \mathbf{v}_1, \mathbf{z}_2 \neq \mathbf{v}_1, \dots, \mathbf{z}_K \neq \mathbf{v}_1 \right].$$

557 We can introduce $L/(B - 1)$ copies $\mathbf{x}_t^{(l)}$ with $l \in [L/(B - 1)]$ for $t \geq 2$, then we have that

$$\begin{aligned} & \left(L/(B - 1) \right) \cdot \mathbb{E} \left[\log \left(\sum_{t \geq 2} \exp([\hat{f}(\mathbf{x}_t, \mathbf{y}_1) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \middle| \mathbf{z}_1 = \mathbf{v}_1, \mathbf{z}_2 \neq \mathbf{v}_1, \dots, \mathbf{z}_K \neq \mathbf{v}_1 \right] \\ &= \mathbb{E} \left[\sum_l \log \left(\sum_{t \geq 2} \exp([\hat{f}(\mathbf{x}_t^{(l)}, \mathbf{y}_1) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \middle| \mathbf{z}_1 = \mathbf{v}_1, \mathbf{z}_2^{(l)}, \dots, \mathbf{z}_K^{(l)} \neq \mathbf{v}_1 \right] \\ &\geq \mathbb{E} \left[\log \left(\sum_l \sum_{t \geq 2} \exp([\hat{f}(\mathbf{x}_t^{(l)}, \mathbf{y}_1) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) + 1 \right) \middle| \mathbf{z}_1 = \mathbf{v}_1, \mathbf{z}_2^{(l)}, \dots, \mathbf{z}_K^{(l)} \neq \mathbf{v}_1 \right] \\ &\geq \mathbb{E} \left[\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}_k, \mathbf{y}) - \hat{f}(\mathbf{x}^*, \mathbf{y})]/\tau) \right) \middle| \mathbf{z} = \mathbf{v}_1 \right]. \end{aligned} \quad (\text{G.15})$$

558 where the first inequality is by Lemma G.5, the second inequality is by the fact that the Exp function
559 is greater than 0, and the $\mathbf{x}_k, \mathbf{x}^*$ in the last line are the ones that defined in Theorem G.1. Plugging
560 (G.15) into (G.12) and applying total expectation completes the proof for the second bullet. The
561 proof for the first bullet is the same.

562 **Third Bullet in Theorem G.1:** By the third equality in (G.10), we have that

$$\begin{aligned} I_1 &\geq \mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp([\hat{f}(\mathbf{x}_t, \mathbf{y}_1) - \hat{f}(\mathbf{x}_1, \mathbf{y}_1)]/\tau) \right) \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \exp(\hat{f}(\mathbf{x}_t, \mathbf{y}_1)/\tau) \right) \middle| \mathbf{z}_1, \dots, \mathbf{z}_B \right] - \mathbb{E}[\hat{f}(\mathbf{x}_1, \mathbf{y}_1)/\tau] \right] \\ &\geq \mathbb{E} \left[\log \left(\left| \{t \in [B] \mid \mathbf{z}_t = \mathbf{z}_1\} \right| \right) \right] + \mathbb{E} \left[\frac{\left| \{t \in [B] \mid \mathbf{z}_t = \mathbf{z}_1\} \right| - 1}{4M^2 \left| \{t \in [B] \mid \mathbf{z}_t = \mathbf{z}_1\} \right|} \text{Var}_{\mathbf{x}_1 | \mathbf{z}_1}(\hat{f}(\mathbf{x}_1, \mathbf{y}_1)) \right]. \end{aligned} \quad (\text{G.16})$$

563 where the inequality is by Lemma G.6. Next we will We analyze the distribution of $\{t \in [B] \mid \mathbf{z}_t =$
564 $\mathbf{z}_1\}$. Without loss of generality, fix $\mathbf{z}_1 = \mathbf{v}_1$. We know that the probability that $\{t \in [B] \mid \mathbf{z}_t =$
565 $\mathbf{z}_1\} \geq 2$ is

$$1 - \mathbb{P}(\mathbf{z}_2 \neq \mathbf{z}_1) \cdot \dots \cdot \mathbb{P}(\mathbf{z}_B \neq \mathbf{z}_1) \geq 1 - (1 - \min p_k)^{B-1} \geq \min\{0.25 * \min p_k \cdot (B - 1), 0.25\},$$

566 the last inequality holds since the strictly increasing function $F(s) = 1 - (1 - \min p_k)^s$ is 0 at $s = 0$
567 and have derivative lower bounded by 0.25 when $s \leq 1/\min p_k$. Therefore we can further lower
568 bound (G.16) as follows,

$$I_1 \geq \mathbb{E} \left[\log \left(\left| \{t \in [B] \mid \mathbf{z}_t = \mathbf{z}_1\} \right| \right) \right] + \mathbb{E} \left[\frac{\min\{0.25 * \min p_k \cdot (B - 1), 0.25\}}{8M^2} \text{Var}_{\mathbf{x}_1 | \mathbf{z}_1}(\hat{f}(\mathbf{x}_1, \mathbf{y}_1)) \right]$$

569 Similarly, we can prove that

$$I_2 \geq \mathbb{E} \left[\log \left(\left| \{t \in [B] \mid \mathbf{z}_t = \mathbf{z}_1\} \right| \right) \right] + \mathbb{E} \left[\frac{\min\{0.25 * \min p_k \cdot (B - 1), 0.25\}}{8M^2} \text{Var}_{\mathbf{y}_1 | \mathbf{z}_1}(\hat{f}(\mathbf{x}_1, \mathbf{y}_1)) \right].$$

570 Plugging the bound of I_1, I_2 into (G.9) completes the proof for the third bullet of Theorem G.1. \square

571 **H Proof of the Results in Section 3**

572 **Corollary H.1.** Suppose the result of Theorem G.1 holds for the learned similarity function \hat{f} . Then
 573 we calculate the similarity score $\hat{f}(\mathbf{x}, \mathbf{y}_k)$ for all $k \in [K]$ and pick the indices of the top- r scores
 574 within the set $\{\hat{f}(\mathbf{x}, \mathbf{y}_k)\}$ as the predictions of the image \mathbf{x} . Then the top- r error is bounded by
 575 $\epsilon' / \log(1 + r)$.

576 *Proof of Corollary H.1.* For $(\mathbf{x}, \mathbf{z}) \sim \mathcal{D}_{\mathbf{x} \times \mathbf{z}}$, $\{\mathbf{y}_k \sim \mathcal{D}_{\mathbf{y}|\mathbf{v}_k}, k \in [K]\}$, let $\mathbf{y}^* = \sum_{k \in [K]} \mathbb{1}(\mathbf{z} =$
 577 $\mathbf{v}_k) \mathbf{y}_k$. Denote \mathcal{E} to be the event that the top- r choice gives the wrong prediction. Then we have that,

$$\begin{aligned} \epsilon' &\geq \mathbb{E} \left[\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) \right] \\ &\geq \mathbb{E} \left[\mathbb{1}(\mathcal{E}) \log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) \right] \\ &\geq \mathbb{E} \left[\mathbb{1}(\mathcal{E}) \log(1 + r) \right] \\ &= \mathbb{P}(\mathcal{E}) \log(1 + r), \end{aligned}$$

578 where the first inequality is by the first bullet of Theorem G.1, the second inequality is due to
 579 the fact that $\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) > 0$, the last inequality is due to
 580 $\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) \geq \log(1 + r)$ since there are at least $r + 1$ number
 581 of $\hat{f}(\mathbf{x}, \mathbf{y}_k)$ are greater than $\hat{f}(\mathbf{x}, \mathbf{y}^*)$ if the prediction is wrong. Therefore, we have that $\mathbb{P}(\mathcal{E}) \leq$
 582 $\epsilon' / \log(1 + r)$ which completes the proof. \square

583 **Remark H.2.** The result in Corollary H.1 can be generalized to out-of-distribution zero-shot transfer.
 584 For example, we can deal with the case where the distribution of the prompts $\mathcal{D}_{\mathbf{y}|\mathbf{v}_k}$ and the image
 585 distribution $\mathcal{D}_{\mathbf{x}}$ are shifted. As long as the χ^2 distance between the shifted distributions is bounded,
 586 we can provide a top- r error guarantee.

587 **Discussion for out-of-distribution zero shot learning.** The result in Corollary H.1 can be generalized
 588 to out-of-distribution zero-shot transfer learning. For example, we can deal with the case where the
 589 distribution of the prompts $\mathcal{D}_{\mathbf{y}|\mathbf{v}_k}$ and the image distribution $\mathcal{D}_{\mathbf{x}}$ are shifted. In particular, let us
 590 consider the case that the distribution of the prompts is shifted to $\mathcal{D}'_{\mathbf{y}|\mathbf{v}_k}$ and the image distribution $\mathcal{D}_{\mathbf{x}}$
 591 is shifted to $\mathcal{D}'_{\mathbf{x}}$. Then the original joint cumulative distribution function function $P(\mathbf{x}, \mathbf{z}, \mathbf{y}_1, \dots, \mathbf{y}_K)$
 592 is shifted to $Q(\mathbf{x}, \mathbf{z}, \mathbf{y}_1, \dots, \mathbf{y}_K)$. Suppose Q is absolutely continuous with respect to P , and the
 593 Pearson χ^2 distance is bounded

$$\int \left(\frac{dQ}{dP} - 1 \right)^2 dP \leq C.$$

594 Then we have that

$$\begin{aligned} &\int \sqrt{\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right)} dQ \\ &= \int \sqrt{\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right)} \left(\frac{dQ}{dP} \right) dP \\ &\leq \sqrt{\int \log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) dP} \cdot \sqrt{\int \left(\frac{dQ}{dP} \right)^2 dP} \\ &= \sqrt{(C + 1)\epsilon'}, \end{aligned}$$

595 where the first inequality is by Cauchy Schwartz inequality and the last equality is due to
 596 $\int \left(\frac{dQ}{dP}\right)^2 dP = \int \left(\frac{dQ}{dP} - 1\right)^2 dP + 1 = C + 1$. Then we can follow a similar analysis in the
 597 proof of Corollary H.1 and have that top-r test error is smaller than $\sqrt{(C+1)\epsilon'/\log(1+r)}$. There-
 598 fore, if the χ^2 distance between the shifted distributions is bounded, we can still provide a top-r error
 599 guarantee. It is worth noting the bound for out-of-distribution zero-shot learning is looser. If we want
 600 to do a more general zero shot analysis, we may need to add more data structure in Assumption 3.1.

601 **Lemma H.3** (Completeness). There exist a score function $f^*(\mathbf{x}, \mathbf{y}) = \langle \mathbf{W}^* \mathbf{x}, \mathbf{y} \rangle$ with $\mathbf{W}^* \in \mathbb{R}^{d_2 \times d_1}$
 602 satisfying

- 603 • $|f^*| \leq 1$,
- 604 • For $(\mathbf{x}, \mathbf{y}, \mathbf{z}) \sim \mathcal{D}_{\mathbf{x} \times \mathbf{y} \times \mathbf{z}}$, variance $\mathbb{E}_{(\mathbf{y}, \mathbf{z})} [\text{Var}_{\mathbf{x}|\mathbf{z}}(f^*(\mathbf{x}, \mathbf{y}))] = \mathbb{E}_{(\mathbf{x}, \mathbf{z})} [\text{Var}_{\mathbf{y}|\mathbf{z}}(f^*(\mathbf{x}, \mathbf{y}))] = 0$,
- 605 • Let $\mathbf{x} \sim \mathcal{D}_{\mathbf{x}|\mathbf{z}}$, $\mathbf{y} \sim \mathcal{D}_{\mathbf{y}|\mathbf{z}}$, $\mathbf{x}' \sim \mathcal{D}_{\mathbf{x}'|\mathbf{z}'}$, $\mathbf{y}' \sim \mathcal{D}_{\mathbf{y}'|\mathbf{z}'}$ where $\mathbf{z} \neq \mathbf{z}'$. With probability 1, we have that
 606 $f^*(\mathbf{x}', \mathbf{y}) \leq f^*(\mathbf{x}, \mathbf{y}) - \gamma$ and $f^*(\mathbf{x}, \mathbf{y}') \leq f^*(\mathbf{x}, \mathbf{y}) - \gamma$.

607 *Proof of Lemma H.3.* We can construct $\mathbf{W}^* = \mathbf{H}(\mathbf{H}^\top \mathbf{H})^{-1} \mathbf{P}(\mathbf{G}^\top \mathbf{G})^{-1} \mathbf{G}^\top$, where $\mathbf{P} \in$
 608 $\mathbb{R}^{(K_1+K_2) \times (K_1+K_3)}$ is the projection matrix $\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ with rank K_1 .

609 It is easy to verify that $\mathbf{H}^\top \mathbf{W}^* \mathbf{G} = \mathbf{P}$. Therefore we have that

$$\langle \mathbf{W}^* \mathbf{x}, \mathbf{y}' \rangle = \langle \mathbf{z}, \mathbf{z}' \rangle.$$

610 Then applying $\|\mathbf{v}_k\|_2 = 1$, $\langle \mathbf{v}_k, \mathbf{v}'_{k'} \rangle \leq 1 - \gamma$, $\forall k \neq k'$ completes the proof. \square

611 **Lemma H.4.** $\|\nabla L_S(f_{\mathbf{W}}, \tau)\|_F \leq L$ where $L = 2\tau^{-1} \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1)$.

612 *Proof.* First, we have that

$$\|\nabla_{\mathbf{W}} \langle \mathbf{W} \mathbf{x}, \mathbf{y} \rangle\|_F = \|\mathbf{x} \mathbf{y}^\top\|_F \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \leq \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1).$$

613 Therefore we have that $\|\nabla L_S(f_{\mathbf{W}}, \tau)\|_F \leq 2\tau^{-1} \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1)$ since LogSumExp function
 614 is an 1-Lipschitz function. \square

615 *Proof of Theorem 3.4.* By the gradient update rule, we have that

$$\begin{aligned} & \|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_F^2 \\ &= 2\eta \langle \nabla \widehat{L}_S(\mathbf{W}^{(t)}, \tau), \mathbf{W}^{(t)} - \mathbf{W}^* \rangle - \eta^2 \|\nabla \widehat{L}_S(\mathbf{W}^{(t)}, \tau)\|_F^2 \\ &\geq 2\eta \widehat{L}_S(\mathbf{W}^{(t)}, \tau) - 2\eta \widehat{L}_S(\mathbf{W}^*, \tau) - \eta^2 L^2. \end{aligned} \quad (\text{H.1})$$

616 Take the telescope sum of (H.1) from 0 to $T - 1$ we have that

$$\begin{aligned} \frac{\sum_{t=0}^{T-1} \widehat{L}_S(\mathbf{W}^{(t)}, \tau)}{T} &\leq \widehat{L}_S(\mathbf{W}^*, \tau) + \eta L^2 + \frac{\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(T)} - \mathbf{W}^*\|_F^2}{2\eta T} \\ &\leq \widehat{L}_S(\mathbf{W}^*, \tau) + \epsilon/4 + \epsilon/4 \\ &= \widehat{L}_S(\mathbf{W}^*, \tau) + \epsilon/2, \end{aligned}$$

617 where the second inequality is by $\eta \leq \epsilon/(4L^2)$ and $T = 4\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2/(\eta\epsilon)$. Therefore, there
 618 exist $t' \leq T - 1$ such that $\widehat{L}_S(\mathbf{W}^{(t')}, \tau) \leq \widehat{L}_S(\mathbf{W}^*, \tau) + \epsilon/2$. Let \widehat{T} to be the first time that
 619 $\widehat{L}_S(\mathbf{W}^{(\widehat{T})}, \tau) \leq \widehat{L}_S(\mathbf{W}^*, \tau) + \epsilon/2$. Again take telescope sum of (H.1) from 0 to $\widehat{T} - 1$, we have that

$$\begin{aligned} \|\mathbf{W}^{(\widehat{T})} - \mathbf{W}^*\|_F^2 &\leq 2\eta \widehat{T} \widehat{L}_S(\mathbf{W}^*, \tau) - 2\eta \widehat{T} \sum_{t=0}^{\widehat{T}-1} \widehat{L}_S(\mathbf{W}^{(t)}, \tau) + 2\eta^2 L^2 \widehat{T} + \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 \\ &\leq -\eta \widehat{T} \epsilon + 0.5\eta \widehat{T} \epsilon + \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2 \\ &\leq \|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2, \end{aligned}$$

620 where the second inequality is due to the definition of \widehat{T} , the last inequality is due to $-0.5\eta \widehat{T} \epsilon \leq 0$.

621 Therefore, within $T = 4\|\mathbf{W}^{(0)} - \mathbf{W}^*\|_F^2/(\eta\epsilon)$ we can find $\widehat{\mathbf{W}} = \mathbf{W}^{(\widehat{T})}$ such that $\widehat{L}_S(\widehat{\mathbf{W}}, \tau) \leq$
 622 $\widehat{L}_S(\mathbf{W}^*, \tau) + \epsilon/2$ and

$$\|\mathbf{W}^{(\widehat{T})}\|_F^2 \leq 2\|\mathbf{W}^*\|_F + \|\mathbf{W}^{(0)}\|_F^2$$

623 where the inequality is by triangle inequality. Therefore, for any \mathbf{x}, \mathbf{y}

$$\begin{aligned}\widehat{f}(\mathbf{x}, \mathbf{y}) &= \langle \mathbf{W}^* \mathbf{x}, \mathbf{y} \rangle + \langle \widehat{\mathbf{W}} - \mathbf{W}^* \mathbf{x}, \mathbf{y} \rangle \\ &\leq 1 + \|\widehat{\mathbf{W}} - \mathbf{W}^*\|_F \|\mathbf{x} \mathbf{y}^\top\|_F \\ &\leq 1 + \|\widehat{\mathbf{W}} - \mathbf{W}^*\|_F \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1) \\ &\leq 1 + \|\mathbf{W}^* - \mathbf{W}^{(0)}\|_F \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1).\end{aligned}$$

624 Therefore the function \widehat{f} is bonded by $M = 1 + \|\mathbf{W}^* - \mathbf{W}^{(0)}\|_F \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1)$. Moreover,
625 the function \widehat{f} must belong to the class $\mathcal{F} = \{\langle \mathbf{W} \mathbf{x}, \mathbf{y} \rangle \mid \|\mathbf{W}\|_F \leq 2\|\mathbf{W}^*\|_F + \|\mathbf{W}^{(0)}\|_F^2\}$. Since the
626 linear function class \mathcal{F} has finite covering the set $\mathcal{N}(\mathcal{F}, \epsilon)$ (Bartlett and Mendelson, 2002; Zhang,
627 2002), by Theorem B.2 we know that when $n \geq (8\tau^{-1}\epsilon^{-2}M \log B) \log(2\mathcal{N}(\mathcal{F}, \epsilon/32M)/\delta)$, with
628 probability at least $1 - \delta$ we have that

$$\begin{aligned}|\widehat{L}_S(\widehat{f}, \tau) - L_{\mathcal{D}^B}(\widehat{f}, \tau)| &\leq \epsilon/4 \\ |\widehat{L}_S(f^*, \tau) - L_{\mathcal{D}^B}(f^*, \tau)| &\leq \epsilon/4.\end{aligned}$$

629 Thus, we can conclude that

$$\begin{aligned}\widehat{L}_{\mathcal{D}^B}(\widehat{f}, \tau) - \widehat{L}_{\mathcal{D}^B}(f^*, \tau) &\leq \widehat{L}_S(\widehat{f}, \tau) - \widehat{L}_S(f^*, \tau) + |\widehat{L}_S(\widehat{f}, \tau) - L_{\mathcal{D}^B}(\widehat{f}, \tau)| \\ &\quad + |\widehat{L}_S(f^*, \tau) - L_{\mathcal{D}^B}(f^*, \tau)| \\ &\leq \epsilon/2 + \epsilon/4 + \epsilon/4 \\ &= \epsilon.\end{aligned}$$

630 where the first inequality is by the triangle inequality, the second inequality is by the bounded gap
631 between empirical and population loss. \square

Proof of Theorem 3.5.

$$\begin{aligned}\mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbf{y}\|_2^2 | \mathbf{z}] &= \mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbb{E}[\mathbf{y} | \mathbf{z}] + \mathbb{E}[\mathbf{y} | \mathbf{z}] - \mathbf{y}\|_2^2 | \mathbf{z}] \\ &= \mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbb{E}[\mathbf{y} | \mathbf{z}]\|_2^2 | \mathbf{z}] + \mathbb{E}[\|\mathbb{E}[\mathbf{y} | \mathbf{z}] - \mathbf{y}\|_2^2 | \mathbf{z}]\end{aligned}$$

632 where the second equality is due to $\mathbf{x} \perp \mathbf{y} | \mathbf{z}$ and $\mathbb{E}[\mathbb{E}[\mathbf{y} | \mathbf{z}] - \mathbf{y} | \mathbf{z}] = \mathbf{0}$. Then taking a total expectation
633 over both sides over \mathbf{z} gives that

$$\mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbf{y}\|_2^2] = \mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbb{E}[\mathbf{y} | \mathbf{z}]\|_2^2] + \mathbb{E}[\|\mathbf{y} - \mathbb{E}[\mathbf{y} | \mathbf{z}]\|_2^2] \geq \mathbb{E}[\|\mathbf{y} - \mathbb{E}[\mathbf{y} | \mathbf{z}]\|_2^2].$$

634 Obviously, $\mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbf{y}\|_2^2]$ achieves global minima when

$$\mathbf{g}(\mathbf{x}) = \mathbb{E}[\mathbf{y} | \mathbf{z}] = \mathbf{H} \begin{bmatrix} \mathbf{z} \\ \mathbb{E}[\boldsymbol{\zeta} | \mathbf{z}] \end{bmatrix}.$$

635 This function \mathbf{g} is also achievable. We can construct function $\mathbf{g}_2(\mathbf{z}) = \mathbf{H} \begin{bmatrix} \mathbf{z} \\ \mathbb{E}[\boldsymbol{\zeta} | \mathbf{z}] \end{bmatrix}$, and projection
636 function $\mathbf{g}_1(\mathbf{x}) = \mathbf{z}$ that is linear. Then we can define $\mathbf{g} = \mathbf{g}_2 \circ \mathbf{g}_1$. \square

637 *Proof of Corollary 3.6.* Since $\boldsymbol{\zeta}$ is independent with \mathbf{z} , we have that

$$\mathbf{g}(\mathbf{x}) = \mathbf{H} \begin{bmatrix} \mathbf{z} \\ \mathbb{E}[\boldsymbol{\zeta} | \mathbf{z}] \end{bmatrix} = 1/3 \cdot \begin{bmatrix} \mathbf{z} \\ \mathbf{e}_1 \\ \mathbf{0} \end{bmatrix} + 2/3 \cdot \begin{bmatrix} \mathbf{z} \\ \mathbf{e}_2 \\ \mathbf{0} \end{bmatrix}.$$

638 Besides, we have that

$$\mathbf{y}' = \mathbf{H} \begin{bmatrix} \mathbf{z}' \\ \boldsymbol{\zeta}' \\ \mathbf{0} \end{bmatrix}$$

639 **Inner product similarity.** We have that $f(\mathbf{x}, \mathbf{y}') = \langle \mathbf{z}, \mathbf{z}' \rangle + 1/3 + 1/3 \cdot \mathbb{1}(\zeta' = \mathbf{e}_2)$. Since margin
 640 $\gamma < 1/3$. There exist j, k such that $\langle \mathbf{v}_j, \mathbf{v}_k \rangle > 2/3$. Then for $\mathbf{z} = \mathbf{v}_j$, we will sample K prompt

641 $\mathbf{y}_1, \dots, \mathbf{y}_K$. When $\mathbf{y}_j = \begin{bmatrix} \mathbf{v}_j \\ \mathbf{e}_1 \\ \mathbf{0} \end{bmatrix}$ and $\mathbf{y}_k = \begin{bmatrix} \mathbf{v}_k \\ \mathbf{e}_2 \\ \mathbf{0} \end{bmatrix}$, we have that

$$f(\mathbf{x}, \mathbf{y}_j) = 4/3 < \langle \mathbf{v}_j, \mathbf{v}_k \rangle + 2/3 = f(\mathbf{x}, \mathbf{y}_k),$$

642 which leads to the wrong top-1 prediction. The key insight behind this consequence is that $f(\mathbf{x}, \mathbf{y}') =$
 643 $\langle \mathbf{z}, \mathbf{z}' \rangle + 1/3 + 1/3 \cdot \mathbb{1}(\zeta' = \mathbf{e}_2)$ is greatly influenced by the unique feature ζ . A similar case also

644 exists for $\mathbf{z} = \mathbf{v}_k$ with $\mathbf{y}_j = \begin{bmatrix} \mathbf{v}_j \\ \mathbf{e}_2 \\ \mathbf{0} \end{bmatrix}$ and $\mathbf{y}_k = \begin{bmatrix} \mathbf{v}_k \\ \mathbf{e}_1 \\ \mathbf{0} \end{bmatrix}$. The probability that the above event occurs is at

645 least $2/K \cdot 1/3 \cdot 2/3 = 4/(9K) \geq 1/(3K)$. Therefore, the test error is at least $1/(3K)$.

646 **Cosine similarity.** Notice that $\|\mathbf{g}(\mathbf{x})\|_2 = \sqrt{1 + 1/9 + 4/9} = \sqrt{14}/3$, and $\|\mathbf{y}\|_2 = 1$, therefore the
 647 cosine similarity is proportional to inner product similarity with factor $\sqrt{14}/3$. Thus, the test error is
 648 still at least $1/(3K)$.

649 **L_2 similarity.** We have that $f(\mathbf{x}, \mathbf{y}') = -\|\mathbf{z} - \mathbf{z}'\|_2^2 - 8/9 + 2/3 \cdot \mathbb{1}(\zeta' = \mathbf{e}_2)$. Since margin
 650 $\gamma < 1/3$. There exist j, k such that $\|\mathbf{v}_j - \mathbf{v}_k\|_2^2 < 2/3$. Then for $\mathbf{z} = \mathbf{v}_j$, we will sample K prompt

651 $\mathbf{y}_1, \dots, \mathbf{y}_K$. When $\mathbf{y}_j = \begin{bmatrix} \mathbf{v}_j \\ \mathbf{e}_1 \\ \mathbf{0} \end{bmatrix}$ and $\mathbf{y}_k = \begin{bmatrix} \mathbf{v}_k \\ \mathbf{e}_2 \\ \mathbf{0} \end{bmatrix}$, we have that

$$f(\mathbf{x}, \mathbf{y}_j) = -8/9 < -\|\mathbf{v}_j, \mathbf{v}_k\|_2^2 + 2/3 = f(\mathbf{x}, \mathbf{y}_k),$$

652 which leads to the wrong top-1 prediction. The key insight behind this consequence is that $f(\mathbf{x}, \mathbf{y}') =$
 653 $-\|\mathbf{z} - \mathbf{z}'\|_2^2 - 8/9 + 2/3 \cdot \mathbb{1}(\zeta' = \mathbf{e}_2)$ is greatly influenced by the unique feature ζ . A similar case

654 also exists for $\mathbf{z} = \mathbf{v}_k$ with $\mathbf{y}_j = \begin{bmatrix} \mathbf{v}_j \\ \mathbf{e}_2 \\ \mathbf{0} \end{bmatrix}$ and $\mathbf{y}_k = \begin{bmatrix} \mathbf{v}_k \\ \mathbf{e}_1 \\ \mathbf{0} \end{bmatrix}$. The probability that the above event occurs

655 is at least $2/K \cdot 1/3 \cdot 2/3 = 4/(9K) \geq 1/(3K)$. Therefore, the test error is at least $1/(3K)$.

656 □

657 I Proof of Results in Section 4

658 *Proof of Corollary 4.1.* For $(\mathbf{x}, \mathbf{z}) \sim \mathcal{D}_{\mathbf{x} \times \mathbf{z}}$, $\{\mathbf{y}_k \sim \mathcal{D}_{\mathbf{y}|\mathbf{v}_k}, k \in [K]\}$, let $\mathbf{y}^* = \sum_{k \in [K]} \mathbb{1}(\mathbf{z} =$
 659 $\mathbf{v}_k) \mathbf{y}_k$. Denote \mathcal{E} to be the event that the top-1 choice gives the wrong prediction or the margin is
 660 smaller than τ . Then we have that,

$$\begin{aligned} \epsilon' &\geq \mathbb{E} \left[\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) \right] \\ &\geq \mathbb{E} \left[\mathbb{1}(\mathcal{E}) \log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) \right] \\ &\geq \mathbb{E} \left[\mathbb{1}(\mathcal{E}) \log(1 + \exp(-1)) \right] \\ &= \mathbb{P}(\mathcal{E}) \log(1 + e^{-1}), \end{aligned}$$

661 where the first inequality is by the first bullet of Theorem G.1, the second inequality is due to
 662 the fact that $\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) > 0$, the last inequality is due to

663 $\log \left(\sum_{k \in [K]} \exp([\hat{f}(\mathbf{x}, \mathbf{y}_k) - \hat{f}(\mathbf{x}, \mathbf{y}^*)]/\tau) \right) \geq \log(1 + e^{-1})$ since there exists at least one

664 similarity score $\hat{f}(\mathbf{x}, \mathbf{y}_k)$ greater than $\hat{f}(\mathbf{x}, \mathbf{y}^*) - \tau$ with $\mathbf{y}_k \neq \mathbf{y}^*$. Therefore, we have that $\mathbb{P}(\mathcal{E}) \leq$
 665 $\epsilon' / \log(1 + e^{-1}) \leq 4\epsilon'$ which completes the proof. □

666 **Discussion of Theorem 4.2.** The reason is that softmax function $L(\mathbf{a}) = \log(\sum_i \exp(a_i))$ is convex
 667 but not strongly convex and has an exponential-decaying tail. Once the score function f with
 668 the features \mathbf{g} and \mathbf{h} achieves the margin of order $\Omega(\tau)$, the gradient will exponentially decrease.
 669 Therefore, the weight will not be updated effectively.

670 *Proof of Theorem 4.2.* Consider the simplest setting where $\boldsymbol{\xi}$ and $\boldsymbol{\zeta}$ are all zero vectors, and we can
671 access to the population loss and its gradient (notice that we are constructing the negative example).
672 We will show that even under this ideal setting, the learned score function with corresponding
673 representations may not achieve a margin greater than $\bar{O}(\tau)$. Notice that

$$\begin{aligned}
\nabla_{\mathbf{W}} \mathbb{E}_{\mathcal{D}^B} L(f, \tau) &= \nabla_{\mathbf{W}} \mathbb{E} \left[\log \left(\sum_{t \in [B]} \exp \left([f(\mathbf{x}_1, \mathbf{y}_t) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right) \right) \right] \\
&\quad + \nabla_{\mathbf{W}} \mathbb{E} \left[\log \left(\sum_{t \in [B]} \exp \left([f(\mathbf{x}_t, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right) \right) \right] \\
&= \mathbb{E} \left[\nabla_{\mathbf{W}} \log \left(\sum_{t \in [B]} \exp \left([f(\mathbf{x}_1, \mathbf{y}_t) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right) \right) \right] \\
&\quad + \mathbb{E} \left[\nabla_{\mathbf{W}} \log \left(\sum_{t \in [B]} \exp \left([f(\mathbf{x}_t, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right) \right) \right] \\
&= \mathbb{E} \left[\sum_{t \in [B]} \frac{\exp \left([f(\mathbf{x}_1, \mathbf{y}_t) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)}{\sum_s \exp \left([f(\mathbf{x}_1, \mathbf{y}_s) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)} (\mathbf{y}_t - \mathbf{y}_1) \mathbf{x}_1^\top \right] \\
&\quad + \mathbb{E} \left[\sum_{t \in [B]} \frac{\exp \left([f(\mathbf{x}_t, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)}{\sum_s \exp \left([f(\mathbf{x}_s, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)} \mathbf{y}_1 (\mathbf{x}_t - \mathbf{x}_1)^\top \right] \\
&= \mathbb{E} \left[\sum_{t \in [B]} \frac{\mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp \left([f(\mathbf{x}_1, \mathbf{y}_t) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)}{\sum_s \exp \left([f(\mathbf{x}_1, \mathbf{y}_s) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)} (\mathbf{y}_t - \mathbf{y}_1) \mathbf{x}_1^\top \right] \\
&\quad + \mathbb{E} \left[\sum_{t \in [B]} \frac{\mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp \left([f(\mathbf{x}_t, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)}{\sum_s \exp \left([f(\mathbf{x}_s, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)} \mathbf{y}_1 (\mathbf{x}_t - \mathbf{x}_1)^\top \right]
\end{aligned}$$

674 where the last inequality is by $\mathbf{x}_t = \mathbf{x}_1$ and $\mathbf{y}_t = \mathbf{y}_1$ when $\mathbf{z}_t = \mathbf{z}_1$. Therefore suppose function f can
675 achieve a margin greater than $\log \left(16 \|\mathbf{G}\|_2^2 \|\mathbf{H}\|_2^2 (R^2 + 1)^2 B \tau^{-1} \eta T \right) \tau$, we have that the gradient

$$\begin{aligned}
&\left\| \nabla_{\mathbf{W}} \mathbb{E}_{\mathcal{D}^B} L(f, \tau) \right\|_F \\
&\leq 2 \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1) \cdot \mathbb{E} \left[\sum_{t \in [B]} \frac{\mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp \left([f(\mathbf{x}_1, \mathbf{y}_t) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)}{\sum_s \exp \left([f(\mathbf{x}_1, \mathbf{y}_s) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)} \right] \\
&\quad + 2 \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1) \cdot \mathbb{E} \left[\sum_{t \in [B]} \frac{\mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp \left([f(\mathbf{x}_t, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)}{\sum_s \exp \left([f(\mathbf{x}_s, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right)} \right] \\
&\leq 2 \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1) \cdot \mathbb{E} \left[\mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \sum_{t \in [B]} \exp \left([f(\mathbf{x}_1, \mathbf{y}_t) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right) \right] \\
&\quad + 2 \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1) \cdot \mathbb{E} \left[\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t \neq \mathbf{z}_1) \exp \left([f(\mathbf{x}_t, \mathbf{y}_1) - f(\mathbf{x}_1, \mathbf{y}_1)] / \tau \right) \right] \\
&\leq 0.25 \tau \|\mathbf{G}\|_2^{-1} \|\mathbf{H}\|_2^{-1} (R^2 + 1)^{-1} \eta^{-1} T^{-1}, \tag{I.1}
\end{aligned}$$

676 is very small. Now suppose the SGD trajectory start at $\mathbf{W}^{(0)} = 2 \log \left(16 \|\mathbf{G}\|_2^2 \|\mathbf{H}\|_2^2 (R^2 + \right.$
677 $\left. 1)^2 B \tau^{-1} \eta T \right) \cdot (\tau / \gamma) \mathbf{W}^*$. Obviously the score function with weight $\mathbf{W}^{(0)}$ achieve a margin
678 $2 \log \left(16 \|\mathbf{G}\|_2^2 \|\mathbf{H}\|_2^2 (R^2 + 1)^2 B \tau^{-1} \eta T \right) \tau$. Suppose there exists a time $t \leq T$ such that $\langle \mathbf{W}^{(t)} \mathbf{x}, \mathbf{y} \rangle$
679 can achieve margin larger than $3 \log \left(16 \|\mathbf{G}\|_2^2 \|\mathbf{H}\|_2^2 (R^2 + 1)^2 B \tau^{-1} \eta T \right) \tau$ or can achieve margin
680 larger than $\log \left(16 \|\mathbf{G}\|_2^2 \|\mathbf{H}\|_2^2 (R^2 + 1)^2 B \tau^{-1} \eta T \right) \tau$. Then there must exist a first time $t < t'$ such

681 that the margin at time t lies outside the range between $\log \left(16 \|\mathbf{G}\|_2^2 \|\mathbf{H}\|_2^2 (R^2 + 1)^2 B \tau^{-1} \eta T \right) \tau$
682 and $3 \log \left(16 \|\mathbf{G}\|_2^2 \|\mathbf{H}\|_2^2 (R^2 + 1)^2 B \tau^{-1} \eta T \right) \tau$. By definition of t (margin gap), we know that there
683 exist \mathbf{x}, \mathbf{y} such that $|\langle \mathbf{W}^{(t)} \mathbf{x}, \mathbf{y} \rangle - \langle \mathbf{W}^{(0)} \mathbf{x}, \mathbf{y} \rangle| > \tau$. On the other hand, we have that

$$\begin{aligned} |\langle \mathbf{W}^{(t)} \mathbf{x}, \mathbf{y} \rangle - \langle \mathbf{W}^{(0)} \mathbf{x}, \mathbf{y} \rangle| &\leq \|\mathbf{W}^{(t)} - \mathbf{W}^{(0)}\|_F \|\mathbf{x}\mathbf{y}^\top\|_F \\ &\leq 2 \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1) \|\mathbf{W}^{(t)} - \mathbf{W}^{(0)}\|_F \\ &\leq 2 \|\mathbf{G}\|_2 \|\mathbf{H}\|_2 (R^2 + 1) \cdot \eta T \cdot 0.25 \tau \|\mathbf{G}\|_2^{-1} \|\mathbf{H}\|_2^{-1} (R^2 + 1)^{-1} \eta^{-1} T^{-1} \\ &\leq 0.5 \tau, \end{aligned}$$

684 a contradiction! Therefore, such a t doesn't exist. The score function learned by SGD within T
685 iterations can't achieve a margin greater than $3 \log \left(16 \|\mathbf{G}\|_2^2 \|\mathbf{H}\|_2^2 (R^2 + 1)^2 B \tau^{-1} \eta T \right) \tau$. \square

686 **Theorem I.1** (Formal statement of Theorem 4.3). Under the same condition as Theorem 3.4, with
687 $\zeta = \mathbf{0}$. (This problem setting includes the special case considered in Theorem 4.2.) Let $\epsilon \leq$
688 $\lambda \gamma^2 \min p_k / (3200 \|\mathbf{H}\|_2^2)$ and $\tau \leq \gamma / \log(\gamma^2 \min p_k / (6400 B \|\mathbf{H}\|_2^2))$, within polynomial iterations,
689 we can find a score function \hat{f} with large margin. In particular, with a probability of at least 0.99, the
690 top-1 result gives the correct label with a margin of at least 0.5γ .

691 *Proof.* For simplicity, consider the case that we can access the population loss and its gradient, i.e.,
692 $n \rightarrow \infty$. The regularized loss then becomes,

$$L^{new} = L_{\mathcal{D}^B}(f, \tau) + \lambda \mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbf{h}(\mathbf{y})\|_2^2].$$

693 Since the new loss is still convex and even strongly convex. By applying the same technique in the
694 proof of the Theorem 3.4, within polynomial iterations, we can find $L^{new}(f, \tau, \lambda) \leq L^{new}(f^*, \tau, \lambda) +$
695 ϵ . Besides,

$$L^{new}(f^*, \tau, \lambda) = L_{\mathcal{D}^B}(f^*, \tau) \leq 2\mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \right) \right] + 2B \exp(-\gamma/\tau)$$

696 where the first equality is by plugging in $\mathbf{W}^* = \mathbf{H}(\mathbf{H}^\top \mathbf{H})^{-1} \mathbf{P}(\mathbf{G}^\top \mathbf{G})^{-1} \mathbf{G}^\top$, $\mathbf{g}(\mathbf{x}) = \mathbf{W}\mathbf{x}$, $\mathbf{h}(\mathbf{y}) =$
697 \mathbf{y} , the inequality is by Lemma G.7. Thus we have that

$$L_{\mathcal{D}^B}(f, \tau) + \lambda \mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbf{h}(\mathbf{y})\|_2^2] \leq 2\mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \right) \right] + \epsilon',$$

698 where $\epsilon' = \epsilon + 2B \exp(-\gamma/\tau)$. By (G.10) and (G.11), we know that $L_{\mathcal{D}^B}(f, \tau) \geq$
699 $2\mathbb{E} \left[\log \left(\sum_{t \in [B]} \mathbf{1}(\mathbf{z}_t = \mathbf{z}_1) \right) \right]$. Therefore, we can conclude that

$$\mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbf{h}(\mathbf{y})\|_2^2] \leq \epsilon' / \lambda \leq \gamma^2 \min p_k / (1600 \|\mathbf{H}\|_2^2),$$

700 where the last inequality is by choose $\epsilon \leq \lambda \gamma^2 \min p_k / (3200 \|\mathbf{H}\|_2^2)$ and $\tau \leq$
701 $\gamma / \log(\gamma^2 \min p_k / (6400 B \|\mathbf{H}\|_2^2))$. Then by Chebyshev's inequality, for any \mathbf{z} , with probability
702 $1 - 0.01$ we have $\|\mathbf{g}(\mathbf{x}) - \mathbf{h}(\mathbf{y})\|_2 \leq \sqrt{100 \max p_k^{-1} \mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbf{h}(\mathbf{y})\|_2^2]} \leq \gamma / (4 \|\mathbf{H}\|_2)$. Then for
703 any \mathbf{y}' that has the different shared feature from \mathbf{y} (i.e., $\mathbf{z}' \neq \mathbf{z}$) we have that

$$\begin{aligned} &\langle \mathbf{g}(\mathbf{x}), \mathbf{h}(\mathbf{y}') \rangle - \langle \mathbf{g}(\mathbf{x}), \mathbf{h}(\mathbf{y}) \rangle \\ &\leq \langle \mathbf{h}(\mathbf{y}), \mathbf{h}(\mathbf{y}') \rangle - \langle \mathbf{h}(\mathbf{y}), \mathbf{h}(\mathbf{y}) \rangle + \|\mathbf{g}(\mathbf{x}) - \mathbf{h}(\mathbf{y})\|_2 \cdot (\|\mathbf{h}(\mathbf{y}')\|_2 + \|\mathbf{h}(\mathbf{y})\|_2) \\ &\leq -\gamma + \gamma/2 \\ &\leq -\gamma/2, \end{aligned}$$

704 where the first inequality is by triangle inequality, the second inequality is by $\|\mathbf{g}(\mathbf{x}) - \mathbf{h}(\mathbf{y})\|_2 \leq$
705 $\gamma / (4 \|\mathbf{H}\|_2)$ and $\|\mathbf{h}(\mathbf{y}')\|_2 = \|\mathbf{h}(\mathbf{y})\|_2 \leq \|\mathbf{H}\|_2$ since $\zeta = \mathbf{0}$. \square