

---

# Multi-Domain Long-Tailed Learning by Augmenting Disentangled Representations

---

Anonymous Author(s)

Affiliation

Address

email

## 1 Introduction

2 Deep classification models can struggle when the number of examples per class varies dramatically [5,  
3 45]. This *long-tailed* setting arises frequently in practice, such as wildlife recognition [5]. Classifiers  
4 tend to be biased towards majority classes and perform poorly on class-balanced test distributions, i.e.  
5 when there is a shift in the label distribution between training and test. Existing approaches focus  
6 on single-domain long-tailed learning, while we study *multi-domain long-tailed learning*, where  
7 each domain has its own long-tailed distribution and the classifiers need to handle distribution shift  
8 amidst class imbalance. Here, we focus on two types of distribution shift: subpopulation shift and  
9 domain shift. In subpopulation shift, we train a model on data from multiple domains and evaluate  
10 the model on a test set with balanced domain-class pairs. A machine learning model trained on the  
11 entire population may fail on the test set when this correlation does not hold anymore. In domain  
12 shift, we expect the trained model to generalize well to completely new test domains.

13 Prior long-tailed classification methods work well in single-domain settings, but may perform poorly  
14 when the test data is from underrepresented domains or novel domains. Meanwhile, invariant learning  
15 approaches alleviate cross-domain performance gaps by learning representations or predictors that  
16 are invariant across different domains [3, 23]. Yet, these approaches are mostly evaluated in class-  
17 balanced settings, where models must be trained on plenty of examples from each class even if  
18 augmentation strategies are applied [40]. With multi-domain long-tailed data, learning a class-  
19 unbiased domain-invariant model is not trivial since the imbalance can exist within a domain or across  
20 domains. We aim to address these challenges in this work, leading to a novel method named **TALLY**.

21 **TALLY** empower augmentation to balance examples over domains and classes by decomposing and  
22 reassembling example pairs, combining the class-relevant semantic information of one example with  
23 the domain-associated nuisances of another. Specifically, **TALLY** first decouples the representation  
24 of each example into semantic information and nuisances with instance normalization. To further  
25 mitigate the effects of nuisances, we first average out domain information over examples of the same  
26 class and construct class prototype representations. Each semantic representation is then linearly  
27 interpolated with a corresponding class prototype, leading to the prototype-enhanced semantic repre-  
28 sentation. The domain-associated factors are similarly interpolated with class-agnostic domain factors  
29 to improve training stability and remove noise. Finally, **TALLY** produces augmented representations  
30 to benefit the training process by reassembling the prototype-enhanced semantic representation  
31 and domain-associated nuisances among examples. To further achieve balanced augmentation, we  
32 additionally propose a selective balanced sampling strategy to draw example pairs for augmentation.

33 In summary, our major contributions are: we investigate and formalize an important yet less explored  
34 problem – multi-domain long-tailed learning, and propose an effective augmentation algorithm called  
35 **TALLY** to simultaneously address the class-imbalance issue and learn domain-invariant predictors.  
36 We empirically demonstrate the effectiveness of **TALLY** under subpopulation shift and domain shift.  
37 We observe that **TALLY** outperforms both prior single-domain long-tailed learning and domain-  
38 invariant learning approaches, with a 5.18% error decrease over all datasets. Furthermore, **TALLY** is  
39 capable of capturing stronger invariant predictors compared with prior invariant learning approaches.  
40

## 2 Preliminaries and Method

### 2.1 Multi-Domain Long-Tailed Learning.

In this paper, we investigate the setting where one predicts the class label  $y \in \mathcal{C}$  based on the input feature  $x \in \mathcal{X}$ , where  $\mathcal{C} = \{1, \dots, C\}$ . Given a machine learning model  $f$  parameterized by parameter  $\theta$  and a loss function  $\ell$ , empirical risk minimization (ERM) trains such a model by minimizing average loss over all training examples as

$$\min_{\theta} \mathbb{E}_{(x,y) \sim P^{tr}} [\ell(f_{\theta}(x), y)], \quad (1)$$

which works well when the label distribution is approximately uniform. In multi-domain long-tailed learning, the overall data distribution is drawn from a set of domains  $\mathcal{D} = \{1, \dots, D\}$  and each domain  $d$  is associated with a class-imbalanced dataset  $\{(x_i, y_i, d)\}_{i=1}^{N_d}$  drawn from domain-specific distribution  $p_d$ . Following [2, 19], both training and test distribution can be formulated as a mixture distribution over domain space  $\mathcal{D}$ , i.e.,  $P^{tr} = \sum_{d=1}^D \eta_d^{tr} P_d^{tr}$  and  $P^{ts} = \sum_{d=1}^D \eta_d^{ts} P_d^{ts}$ . The corresponding training and test domains are  $\mathcal{D}^{tr} = \{d \in \mathcal{D} | \eta_d^{tr} > 0\}$  and  $\mathcal{D}^{ts} = \{d \in \mathcal{D} | \eta_d^{ts} > 0\}$ , respectively, where  $\eta_d^{tr}$  and  $\eta_d^{ts}$  represent the mixture probability. For each domain  $d$ , we define the number of training examples in each class as  $\{n_{1,d}^{tr}, \dots, n_{C,d}^{tr}\}$ , sorted by cardinality. The imbalance ratio  $\rho^{tr}$  is extended to domain-level ratio as  $\rho_d^{tr} = n_{C,d}^{tr} / n_{1,d}^{tr}$ . During test time, we consider two kinds of test distributions, corresponding to two categories of distribution shifts – subpopulation shift and domain shift. In subpopulation shift, the test domains have been observed during training time, but the test distribution is class-balanced and domain-balanced, i.e.,  $\mathcal{D}^{ts} \subseteq \mathcal{D}^{tr}$  and  $\{\eta_d^{ts} = 1/|\mathcal{D}^{ts}| \forall d \in \mathcal{D}^{ts}\}$ . In domain shift, the test domains are disjoint from the training domains, i.e.,  $\mathcal{D}^{tr} \cap \mathcal{D}^{ts} = \emptyset$ .

### 2.2 Detailed Descriptions of TALLY

To improve robustness in multi-domain long-tailed learning, we would like method that can learn class-unbiased domain-invariant representations. To accomplish this, we introduce TALLY to do balanced augmentation over classes and domains.

#### Representation Disentanglement and Reassembly

As described above, TALLY reassembles augmented examples from pairs of examples by combining the semantic representation of one with the domain-related nuisance factors of the other. Motivated by style transfer [15], we use instance normalization (InstanceNorm) to perform the required disentanglement of semantic and nuisance information. Concretely, given an example  $(x, y, d)$  we denote the hidden representation at layer  $r$  as  $s = f^r(x) \in \mathbb{R}^{C \times H \times W}$ , where  $C$ ,  $H$ , and  $W$  denote channel, height, and width dimensions, respectively. Ignoring affine parameters, InstanceNorm normalizes the example as:

$$z(s) = \text{InstanceNorm}(s) = \frac{s - \mu(s)}{\sigma(s)}, \quad \text{where } z(s), \mu(s), \sigma(s) \in \mathbb{R}^C \quad (2)$$

Following Huang and Belongie [15], we treat the normalized example  $z(s)$  as the semantic representation, and regard  $\mu(s)$  and  $\sigma(s)$  as the domain-associated nuisances.

After decoupling representations, we produce an augmented representation from a pair of examples  $(x_i, y_i, d_i)$  and  $(x_j, y_j, d_j)$  by swapping semantic representations and domain-associated nuisances:

$$\tilde{s} = \sigma(s_j) \left( \frac{s_i - \mu(s_i)}{\sigma(s_i)} \right) + \mu(s_j), \quad \tilde{y} = y_i. \quad (3)$$

Since the semantic content of the augmented representation  $\tilde{s}$  is from example  $(x_i, y_i, d_i)$ , we label our augmented example with  $\tilde{y} = y_i$ . By reassembling disentangled representations, we can augment representations for minority domains or minority classes.

**Selective Balanced Sampling.** In the process of representation disentanglement and reassembly, finding a suitable strategy of sampling examples from the training distribution is crucial to solving the class-domain imbalance problem. In multi-domain long-tailed learning, the most straightforward way is up-sampling examples from minority domain-class groups, which is named *balanced sampling*

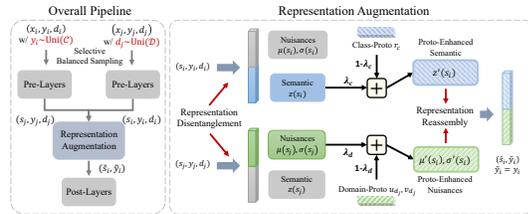


Figure 1: An illustration of TALLY.

88 here. In practice, for each example  $(x_i, y_i, d_i)$ , the label  $y_i$  and domain  $d_j$  are uniformly sampled a  
 89 joint uniform distribution over all domain-class combinations, i.e.,  $(y_i, d_i) \sim \text{Uniform}(\mathcal{C}, \mathcal{D})$ .

90 However, to transfer the knowledge between different domain-class groups in TALLY, using such a  
 91 sampling strategy may overemphasize the importance of minority domain-class groups. To augment  
 92 minority groups, balanced sampling tends to repeatedly draw examples from the same minority  
 93 group. We do not expect this because of two reasons: first, it limits the sample diversity in knowledge  
 94 transfer; second, minority groups typically perform worse than majority groups, which may make  
 95 the knowledge transfer less reliable. Hence, we propose a selective balanced sampling strategy in  
 96 TALLY. Concretely, for a pair of examples  $(x_i, y_i, d_i)$  and  $(x_j, y_j, d_j)$ , the label  $y_i$  of example  $i$  is  
 97 uniformly sampled from all classes ( $y_i \sim \text{Uniform}(\mathcal{C})$ ) and the domain  $d_j$  of example  $j$  is uniformly  
 98 sampled from all domains ( $d_j \sim \text{Uniform}(\mathcal{D})$ ).

99 **Prototype-guided Invariant Learning.** Since the semantic representation  $z(s)$  (Eqn. 2) should  
 100 contain only class-relevant information, it should ideally be *domain-invariant*. However, per-instance  
 101 statistics can be noisy and instance normalization may not perfectly disentangle the semantic infor-  
 102 mation from the domain-related nuisances. To improve robustness, we can “average out” domain in-  
 103 formation over many examples of the same class from different domains. However, merely averaging  
 104 over examples would remove the diversity that distinguishes different examples of the same class. We  
 105 balance diversity and domain-invariance by interpolating  $z(s)$  with the corresponding *class prototype*  
 106 *representation*. We define the class prototype representation  $r_c$  as the average *semantic* representation  
 107 over examples belonging to class  $c$  regardless of domain:  $r_c = \frac{1}{n_c^{tr}} \sum_{i=1}^{n_c^{tr}} z(s_i) = \frac{1}{n_c^{tr}} \sum_{i=1}^{n_c^{tr}} \frac{s_i - \mu(s_i)}{\sigma(s_i)}$ .

108 For each example  $(x_i, y_i, d_i)$  with  $y_i = c$ , we obtain the prototype-enhanced semantic representation  
 109 by linearly interpolating  $z(s_i)$  with the corresponding class prototype  $r_c$ :

$$z'(s_i) = \lambda_c z(s_i) + (1 - \lambda_c) r_c, \quad (4)$$

110 where  $\lambda_c \sim \text{Beta}(\alpha_c, \alpha_c)$  is the interpolation coefficient. By applying this class prototype-based  
 111 interpolation strategy, we are capable of capturing invariant knowledge and keeping the diversity of  
 112 instance-level semantic representation when swapping information.

113 We also desire that the disentangled  $\mu(s)$  and  $\sigma(s)$  (Eqn. 2) contain only domain-related nuisance  
 114 information. However, for similar reasons as with  $z(s)$ , they may still contain some class-related  
 115 semantic information which we would like to remove by “averaging out.” In this case, we remove  
 116 semantic information by averaging over examples from *different classes* within the *same domain*:  
 117  $u_d = \frac{1}{n_d^{tr}} \sum_{i=1}^{n_d^{tr}} \mu(s_i)$ ,  $v_d = \frac{1}{n_d^{tr}} \sum_{i=1}^{n_d^{tr}} \sigma(s_i)$ , where  $n_d^{tr}$  represents the number of training examples in  
 118 domain  $d$ . Then, for each example, we linearly interpolate its domain-associated nuisances with the  
 119 above class-agnostic nuisances as:

$$\mu'(x_i) = \lambda_d \mu(x) + (1 - \lambda_d) u_d, \quad \sigma'(x_i) = \lambda_d \sigma(x) + (1 - \lambda_d) v_d, \quad (5)$$

120 where the interpolation ratio is  $\lambda_d \sim \text{Beta}(\alpha_d, \alpha_d)$ .

121 By replacing the original semantic representation and domain-associated nuisances in Eqn. 3 with the  
 122 prototype-guided ones, we obtain the enhanced augmented representation as follows:

$$\tilde{s}' = \sigma'(s_j) z'(s_i) + \mu'(s_j), \quad \tilde{y}' = y_i. \quad (6)$$

123 Finally, we replace the original training data with the augmented ones. We summarize the overall  
 124 framework of TALLY in Algorithm 1 in Appendix.

### 125 3 Experiments

126 In this section, we conduct extensive experiments to evaluate how TALLY performs. To answer  
 127 Q1, we compare TALLY to two categories of algorithms. The first category includes single-domain  
 128 long-tailed learning methods such as Focal [24], LDAM [6], CRT [17], MiSLAS [47], and Remix [8].  
 129 Focal, LDAM, and CRT are up-weighting or up-sampling approaches, while MiSLAS and Remix are  
 130 data augmentation strategies. The second category includes approaches for improving robustness to  
 131 distribution shift: IRM [3], GroupDRO [29], LISA [40], MixStyle [50], DDG [43], and BODA [39].  
 132 Follow Yang et al. [39], we use a ResNet-50 architecture for all algorithms, and detail the baselines in  
 133 Appendix B. All hyperparameters are selected via cross-validation. Due to space limitation, we only  
 134 show ERM and top-5 baselines here and put full results in Appendix. We also provide comprehensive  
 135 analysis to understand the results in Appendix E.

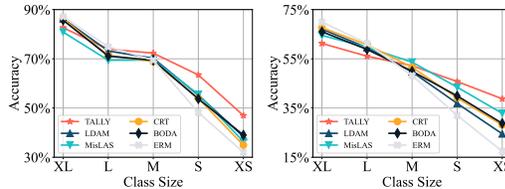
Table 1: Results of subpopulation shifts and domain shifts on synthetic data (Full table Appx. C.3).

	Subpopulation Shift				Domain Shift			
	VLCS-LT	PACS-LT	OH-LT	DN-LT	VLCS-LT	PACS-LT	OH-LT	DN-LT
ERM	73.33%	90.40%	61.07%	44.33%	67.62%	76.27%	51.95%	33.21%
Focal	74.83%	90.44%	62.57%	47.35%	69.38%	75.29%	54.03%	35.23%
MiSLAS	71.83%	90.99%	61.38%	49.15%	68.64%	77.94%	52.86%	36.18%
CORAL	71.67%	88.22%	59.10%	43.92%	66.54%	75.62%	50.74%	33.44%
MixStyle	74.30%	91.55%	62.26%	43.59%	67.75%	79.78%	52.47%	33.71%
BODA	74.83%	91.03%	62.79%	47.61%	69.63%	78.81%	53.32%	35.85%
<b>TALLY (ours)</b>	<b>76.83%</b>	<b>92.38%</b>	<b>67.00%</b>	<b>50.15%</b>	<b>70.60%</b>	<b>81.55%</b>	<b>55.69%</b>	<b>36.45%</b>

136 **3.1 Evaluation on Long-Tailed Variants of Domain Generalization Benchmarks**

137 **Datasets.** We curate four *multi-domain long-tailed* datasets by modifying four existing domain-  
 138 generalization benchmarks: VLCS [11], PACS [22], OfficeHome [33], and DomainNet [27]. We  
 139 modify the prior datasets by removing training examples so that each domain has a long-tailed  
 140 label distribution (overall imbalance ratio: 50) and call the resulting datasets **VLCS-LT**, **PACS-LT**,  
 141 **OfficeHome-LT**, and **DomainNet-LT**. See Appendix C for more details and evaluation protocol.

142 **Results.** The overall performance of TALLY and  
 143 prior methods for tackling subpopulation shift  
 144 and domain shift is reported in Table 1. For  
 145 subpopulation shift, we report the average per-  
 146 formance over all domains. We observe that  
 147 TALLY consistently outperforms all methods,  
 148 verifying its effectiveness in improving the ro-  
 149 bustness to subpopulation shifts. In addition, Fig-  
 150 ure 2 shows performance broken down by class  
 151 size for OfficeHome-LT and DomainNet-LT un-  
 152 der subpopulation shift, where we split all classes  
 153 into five levels according to their cardinality. We compare TALLY with ERM, and four strongest  
 154 baselines. The results show that TALLY’s performance improvements arise from larger improvements  
 155 on smaller classes rather than performance improvements across the board, hence indicating that it is  
 156 particularly well-suited for class-imbalanced problems.



(a) : OfficeHome-LT (b) : DomainNet-LT

Figure 2: Performance w.r.t. Class Size. XL and XS represent the largest and smallest classes.

157 **3.2 Evaluation on Naturally Imbalanced Multi-Domain Data**

158 **Datasets.** To further evaluate TALLY and prior methods,  
 159 we study two multi-domain datasets that are naturally  
 160 imbalanced: Terra Incognita (TerraInc) [4] and iWild-  
 161 Cam [5], both of which aim to classify wildlife across  
 162 different camera traps. More details of these datasets  
 163 and class distribution are described in Appendix D. To  
 164 better capture performance on rare species, we use macro  
 165 F1 score as the primary evaluation metric following Koh  
 166 et al. [19], but we also report average accuracy. We list  
 167 all hyperparameters in Appendix D.2.

168 **Results.** We report the results over all test domains in  
 169 Table 2. The conclusions are largely consistent with the  
 170 results from Sec. 3.1, where TALLY consistently improves the performance over all baselines and  
 171 enhances the robustness of multi-domain long-tailed learning. The superiority of TALLY over prior  
 172 augmentation techniques is further evidence of the effectiveness of balanced augmentation.

Table 2: Results of Domain Shifts on Real-world Data (full results: Appx. D.3)

	TerraInc		iWildCam	
	Macro F1	Acc	Macro F1	Acc
ERM	42.35%	54.81%	32.0%	69.0%
Focal	43.54%	56.62%	33.2%	74.7%
MiSLAS	40.68%	52.96%	30.5%	59.8%
CORAL	45.43%	58.10%	32.8%	73.3%
MixStyle	44.73%	57.55%	32.4%	74.9%
BODA	44.47%	57.52%	32.9%	70.5%
<b>TALLY (ours)</b>	<b>46.23%</b>	<b>59.89%</b>	<b>34.4%</b>	<b>73.4%</b>

173 **4 Conclusion**

174 In this paper we investigate multi-domain imbalanced learning, a natural extension of classical single-  
 175 domain imbalanced learning. We propose a novel balanced augmentation algorithm called TALLY  
 176 to achieve robust imbalanced learning that can overcome distribution shifts. TALLY introduces a  
 177 prototype enhanced disentanglement procedure for separating semantic and nuisance information,  
 178 and then mixes the enhanced semantic and domain-associated nuisance information among examples.  
 179 The results demonstrate its effectiveness of TALLY.

## References

- 180
- 181 [1] Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Yoshua Bengio, Ioannis Mitliagkas, and Irina  
182 Rish. Invariance principle meets information bottleneck for out-of-distribution generalization.  
183 2021.
- 184 [2] Isabela Albuquerque, João Monteiro, Mohammad Darvishi, Tiago H Falk, and Ioannis  
185 Mitliagkas. Generalizing to unseen domains via distribution matching. *arXiv preprint*  
186 *arXiv:1911.00804*, 2019.
- 187 [3] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk mini-  
188 mization. *arXiv preprint arXiv:1907.02893*, 2019.
- 189 [4] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In *Proceedings*  
190 *of the European conference on computer vision (ECCV)*, pages 456–473, 2018.
- 191 [5] Sara Beery, Elijah Cole, and Arvi Gjoka. The iwildcam 2020 competition dataset. *arXiv*  
192 *preprint arXiv:2004.10340*, 2020.
- 193 [6] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Archiga, and Tengyu Ma. Learning imbalanced  
194 datasets with label-distribution-aware margin loss. *Advances in neural information processing*  
195 *systems*, 32, 2019.
- 196 [7] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote:  
197 synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:  
198 321–357, 2002.
- 199 [8] Hsin-Ping Chou, Shih-Chieh Chang, Jia-Yu Pan, Wei Wei, and Da-Cheng Juan. Remix:  
200 rebalanced mixup. In *European Conference on Computer Vision*, pages 95–110. Springer, 2020.
- 201 [9] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based  
202 on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer*  
203 *vision and pattern recognition*, pages 9268–9277, 2019.
- 204 [10] Andrew Estabrooks, Taeho Jo, and Nathalie Japkowicz. A multiple resampling method for  
205 learning from imbalanced data sets. *Computational intelligence*, 20(1):18–36, 2004.
- 206 [11] Chen Fang, Ye Xu, and Daniel N Rockmore. Unbiased metric learning: On the utilization of  
207 multiple datasets and web images for softening bias. In *Proceedings of the IEEE International*  
208 *Conference on Computer Vision*, pages 1657–1664, 2013.
- 209 [12] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. *arXiv preprint*  
210 *arXiv:2007.01434*, 2021.
- 211 [13] Ruocheng Guo, Pengchuan Zhang, Hao Liu, and Emre Kiciman. Out-of-distribution pre-  
212 diction with invariant risk minimization: The limitation and an effective fix. *arXiv preprint*  
213 *arXiv:2101.07732*, 2021.
- 214 [14] Youngkyu Hong, Seungju Han, Kwanghee Choi, Seokjun Seo, Beomsu Kim, and Buru Chang.  
215 Disentangling label distribution for long-tailed visual recognition. In *Proceedings of the*  
216 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6626–6636, 2021.
- 217 [15] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance  
218 normalization. In *Proceedings of the IEEE international conference on computer vision*, pages  
219 1501–1510, 2017.
- 220 [16] Muhammad Abdullah Jamal, Matthew Brown, Ming-Hsuan Yang, Liqiang Wang, and Boqing  
221 Gong. Rethinking class-balanced methods for long-tailed visual recognition from a domain  
222 adaptation perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
223 *Pattern Recognition*, pages 7610–7619, 2020.
- 224 [17] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and  
225 Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. 2020.

- 226 [18] Kia Khezeli, Arno Blaas, Frank Soboczenski, Nicholas Chia, and John Kalantari. On invariance  
227 penalties for risk minimization. *arXiv preprint arXiv:2106.09777*, 2021.
- 228 [19] Pang Wei Koh, Shiori Sagawa, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani,  
229 Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, et al. Wilds: A  
230 benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*,  
231 pages 5637–5664. PMLR, 2021.
- 232 [20] Masanori Koyama and Shoichiro Yamaguchi. Out-of-distribution generalization with maximal  
233 invariant predictor. *arXiv preprint arXiv:2008.01883*, 2020.
- 234 [21] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui  
235 Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrap-  
236 olation (rex). In *International Conference on Machine Learning*, pages 5815–5826. PMLR,  
237 2021.
- 238 [22] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier  
239 domain generalization. In *Proceedings of the IEEE international conference on computer vision*,  
240 pages 5542–5550, 2017.
- 241 [23] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with  
242 adversarial feature learning. In *Proceedings of the IEEE Conference on Computer Vision and*  
243 *Pattern Recognition*, pages 5400–5409, 2018.
- 244 [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense  
245 object detection. In *Proceedings of the IEEE international conference on computer vision*,  
246 pages 2980–2988, 2017.
- 247 [25] Jialun Liu, Yifan Sun, Chuchu Han, Zhaopeng Dou, and Wenhui Li. Deep representation  
248 learning on long-tailed data: A learnable embedding augmentation perspective. In *Proceedings*  
249 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2970–2979,  
250 2020.
- 251 [26] Xu-Ying Liu, Jianxin Wu, and Zhi-Hua Zhou. Exploratory undersampling for class-imbalance  
252 learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 39(2):  
253 539–550, 2008.
- 254 [27] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment  
255 matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international*  
256 *conference on computer vision*, pages 1406–1415, 2019.
- 257 [28] Alexander Robey, George J Pappas, and Hamed Hassani. Model-based domain generalization.  
258 *Advances in Neural Information Processing Systems*, 34:20210–20229, 2021.
- 259 [29] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally  
260 robust neural networks for group shifts: On the importance of regularization for worst-case  
261 generalization. In *ICLR*, 2020.
- 262 [30] Yang Shu, Zhangjie Cao, Chenyu Wang, Jianmin Wang, and Mingsheng Long. Open do-  
263 main generalization with domain-augmented meta-learning. In *Proceedings of the IEEE/CVF*  
264 *Conference on Computer Vision and Pattern Recognition*, pages 9624–9633, 2021.
- 265 [31] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation.  
266 In *European conference on computer vision*, pages 443–450. Springer, 2016.
- 267 [32] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain  
268 confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- 269 [33] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan.  
270 Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE*  
271 *conference on computer vision and pattern recognition*, pages 5018–5027, 2017.
- 272 [34] Yu-Xiong Wang, Deva Ramanan, and Martial Hebert. Learning to model the tail. *Advances in*  
273 *Neural Information Processing Systems*, 30, 2017.

- 274 [35] Yufei Wang, Haoliang Li, and Alex C Kot. Heterogeneous domain generalization via domain  
275 mixup. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal*  
276 *Processing (ICASSP)*, pages 3622–3626. IEEE, 2020.
- 277 [36] Liuyu Xiang, Guiguang Ding, and Jungong Han. Learning from multiple experts: Self-paced  
278 knowledge distillation for long-tailed classification. In *European Conference on Computer*  
279 *Vision*, pages 247–263. Springer, 2020.
- 280 [37] Minghao Xu, Jian Zhang, Bingbing Ni, Teng Li, Chengjie Wang, Qi Tian, and Wenjun Zhang.  
281 Adversarial domain adaptation with domain mixup. In *Proceedings of the AAAI Conference on*  
282 *Artificial Intelligence*, volume 34, pages 6502–6509, 2020.
- 283 [38] Shen Yan, Huan Song, Nanxiang Li, Lincan Zou, and Liu Ren. Improve unsupervised domain  
284 adaptation with mixup training. *arXiv preprint arXiv:2001.00677*, 2020.
- 285 [39] Yuzhe Yang, Hao Wang, and Dina Katabi. On multi-domain long-tailed recognition, generaliza-  
286 tion and beyond. *arXiv preprint arXiv:2203.09513*, 2022.
- 287 [40] Huaxiu Yao, Yu Wang, Sai Li, Linjun Zhang, Weixin Liang, James Zou, and Chelsea  
288 Finn. Improving out-of-distribution robustness via selective augmentation. *arXiv preprint*  
289 *arXiv:2201.00299*, 2022.
- 290 [41] Xi Yin, Xiang Yu, Kihyuk Sohn, Xiaoming Liu, and Manmohan Chandraker. Feature transfer  
291 learning for face recognition with under-represented data. In *Proceedings of the IEEE/CVF*  
292 *conference on computer vision and pattern recognition*, pages 5704–5713, 2019.
- 293 [42] Xiangyu Yue, Yang Zhang, Sicheng Zhao, Alberto Sangiovanni-Vincentelli, Kurt Keutzer, and  
294 Boqing Gong. Domain randomization and pyramid consistency: Simulation-to-real general-  
295 ization without accessing target domain data. In *Proceedings of the IEEE/CVF International*  
296 *Conference on Computer Vision*, pages 2100–2110, 2019.
- 297 [43] Hanlin Zhang, Yi-Fan Zhang, Weiyang Liu, Adrian Weller, Bernhard Schölkopf, and Eric P  
298 Xing. Towards principled disentanglement for domain generalization. In *CVPR*, 2022.
- 299 [44] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond  
300 empirical risk minimization. 2018.
- 301 [45] Yifan Zhang, Bingyi Kang, Bryan Hooi, Shuicheng Yan, and Jiashi Feng. Deep long-tailed  
302 learning: A survey. *arXiv preprint arXiv:2110.04596*, 2021.
- 303 [46] Zizhao Zhang and Tomas Pfister. Learning fast sample re-weighting without reward data. In  
304 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 725–734,  
305 2021.
- 306 [47] Zhisheng Zhong, Jiequan Cui, Shu Liu, and Jiaya Jia. Improving calibration for long-tailed  
307 recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
308 *Recognition*, pages 16489–16498, 2021.
- 309 [48] Allan Zhou, Fahim Tajwar, Alexander Robey, Tom Knowles, George J Pappas, Hamed Hassani,  
310 and Chelsea Finn. Do deep networks transfer invariances across classes? 2022.
- 311 [49] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. Bbn: Bilateral-branch network  
312 with cumulative learning for long-tailed visual recognition. In *Proceedings of the IEEE/CVF*  
313 *conference on computer vision and pattern recognition*, pages 9719–9728, 2020.
- 314 [50] Fan Zhou, Zhuqing Jiang, Changjian Shui, Boyu Wang, and Brahim Chaib-draa. Domain  
315 generalization with optimal transport and metric learning. *arXiv preprint arXiv:2007.10573*,  
316 2020.
- 317 [51] Kaiyang Zhou, Yongxin Yang, Timothy Hospedales, and Tao Xiang. Deep domain-adversarial  
318 image generation for domain generalisation. In *Proceedings of the AAAI Conference on Artificial*  
319 *Intelligence*, volume 34, pages 13025–13032, 2020.

---

**Algorithm 1** TALLY Training Process

---

**Require:** learning rates  $\eta$ ; warm start epochs  $T_0$ ; prototype momentum  $\gamma$ ; model  $f_\theta(\cdot)$  with hidden representation  $f_\theta^r(\cdot)$  at layer  $r$ ; Dataset  $\mathcal{D}^{tr} = \{(x, y, d)\}$

- 1: Initialize domain-agnostic prototypes  $\{r_c^{(0)}\}_{c=1}^C$  and class-agnostic statistics  $\{(u_d^{(0)}, v_d^{(0)})\}_{d=1}^D$
- 2: Train  $f_\theta$  with ERM for  $t < T_0$
- 3: **for**  $t = T_0$  to  $T$  **do**
- 4:    $d_i, d_j \sim \text{Uniform}(\mathcal{D}), y_i, y_j \sim \text{Uniform}(\mathcal{C})$     $\triangleright$  Randomly sample domains and classes
- 5:    $(x_i, y_i, d_i) \sim \{\mathcal{D}^{tr} | y = y_i, d = d_i\}, (x_j, y_j, d_j) \sim \{\mathcal{D}^{tr} | y = y_j, d = d_j\}$
- 6:    $(s_i, s_j) \leftarrow (f^r(x_i), f^r(x_j))$     $\triangleright$  Compute hidden representations
- 7:    $z(s_i) \leftarrow \text{InstanceNorm}(s_i)$     $\triangleright$  Disentangle semantic factor (Eqn. 2)
- 8:    $z'(s_i) \leftarrow \lambda_{y_i} z(s_i) + (1 - \lambda_{y_i}) r_{y_i}^t$     $\triangleright$  Enhance semantic factor (Eqn. 4)
- 9:    $(\mu'(s_j), \sigma'(s_j)) \leftarrow \lambda_{y_j} (\mu(s_j), \sigma(s_j)) + (1 - \lambda_{y_j}) (u_{y_j}^t, v_{y_j}^t)$     $\triangleright$  Enhance nuisances (Eqn. 5)
- 10:    $(\tilde{s}', \tilde{y}') \leftarrow (\sigma'(s_j) z'(s_i) + \mu'(s_j), y_i)$     $\triangleright$  Generate augmented example (Eqn. 6)
- 11:   Optimize  $[\ell(f_\theta^{L-r}(\tilde{s}'), \tilde{y}')]$     $\triangleright$  Train on augmented example
- 12:   Estimate the current prototypes and feature statistics  $\{r_c\}_{c=1}^C, \{(u_d, v_d)\}_{d=1}^D$
- 13:   **for**  $c = 1$  to  $C$  **do**
- 14:      $r_c^{(t+1)} \leftarrow \gamma r_c^t + (1 - \gamma) r_c$     $\triangleright$  Update domain-agnostic prototypes
- 15:   **for**  $d = 1$  to  $D$  **do**
- 16:      $(u_d^{(t+1)}, v_d^{(t+1)}) \leftarrow \gamma (u_d^t, v_d^t) + (1 - \gamma) (u_d, v_d)$     $\triangleright$  Update class-agnostic statistics

---

## 320 A Related Work

### 321 A.1 Long-Tailed Learning

322 Training a well-performed machine learning model on class-imbalanced data has been widely studied.  
323 A typical setting of imbalanced learning is the long-tailed class distribution, where the model can be  
324 easily biased towards majority classes [45]. A lot of approaches have been proposed under this setting,  
325 including over-sampling minority classes or under-sampling majority classes [7, 10, 17, 26, 46],  
326 adjusting loss functions or logits for different classes during training [6, 9, 14, 16, 24], transferring  
327 knowledge from head classes to tail classes [34, 25, 41, 48], directly augmenting tail classes [8, 17, 47],  
328 and ensembling models with different sampling or loss weighting strategies [36, 49]. Unlike single-  
329 domain imbalanced learning, Yang et al. [39] targets on the multi-domain imbalanced learning  
330 scenario by encouraging invariant representation learning with a domain-class calibrated regularizer.  
331 However, BODA focuses on subpopulation shift with the imbalanced distribution for each domain,  
332 while the overall distribution among all classes are relatively balanced. TALLY instead studies  
333 more kinds of distribution shifts with conceptually different direction to alleviate domain-associated  
334 nuisances via balanced augmentation. It relaxes the explicit constraint on internal representations and  
335 leads to stronger empirical performance.

### 336 A.2 Domain Generalization and Out-of-Distribution Robustness

337 To improve out-of-distribution robustness, one line of works aims to learn domain-invariant rep-  
338 resentations by 1) minimizing the discrepancy of feature representations across all training do-  
339 mains [23, 31, 32, 50]; 2) leveraging domain augmentation methods to generate more training  
340 domains and improve the consistency of feature representations between the original and augmented  
341 domains [30, 35, 37, 38, 42, 51]; 3) disentangling feature representations to semantic and domain-  
342 varying ones and minimizing the semantic differences across training domains [28, 43]. Another line  
343 of works focuses on strengthening the correlations between representations and labels, leading to  
344 stronger invariant predictors. These works introduce various regularizers in learning invariant predic-  
345 tors, including minimizing the variances of risks across domains [21], encouraging a predictor that  
346 performs well over all domains [1, 3, 13, 18], and matching the gradient across different domains [20].  
347 Besides explicitly involving regularizers, data interpolation is also a promising approach for learning  
348 invariant predictors [40, 50]. Unlike previous augmentation methods that require sufficient training  
349 examples for each class to learn invariance, TALLY tackles the class-imbalanced issue in domain gen-  
350 eralization and employs a domain-balanced augmentation strategy to learn class-unbiased invariant  
351 representation.

## 352 **B Detailed Description of Baselines**

353 In this paper, we compare TALLY with two types of approaches: long-tailed classification methods  
354 and invariant learning approaches. We detail these methods here:

### 355 **B.1 Long-tailed Classification Methods**

356 We compare TALLY with Focal [24], LDAM [6], CRT [17], MiSLAS [47], and Remix [8]. Here,  
357 Focal and LDAM up-weight the loss for minority classes. CRT uses up-sampling strategy to fine-tune  
358 the classifier. MiSLAS and Remix modify the vanilla mixup [44] and make it suitable to long-tailed  
359 distribution.

### 360 **B.2 Invariant Learning**

361 We further compare TALLY with invariant learning approaches, i.e., IRM [3], GroupDRO [29],  
362 LISA [40], MixStyle [50], DDG [43], and BODA [39]. IRM learns invariant predictors that perform  
363 well across different domains. GroupDRO optimizes the worst-domain loss. LISA cancels out  
364 domain-associated information by mixing examples with the same label but different domains.  
365 MixStyle decomposes the feature representation into content information and style information. It  
366 then mixes the style information and generates new examples. Unlike MixStyle, TALLY generates  
367 examples of minority classes or domains, and uses prototypes to improve the model robustness, which  
368 is more suitable for long-tailed multi-domain learning. DDG uses an extra network to disentangle  
369 original examples and generate more. Finally, BODA is a concurrent work for long-tailed multi-  
370 domain learning with an explicit regularizer. Unlike BODA, TALLY studies a conceptually different  
371 direction to cancel out domain-associated nuisances by domain-class balanced augmentation, leading  
372 to stronger empirical performance.

## 373 **C Additional Results of Synthetic Data**

### 374 **C.1 Detailed Dataset Description**

375 **VLCS-LT** contains examples from 4 different domains, including Caltech101, LabelMe, SUN09,  
376 VOC2007. To create the long-tailed class distribution, we modify the original dataset by removing  
377 training examples. The dataset contains 5 classes with 6,361 images of dimension (224,224,3). The  
378 long-tailed training distribution is visualized in Figure 3a. In subpopulation shift, the number of  
379 examples of each class per domain for validation and testing is 5, 10, respectively.

380  
381 **PACS-LT** includes 3,097 images collected from 4 domains (Art painting, Cartoon, Photo, Sketch)  
382 and 7 classes. Similar to VLCS-LT, we construct PACS-LT with long-tailed training distribution  
383 illustrated in Figure 3b. The validation set size and test set size of each class per domain in  
384 subpopulation shift are 15 and 30 respectively.

385  
386 **OfficeHome-LT** is built upon the original OfficeHome dataset, including 3280 images of 65 classes  
387 collected from four domains – Art, Clipart, Product, Real. The long-tailed training distribution is  
388 shown in Figure 3c and the number of examples for each class per domain in validation and test sets  
389 are 4, 8, respectively.

390  
391 **DomainNet-LT**. Similar to the other three datasets, DomainNet-LT covers 173,200 examples from  
392 Sketch, Infograph, Painting, Quickdraw, Real, Clipart. There are 345 classes in DomainNet-LT.  
393 In subpopulation shift, the number of examples of each class per domain is 3, 6, respectively. We  
394 illustrate the long-tailed training distribution in Figure 3d.

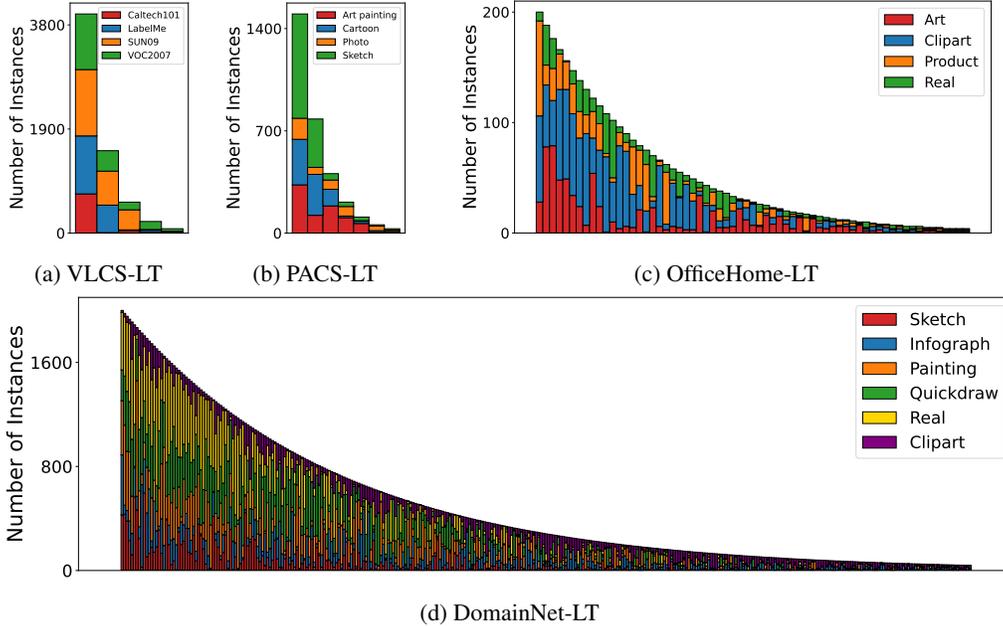


Figure 3: Long-tailed training distributions for all synthetic datasets. Here, the x-axis represents sorted class indices.

## 395 C.2 Detailed Hyperparameters

396 We evaluate performance under both subpopulation shift and domain shift. In subpopulation shift,  
 397 the test set is balanced across both domains and classes, which means that each domain-class pair  
 398 contains the same number of test examples. In domain shift, we use the classical domain generalization  
 399 setting [43]. More specifically, we alternately use one domain as the test domain, and the rest as the  
 400 training domains. Results are averaged over all combinations. We list the hyperparameters in Table 3  
 401 for the above four synthetic datasets.

Table 3: Hyperparameters for experiments on synthetic data.

Hyperparameters	VLCS-LT	PACS-LT	OfficeHome-LT	DomainNet-LT
Learning Rate	1e-5	1e-5	3e-5	3e-5
Weight Decay	1e-6	1e-6	1e-6	1e-6
Batch Size	18	18	18	18
Epochs	15	15	15	15
Steps	200	500	500	1000
Warm Start Epochs	7	7	7	7
$\gamma$ in feat. estimation	0.8	0.8	0.8	0.8
class prototype mixup parameter $\alpha_c$	0.2	0.5	0.5	0.5
domain prototype mixup parameter $\alpha_d$	0.2	0.5	0.5	0.5

## 402 C.3 Full Results

403 The full results of subpopulation shift are reported in Table 4. In domain shift, we report the results  
 404 of each domain for VLCS-LT, PACS-LT, OfficeHome-LT and DomainNet-LT in Table 5, 6, 7, 8,  
 405 respectively. In the domain shift scenario of VLCS-LT, though TALLY only performs best in  
 406 VOC2007 (VLCS), the results of TALLY is relatively more stable compared to other approaches,  
 407 leading to the best averaged performance.

Table 4: Full results of subpopulation shifts on long-tailed variants of domain generalization benchmarks. The standard deviation is computed across three seeds.

	VLCS-LT	PACS-LT	OfficeHome-LT	DomainNet-LT
ERM	73.33 ± 0.76%	90.40 ± 0.88%	61.07 ± 0.73%	44.33 ± 0.14%
Focal	74.83 ± 0.29%	90.44 ± 0.06%	62.57 ± 0.50%	47.35 ± 0.09%
LDAM	73.83 ± 1.04%	90.91 ± 0.15%	63.57 ± 0.08%	46.71 ± 0.33%
CRT	73.83 ± 1.89%	89.17 ± 1.47%	61.92 ± 0.54%	47.37 ± 0.83%
MiSLAS	71.83 ± 1.25%	90.99 ± 0.90%	61.38 ± 0.19%	49.15 ± 0.69%
Remix	74.16 ± 0.76%	90.83 ± 0.77%	61.59 ± 0.44%	47.56 ± 0.25%
Avg.				
IRM	50.50 ± 8.18%	65.24 ± 7.57%	45.48 ± 4.30%	35.57 ± 5.76%
GroupDRO	72.50 ± 0.50%	89.80 ± 0.70%	59.79 ± 0.43%	43.86 ± 0.33%
CORAL	71.67 ± 0.28%	88.22 ± 0.67%	59.10 ± 0.20%	43.92 ± 0.36%
LISA	74.67 ± 0.76%	90.08 ± 0.45%	57.39 ± 0.59%	43.17 ± 0.53%
MixStyle	74.30 ± 1.04%	91.55 ± 0.25%	62.26 ± 0.22%	43.59 ± 0.57%
DDG	73.00 ± 1.63%	89.60 ± 0.40%	58.80 ± 0.57%	44.46 ± 0.06%
BODA	74.83 ± 1.84%	91.03 ± 0.31%	62.79 ± 0.45%	47.61 ± 0.04%
<b>TALLY (ours)</b>	<b>76.83 ± 1.04%</b>	<b>92.38 ± 0.26%</b>	<b>67.00 ± 0.47%</b>	<b>50.15 ± 0.46%</b>
ERM	52.67 ± 2.31%	83.81 ± 2.43%	54.48 ± 0.89%	25.36 ± 0.63%
Focal	52.67 ± 1.15%	84.44 ± 0.81%	56.41 ± 1.34%	27.68 ± 0.13%
LDAM	51.33 ± 2.31%	85.24 ± 1.03%	58.07 ± 0.82%	27.23 ± 0.26%
CRT	52.00 ± 0.00%	83.02 ± 1.12%	55.51 ± 0.79%	27.55 ± 0.49%
MiSLAS	52.00 ± 3.46%	86.03 ± 0.90%	52.82 ± 0.39%	29.42 ± 0.15%
Remix	51.33 ± 3.05%	86.98 ± 0.59%	53.85 ± 0.54%	28.13 ± 0.99%
Worst				
IRM	32.63 ± 7.03%	59.38 ± 5.93%	40.58 ± 4.42%	20.48 ± 3.71%
GroupDRO	51.33 ± 1.15%	83.02 ± 0.59%	54.04 ± 0.30%	25.02 ± 0.73%
CORAL	49.33 ± 1.15%	81.59 ± 0.81%	53.53 ± 0.60%	24.50 ± 0.68%
LISA	53.33 ± 1.15%	83.01 ± 0.81%	49.04 ± 0.40%	24.05 ± 0.48%
MixStyle	54.00 ± 2.00%	86.98 ± 0.98%	55.19 ± 1.10%	22.65 ± 0.22%
DDG	51.33 ± 0.94%	82.70 ± 2.38%	51.99 ± 0.55%	24.35 ± 0.20%
BODA	54.00 ± 2.83%	85.08 ± 1.37%	55.70 ± 0.50%	26.94 ± 0.44%
<b>TALLY (ours)</b>	<b>56.00 ± 2.00%</b>	<b>89.21 ± 0.22%</b>	<b>60.45 ± 0.09%</b>	<b>29.55 ± 0.19%</b>

Table 5: Domain shift results on VLCS-LT.

	Caltech101	LabelMe	SUN09	VOC2007	Avg
ERM	92.39 ± 0.35%	47.74 ± 1.13%	59.79 ± 2.70%	70.55 ± 1.51%	67.62%
Focal	<b>97.12 ± 1.06%</b>	48.83 ± 0.38%	58.66 ± 2.31%	72.91 ± 1.51%	69.38%
LDAM	95.55 ± 1.65%	47.61 ± 1.12%	61.34 ± 3.02%	73.17 ± 1.33%	69.41%
CRT	92.39 ± 1.82%	47.74 ± 1.52%	55.10 ± 2.61%	67.45 ± 1.54%	65.67%
MiSLAS	95.24 ± 1.51%	47.00 ± 0.99%	56.03 ± 1.29%	76.28 ± 1.66%	68.64%
Remix	92.66 ± 1.42%	48.77 ± 1.31%	57.98 ± 2.62%	71.45 ± 0.87%	67.71%
IRM	74.10 ± 2.67%	37.07 ± 1.87%	34.33 ± 1.77%	47.78 ± 3.59%	48.32%
GroupDRO	93.79 ± 1.01%	49.63 ± 1.09%	<b>62.25 ± 1.89%</b>	71.06 ± 0.55%	69.18%
CORAL	93.94 ± 1.53%	48.29 ± 1.08%	56.12 ± 1.84%	67.82 ± 1.43%	66.54%
LISA	90.28 ± 0.68%	48.51 ± 1.58%	58.82 ± 2.41%	68.09 ± 1.53%	66.42%
MixStyle	96.58 ± 0.84%	48.15 ± 1.20%	58.82 ± 1.94%	68.09 ± 1.98%	67.75%
DDG	95.46 ± 1.19%	50.42 ± 1.45%	57.44 ± 2.07%	70.21 ± 1.33%	68.38%
BODA	95.60 ± 1.37%	<b>51.42 ± 1.31%</b>	59.93 ± 1.97%	71.57 ± 1.18%	69.63%
<b>TALLY (ours)</b>	95.22 ± 0.92%	50.07 ± 1.17%	60.13 ± 2.17%	<b>76.98 ± 0.57%</b>	<b>70.60%</b>

Table 6: Domain shift results on PACS-LT.

	Art painting	Cartoon	Photo	Sketch	Avg
ERM	80.41 ± 1.21%	70.21 ± 1.14%	94.46 ± 0.19%	60.00 ± 5.04%	76.27%
Focal	80.92 ± 0.51%	69.58 ± 0.64%	93.81 ± 0.80%	56.83 ± 2.04%	75.29%
LDAM	81.82 ± 1.14%	71.64 ± 0.66%	95.34 ± 0.32%	61.30 ± 4.83%	77.53%
CRT	78.14 ± 0.99%	67.17 ± 0.73%	94.33 ± 0.78%	55.62 ± 6.57%	73.82%
MiSLAS	81.31 ± 0.49%	71.15 ± 0.28%	93.51 ± 1.40%	65.78 ± 2.13%	77.94%
Remix	82.79 ± 1.21%	69.10 ± 1.13%	92.09 ± 1.01%	57.00 ± 4.13%	75.25%
IRM	51.87 ± 4.93%	50.27 ± 5.77%	69.11 ± 4.43%	39.13 ± 9.65%	52.60%
GroupDRO	80.20 ± 0.57%	70.61 ± 1.41%	94.58 ± 0.90%	61.61 ± 1.48%	76.75%
CORAL	77.60 ± 0.79%	68.19 ± 0.73%	93.88 ± 0.40%	62.82 ± 2.67%	75.62%
LISA	81.09 ± 0.61%	65.68 ± 0.87%	94.40 ± 0.33%	56.69 ± 1.99%	74.47%
MixStyle	83.45 ± 0.90%	72.84 ± 0.59%	95.20 ± 0.49%	67.61 ± 0.83%	79.78%
DDG	79.67 ± 0.77%	68.30 ± 0.34%	94.72 ± 0.50%	61.20 ± 0.82%	75.97%
BODA	81.13 ± 0.59%	72.03 ± 0.65%	95.73 ± 0.56%	66.34 ± 1.55%	78.81%
<b>TALLY (ours)</b>	<b>85.86 ± 0.40%</b>	<b>74.20 ± 0.30%</b>	<b>96.56 ± 0.20%</b>	<b>69.58 ± 0.62%</b>	<b>81.55%</b>

Table 7: Domain shift results on OfficeHome-LT.

	Art	Clipart	Product	Real	Avg
ERM	45.20 ± 0.73%	41.94 ± 0.17%	59.21 ± 0.44%	61.44 ± 0.27%	51.95%
Focal	47.06 ± 0.24%	43.29 ± 0.71%	62.34 ± 0.16%	63.45 ± 0.19%	54.03%
LDAM	47.08 ± 0.37%	42.89 ± 0.18%	61.48 ± 0.55%	62.93 ± 0.24%	53.60%
CRT	47.17 ± 0.26%	42.62 ± 0.55%	61.37 ± 0.12%	63.31 ± 0.25%	53.62%
MiSLAS	45.22 ± 0.52%	41.36 ± 0.09%	62.28 ± 0.49%	62.56 ± 0.25%	52.86%
Remix	44.26 ± 0.49%	39.18 ± 0.34%	60.70 ± 0.28%	61.58 ± 0.42%	51.43%
IRM	33.55 ± 4.21%	34.34 ± 3.74%	49.54 ± 5.30%	51.95 ± 4.64%	42.34%
GroupDRO	44.62 ± 0.51%	41.84 ± 0.68%	58.40 ± 0.43%	59.63 ± 0.53%	51.12%
CORAL	43.93 ± 0.56%	42.71 ± 0.59%	56.91 ± 0.45%	59.40 ± 0.94%	50.74%
LISA	41.80 ± 0.36%	36.96 ± 0.45%	56.51 ± 0.16%	57.62 ± 0.39%	48.22%
MixStyle	45.11 ± 0.18%	<b>45.52 ± 0.20%</b>	58.32 ± 0.64%	60.92 ± 0.22%	52.47%
DDG	43.89 ± 0.39%	42.79 ± 0.91%	57.92 ± 0.15%	59.69 ± 0.30%	51.07%
BODA	47.08 ± 0.25%	44.38 ± 0.77%	59.58 ± 0.26%	62.25 ± 0.10%	53.32%
<b>TALLY (ours)</b>	<b>49.79 ± 0.76%</b>	44.22 ± 0.45%	<b>63.02 ± 0.52%</b>	<b>65.71 ± 0.26%</b>	<b>55.69%</b>

Table 8: Domain shift results on DomainNet-LT.

	Sketch	Infograph	Painting	Quickdraw	Real	Clipart	Avg
ERM	39.22 ± 0.24%	18.96 ± 0.37%	34.71 ± 0.49%	10.70 ± 0.06%	50.87 ± 0.43%	44.83 ± 0.25%	33.21%
Focal	41.01 ± 0.61%	19.99 ± 0.14%	36.42 ± 0.45%	10.29 ± 0.23%	<b>55.63 ± 0.56%</b>	48.09 ± 0.72%	35.23%
LDAM	40.44 ± 0.26%	19.06 ± 0.20%	36.32 ± 0.52%	11.38 ± 0.29%	52.91 ± 0.52%	46.38 ± 0.44%	34.42%
CRT	40.78 ± 0.35%	20.41 ± 0.42%	39.01 ± 0.35%	<b>11.41 ± 0.17%</b>	55.26 ± 0.69%	50.00 ± 0.83%	36.14%
MiSLAS	41.34 ± 0.39%	19.89 ± 0.25%	39.85 ± 0.56%	11.00 ± 0.13%	55.50 ± 0.36%	49.49 ± 0.29%	36.18%
Remix	40.01 ± 0.51%	19.17 ± 0.46%	38.93 ± 0.32%	11.39 ± 0.20%	53.43 ± 0.60%	47.94 ± 0.71%	35.14%
IRM	34.65 ± 0.75%	15.41 ± 0.83%	28.18 ± 1.26%	7.69 ± 0.40%	40.83 ± 0.72%	42.36 ± 1.25%	28.19%
GroupDRO	38.47 ± 0.27%	18.63 ± 0.07%	34.23 ± 0.15%	10.26 ± 0.35%	50.80 ± 0.47%	42.85 ± 0.63%	32.54%
CORAL	39.42 ± 0.42%	19.30 ± 0.33%	35.15 ± 0.70%	10.61 ± 0.22%	51.05 ± 0.28%	45.15 ± 0.38%	33.44%
LISA	40.75 ± 0.46%	18.47 ± 0.14%	37.99 ± 0.19%	9.98 ± 0.09%	54.33 ± 0.49%	48.42 ± 0.42%	34.99%
MixStyle	40.99 ± 0.60%	18.64 ± 0.32%	35.86 ± 0.39%	11.03 ± 0.15%	50.26 ± 0.53%	45.49 ± 0.84%	33.71%
DDG	40.66 ± 0.58%	19.08 ± 0.34%	35.61 ± 0.63%	11.39 ± 0.29%	50.93 ± 0.41%	45.95 ± 0.47%	33.94%
BODA	41.95 ± 0.45%	<b>20.65 ± 0.58%</b>	37.98 ± 0.27%	11.02 ± 0.23%	55.22 ± 0.65%	48.26 ± 0.33%	35.85%
<b>TALLY (ours)</b>	<b>42.66 ± 0.32%</b>	19.26 ± 0.09%	<b>40.49 ± 0.34%</b>	11.15 ± 0.21%	54.79 ± 0.62%	<b>50.36 ± 0.41%</b>	<b>36.45%</b>

408 **D Additional Results of Real-world Data**

409 **D.1 Detailed Dataset Description**

410 **TerraInc.** Building upon the original Terra Incognita [4], we select images from 10 classes and split  
 411 the entire dataset to training, validation and test domains, which includes images from 38,042, 6,783,  
 412 7,303 camera traps, respectively.

413 **iWildCam** is a wildlife recognition datasets. It is a multi-class species classification, where the  
 414 training data are collected from 243 domains and the test data includes images from 164 domains.  
 415 We follow Koh et al. [19] to split the data and construct training, validation and test sets.

416 **D.2 Detailed Hyperparameters**

417 We list the hyperparameters in Table 9 for both TerraInc and iWildCam datasets.

Table 9: Hyperparameters for experiments on real-world data.

Hyperparameters	TerraInc	iWildCam
Learning Rate	3e-5	3e-5
Weight Decay	1e-6	0
Batch Size	18	16
Epochs	15	15
Steps	1000	1000
Warm Start Epochs	7	7
$\gamma$ in feat. estimation	0.8	0.8
class prototype mixup parameter $\alpha_c$	0.5	0.5
domain prototype mixup parameter $\alpha_d$	0.5	0.5

418 **D.3 Full Results**

419 The full results on Real-world Data are reported in Table 10.

Table 10: Full Results of Domain Shifts on Real-world Data.

	TerraInc		iWildCam	
	Macro F1	Acc	Macro F1	Acc
ERM	42.35 ± 1.25%	54.81 ± 0.83%	32.0 ± 1.5%	69.0 ± 0.4%
Focal	43.54 ± 0.81%	56.62 ± 1.49%	33.2 ± 1.2%	74.7 ± 1.9%
LDAM	44.29 ± 1.41%	57.22 ± 0.92%	32.7 ± 0.9%	<b>75.2 ± 2.0%</b>
CRT	43.09 ± 0.79%	58.27 ± 1.35%	32.5 ± 1.8%	67.3 ± 1.3%
MiSLAS	40.68 ± 1.33%	52.96 ± 2.58%	30.5 ± 1.1%	59.8 ± 2.8%
Remix	43.72 ± 1.87%	58.40 ± 2.57%	28.4 ± 0.8%	65.8 ± 1.6%
IRM	31.17 ± 3.52%	49.27 ± 5.22%	15.1 ± 4.9%	59.8 ± 3.7%
GroupDRO	42.22 ± 0.87%	56.43 ± 1.63%	23.9 ± 2.1%	72.7 ± 2.0%
CORAL	45.43 ± 0.92%	58.10 ± 1.38%	32.8 ± 0.1%	73.3 ± 4.3%
LISA	39.27 ± 0.69%	54.92 ± 1.04%	27.6 ± 1.2%	64.9 ± 2.2%
MixStyle	44.73 ± 0.99%	57.55 ± 2.05%	32.4 ± 1.1%	74.9 ± 2.7%
DDG	40.47 ± 1.93%	53.61 ± 1.71%	29.8 ± 0.2%	69.7 ± 2.3%
BODA	44.47 ± 0.84%	57.52 ± 1.13%	32.9 ± 0.3%	70.5 ± 2.3%
<b>TALLY (ours)</b>	<b>46.23 ± 0.56%</b>	<b>59.89 ± 1.32%</b>	<b>34.4 ± 0.4%</b>	73.4 ± 1.8%

## 420 E Analysis of Performance

### 421 E.1 Can we simply combine invariant learning approaches with long-tailed learning 422 techniques?

423 To further understand the performance gains of TALLY, we investigate whether combining existing  
424 invariant learning and long-tailed learning approaches can tackle multi-domain long-tailed distribution  
425 shifts. Specifically, we incorporate four up-weighting or up-sampling approaches (UW, Focal, LDAM,  
426 CRT) with two representative invariant learning methods (CORAL, MixStyle). We report the relative  
427 improvement of each combination over the vanilla methods in Figure 4. Here, we use Officehome-LT  
428 and DomainNet-LT to evaluate subpopulation shift and TerraInc and iWildCam to evaluate domain  
429 shift performance. We see that applying loss up-weighting or up-sampling approaches on performant  
430 invariant learning approaches does improve their performance, as evidenced by Figure 4. Nonetheless,  
431 the consistent improvements from TALLY indicates the importance of considering domain-class pair  
432 information to achieve balanced augmentation.

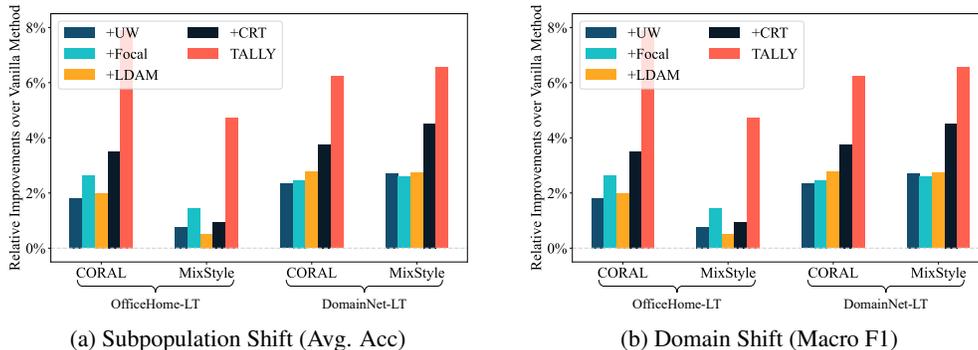


Figure 4: Comparison between TALLY and variants of two domain generalization approaches (CORAL, MixStyle), where we replace the losses of them with class re-weighting or re-sampling ones.

### 433 E.2 How do prototypes benefit invariant learning?

434 We analyze the effects of prototypes in alleviating domain-associated nuisances. Specifically, we  
435 compare TALLY with three variants: (1) without using any prototype information (None); (2)  
436 only applying class prototype (C Only); (3) only applying class-agnostic nuisances (D Only). We  
437 report the results in Figure 5. We observe that adding class prototype does improve the perfor-  
438 mance, especially the worst-domain accuracy in Officehome-LT. The class-agnostic domain factors  
439 also benefits the performance to some extent. In summary, TALLY outperforms its vari-  
440 ants, demonstrating the effectiveness of prototype representation in mitigating domain-associated  
441 nuisances, especially the worst-domain accuracy in Officehome-LT. The class-agnostic domain factors  
442 also benefits the performance to some extent. In summary, TALLY outperforms its vari-  
443 ants, demonstrating the effectiveness of prototype representation in mitigating domain-associated  
444 nuisances. In summary, TALLY outperforms its variants, demonstrating the effectiveness of prototype  
445 representation in mitigating domain-associated nuisances. In summary, TALLY outperforms its vari-  
446 ants, demonstrating the effectiveness of prototype representation in mitigating domain-associated  
447 nuisances.

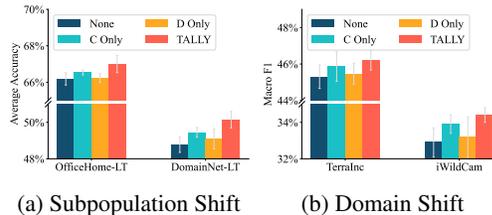


Figure 5: Analysis of prototype-guided invariant learning. C Only and D Only represent only using class prototype representation or class-agnostic domain factors, respectively.

### 448 E.3 Does TALLY lead to stronger domain invariance?

449 We analyze and compare the domain invariance of classifiers trained by ERM, TALLY, and other  
450 invariant learning approaches. Following [39, 40], we measure the lack of domain invariance as  
451 the accuracy of domain prediction ( $I_{acc}$ ) and as the pairwise divergence of unscaled logits ( $I_{kl}$ ).  
452 Specifically, for the accuracy of domain prediction, we perform logistic regression on top of the  
453 unscaled logits to predict the domain. For the pairwise divergence, we use kernel density estimation  
454 to estimate the probability density function  $P(h^{c,d})$  of logits from domain-class pair  $(c, d)$  and  
455 calculate the KL divergence of the distribution of logits from different pairs. Formally,  $I_{kl}$  is defined  
456 as  $I_{kl} = \frac{1}{|C||D|^2} \sum_{c \in C} \sum_{d', d \in D} \text{KL}(P(h^{c,d}) || P(h^{c,d'}))$ . We report the results of Officehome-LT and

Table 11: Invariance Analysis of TALLY. OH-LT and DN-LT represents Officehome-LT and DomainNet-LT, respectively.

Model	OH-LT		DN-LT	
	$I_{acc} \downarrow$	$I_{kl} \downarrow$	$I_{acc} \downarrow$	$I_{kl} \downarrow$
ERM	46.35%	2.030	70.00%	4.852
MixStyle	44.42%	2.169	67.11%	5.661
CORAL	42.21%	1.248	66.79%	4.593
BODA	40.10%	2.052	65.15%	6.810
<b>TALLY</b>	<b>39.52%</b>	<b>1.179</b>	<b>63.80%</b>	<b>3.956</b>

457 DomainNet-LT in Table 11. Smaller  $I_{acc}$  and  $I_{kl}$  values indicate more invariant representations with  
 458 respect to the labels. The results show that TALLY does lead to greater domain-invariance compared to  
 459 to prior invariant learning approaches (e.g., BODA).

#### 460 E.4 Analysis of Sampling Strategies

461 Finally, we compare the proposed selective balanced sampling in TALLY with domain-class balanced  
 462 sampling. For an example pair  $(x_i, y_i, d_i)$  and  $(x_j, y_j, d_j)$ , selective balanced sampling gets  
 463  $y_i \sim \text{Uniform}(\mathcal{C})$  and  $d_j \sim \text{Uniform}(\mathcal{D})$ , while traditional balanced sampling get  $(y_i, d_i), (y_j, d_j) \sim$   
 464  $\text{Uniform}(\mathcal{C}, \mathcal{D})$ . The results of subpopulation shifts in OfficeHome-LT, DomainNet-LT and of domain  
 465 shifts (Macro-F1) in TerraInc, iWildCam are reported in Table 12, indicating the effectiveness of  
 466 selective balanced sampling in transferring knowledge over domains and classes.

Table 12: Comparison between sampling strategies.

	OfficeHome-LT		DomainNet-LT	
	Avg.	Worst	Avg.	Worst
Balanced Sampling	$65.03 \pm 0.91\%$	$58.33 \pm 0.24\%$	$49.35 \pm 0.21\%$	$27.97 \pm 0.25\%$
<b>TALLY(Selective)</b>	<b><math>67.00 \pm 0.47\%</math></b>	<b><math>60.45 \pm 0.09\%</math></b>	<b><math>50.15 \pm 0.46\%</math></b>	<b><math>29.55 \pm 0.19\%</math></b>
	TerraInc		iWildCam	
	Macro F1	Acc	Macro F1	Acc
Balanced Sampling	$44.79 \pm 0.62\%$	$57.78 \pm 0.36\%$	$33.1 \pm 0.4\%$	$71.6 \pm 1.8\%$
<b>TALLY(Selective)</b>	<b><math>46.23 \pm 0.56\%</math></b>	<b><math>59.89 \pm 1.32\%</math></b>	<b><math>34.4 \pm 0.4\%</math></b>	$73.4 \pm 1.8\%$

## 467 F Results on Standard Domain Generalization Benchmarks

468 In this section, we present the additional comparison on standard domain generalization benchmarks.  
 469 Notice that the data distributions in these standard benchmarks are not long-tailed, which is thus  
 470 not our focus in this paper. The goal is to compare our approach with other domain generalization  
 471 methods. In Table 13-16, we present results on four standard benchmarks: VLCS, PACS, OfficeHome,  
 472 DomainNet, respectively. Results for all algorithms except TALLY are directly copied from [12] and  
 473 [39]. In Table 17, we summarize all results and show the comparison between different approaches.  
 474 According to the results, TALLY can achieve comparable performance compared with state-of-the-art  
 475 domain generalization approaches.

Table 13: Comparison on the standard VLCS benchmark.

	Caltech101	LabelMe	SUN09	VOC2007	Avg
ERM	97.7 ± 0.4	64.3 ± 0.9	73.4 ± 0.5	74.6 ± 1.3	77.5
IRM	98.6 ± 0.1	64.9 ± 0.9	73.4 ± 0.6	77.3 ± 0.9	78.5
GroupDRO	97.3 ± 0.3	63.4 ± 0.9	69.5 ± 0.8	76.7 ± 0.7	76.7
Mixup	98.3 ± 0.6	64.8 ± 1.0	72.1 ± 0.5	74.3 ± 0.8	77.4
MLDG	97.4 ± 0.2	65.2 ± 0.7	71.0 ± 1.4	75.3 ± 1.0	77.2
CORAL	98.3 ± 0.1	66.1 ± 1.2	73.4 ± 0.3	77.5 ± 1.2	<b>78.8</b>
MMD	97.7 ± 0.1	64.0 ± 1.1	72.8 ± 0.2	75.3 ± 3.3	77.5
DANN	<b>99.0 ± 0.3</b>	65.1 ± 1.4	73.1 ± 0.3	77.2 ± 0.6	78.6
CDANN	97.1 ± 0.3	65.1 ± 1.2	70.7 ± 0.8	77.1 ± 1.5	77.5
MTL	97.8 ± 0.4	64.3 ± 0.3	71.5 ± 0.7	75.3 ± 1.7	77.2
SagNet	97.9 ± 0.4	64.5 ± 0.5	71.4 ± 1.3	77.5 ± 0.5	77.8
ARM	98.7 ± 0.2	63.6 ± 0.7	71.3 ± 1.2	76.7 ± 0.6	77.6
VREx	98.4 ± 0.3	64.4 ± 1.4	74.1 ± 0.4	76.2 ± 1.3	78.3
RSC	97.9 ± 0.1	62.5 ± 0.7	72.3 ± 1.2	75.6 ± 0.8	77.1
BODA	98.1 ± 0.3	64.5 ± 0.4	<b>74.3 ± 0.3</b>	78.0 ± 0.6	78.5
<b>TALLY (ours)</b>	97.5 ± 0.5	<b>67.2 ± 1.1</b>	73.8 ± 0.5	<b>79.2 ± 0.9</b>	<b>78.8</b>

Table 14: Comparison on the standard PACS benchmark.

	Art painting	Cartoon	Photo	Sketch	Avg
ERM	84.7 ± 0.4	80.8 ± 0.6	97.2 ± 0.3	79.3 ± 1.0	85.5
IRM	84.8 ± 1.3	76.4 ± 1.1	96.7 ± 0.6	76.1 ± 1.0	83.5
GroupDRO	83.5 ± 0.9	79.1 ± 0.6	96.7 ± 0.3	78.3 ± 2.0	84.4
Mixup	86.1 ± 0.5	78.9 ± 0.8	97.6 ± 0.1	75.8 ± 1.8	84.6
MLDG	85.5 ± 1.4	80.1 ± 1.7	97.4 ± 0.3	76.6 ± 1.1	84.9
CORAL	88.3 ± 0.2	80.0 ± 0.5	97.5 ± 0.3	78.8 ± 1.3	86.2
MMD	86.1 ± 1.4	79.4 ± 0.9	96.6 ± 0.2	76.5 ± 0.5	84.6
DANN	86.4 ± 0.8	77.4 ± 0.8	97.3 ± 0.4	73.5 ± 2.3	83.6
CDANN	84.6 ± 1.8	75.5 ± 0.9	96.8 ± 0.3	73.5 ± 0.6	82.6
MTL	87.5 ± 0.8	77.1 ± 0.5	96.4 ± 0.8	77.3 ± 1.8	84.6
SagNet	87.4 ± 1.0	80.7 ± 0.6	97.1 ± 0.1	80.0 ± 0.4	86.3
ARM	86.8 ± 0.6	76.8 ± 0.5	97.4 ± 0.3	79.3 ± 1.2	85.1
VREx	86.0 ± 1.6	79.1 ± 0.6	96.9 ± 0.5	77.7 ± 1.7	84.9
RSC	85.4 ± 0.8	79.7 ± 1.8	97.6 ± 0.3	78.2 ± 1.2	85.2
BODA	88.2 ± 0.2	<b>81.7 ± 0.3</b>	<b>97.8 ± 0.2</b>	80.2 ± 0.3	86.9
<b>TALLY (ours)</b>	<b>89.5 ± 0.8</b>	81.2 ± 0.7	97.0 ± 0.1	<b>81.7 ± 0.9</b>	<b>87.4</b>

Table 15: Comparison on the standard OfficeHome benchmark.

	Art	Clipart	Product	Real	Avg
ERM	61.3 ± 0.7	52.4 ± 0.3	75.8 ± 0.1	76.6 ± 0.3	66.5
IRM	58.9 ± 2.3	52.2 ± 1.6	72.1 ± 2.9	74.0 ± 2.5	64.3
GroupDRO	60.4 ± 0.7	52.7 ± 1.0	75.0 ± 0.7	76.0 ± 0.7	66.0
Mixup	62.4 ± 0.8	54.8 ± 0.6	76.9 ± 0.3	78.3 ± 0.2	68.1
MLDG	61.5 ± 0.9	53.2 ± 0.6	75.0 ± 1.2	77.5 ± 0.4	66.8
CORAL	65.3 ± 0.4	54.4 ± 0.5	76.5 ± 0.1	78.4 ± 0.5	68.7
MMD	60.4 ± 0.2	53.3 ± 0.3	74.3 ± 0.1	77.4 ± 0.6	66.3
DANN	59.9 ± 1.3	53.0 ± 0.3	73.6 ± 0.7	76.9 ± 0.5	65.9
CDANN	61.5 ± 1.4	50.4 ± 2.4	74.4 ± 0.9	76.6 ± 0.8	65.8
MTL	61.5 ± 0.7	52.4 ± 0.6	74.9 ± 0.4	76.8 ± 0.4	66.4
SagNet	63.4 ± 0.2	54.8 ± 0.4	75.8 ± 0.4	78.3 ± 0.3	68.1
ARM	58.9 ± 0.8	51.0 ± 0.5	74.1 ± 0.1	75.2 ± 0.3	64.8
VREx	60.7 ± 0.9	53.0 ± 0.9	75.3 ± 0.1	76.6 ± 0.5	66.4
RSC	60.7 ± 1.4	51.4 ± 0.3	74.8 ± 1.1	75.1 ± 1.3	65.5
BODA	<b>65.4 ± 0.1</b>	<b>55.4 ± 0.3</b>	77.1 ± 0.1	<b>79.5 ± 0.3</b>	<b>69.3</b>
<b>TALLY (ours)</b>	64.2 ± 0.5	55.1 ± 0.8	<b>78.0 ± 1.1</b>	79.2 ± 0.5	69.1

Table 16: Comparison on the standard DomainNet benchmark.

	Sketch	Infograph	Painting	Quickdraw	Real	Clipart	Avg
ERM	49.8 ± 0.4	18.8 ± 0.3	46.7 ± 0.3	12.2 ± 0.4	59.6 ± 0.1	58.1 ± 0.3	40.9
IRM	42.3 ± 3.1	15.0 ± 1.5	38.3 ± 4.3	10.9 ± 0.5	48.2 ± 5.2	48.5 ± 2.8	33.9
GroupDRO	40.1 ± 0.6	17.5 ± 0.4	33.8 ± 0.5	9.3 ± 0.3	51.6 ± 0.4	47.2 ± 0.5	33.3
Mixup	48.2 ± 0.5	18.5 ± 0.5	44.3 ± 0.5	12.5 ± 0.4	55.8 ± 0.3	55.7 ± 0.3	39.2
MLDG	50.2 ± 0.4	19.1 ± 0.3	45.8 ± 0.7	13.4 ± 0.3	59.6 ± 0.2	59.1 ± 0.2	41.2
CORAL	50.1 ± 0.6	19.7 ± 0.2	46.6 ± 0.3	13.4 ± 0.4	59.8 ± 0.2	59.2 ± 0.1	41.5
MMD	28.9 ± 11.9	11.0 ± 4.6	26.8 ± 11.3	8.7 ± 2.1	32.7 ± 13.8	32.1 ± 13.3	23.4
DANN	46.8 ± 0.6	18.3 ± 0.1	44.2 ± 0.7	11.8 ± 0.1	55.5 ± 0.4	53.1 ± 0.2	38.3
CDANN	45.9 ± 0.5	17.3 ± 0.1	43.7 ± 0.9	12.1 ± 0.7	56.2 ± 0.4	54.6 ± 0.4	38.3
MTL	49.2 ± 0.1	18.5 ± 0.4	46.0 ± 0.1	12.5 ± 0.1	59.5 ± 0.3	57.9 ± 0.5	40.6
SagNet	48.8 ± 0.2	19.0 ± 0.2	45.3 ± 0.3	12.7 ± 0.5	58.1 ± 0.5	57.7 ± 0.3	40.3
ARM	43.5 ± 0.4	16.3 ± 0.5	40.9 ± 1.1	9.4 ± 0.1	53.4 ± 0.4	49.7 ± 0.3	35.5
VREx	42.0 ± 3.0	16.0 ± 1.5	35.8 ± 4.6	10.9 ± 0.3	49.6 ± 4.9	47.3 ± 3.5	33.6
RSC	47.8 ± 0.9	18.3 ± 0.5	44.4 ± 0.6	12.2 ± 0.2	55.7 ± 0.7	55.0 ± 1.2	38.9
BODA	<b>51.3 ± 0.3</b>	<b>20.5 ± 0.7</b>	<b>48.0 ± 0.1</b>	<b>13.8 ± 0.6</b>	<b>60.6 ± 0.4</b>	<b>62.1 ± 0.4</b>	<b>42.7</b>
<b>TALLY (ours)</b>	<u>50.5 ± 0.2</u>	<u>19.7 ± 0.1</u>	<u>47.7 ± 0.6</u>	<b>14.1 ± 0.3</b>	<u>60.0 ± 0.2</u>	<u>60.1 ± 0.5</u>	<u>42.0</u>

Table 17: Domain shift results over all four benchmarks.

	VLCS	PACS	OfficeHome	DomainNet	Avg
ERM	77.5 ± 0.4	85.5 ± 0.2	66.5 ± 0.3	40.9 ± 0.1	67.6
IRM	78.5 ± 0.5	83.5 ± 0.8	64.3 ± 2.2	33.9 ± 2.8	65.1
GroupDRO	76.7 ± 0.6	84.4 ± 0.8	66.0 ± 0.7	33.3 ± 0.2	65.1
Mixup	77.4 ± 0.6	84.6 ± 0.6	68.1 ± 0.3	39.2 ± 0.1	67.3
MLDG	77.2 ± 0.4	84.9 ± 1.0	66.8 ± 0.6	41.2 ± 0.1	67.5
CORAL	<b>78.8 ± 0.6</b>	86.2 ± 0.3	68.7 ± 0.3	41.5 ± 0.1	68.8
MMD	77.5 ± 0.9	84.6 ± 0.5	66.3 ± 0.1	23.4 ± 9.5	63.0
DANN	78.6 ± 0.4	83.6 ± 0.4	65.9 ± 0.6	38.3 ± 0.1	66.6
CDANN	77.5 ± 0.1	82.6 ± 0.9	65.8 ± 1.3	38.3 ± 0.3	66.3
MTL	77.2 ± 0.4	84.6 ± 0.5	66.4 ± 0.5	40.6 ± 0.1	67.2
SagNet	77.8 ± 0.5	86.3 ± 0.2	68.1 ± 0.1	40.3 ± 0.1	68.1
ARM	77.6 ± 0.3	85.1 ± 0.4	64.8 ± 0.3	35.5 ± 0.2	65.8
VREx	78.3 ± 0.2	84.9 ± 0.6	66.4 ± 0.6	33.6 ± 2.9	65.8
RSC	77.1 ± 0.5	85.2 ± 0.9	65.5 ± 0.9	38.9 ± 0.5	66.7
BODA	78.5 ± 0.3	86.9 ± 0.4	<b>69.3 ± 0.1</b>	<b>42.7 ± 0.1</b>	<b>69.4</b>
<b>TALLY (ours)</b>	<b>78.8 ± 0.4</b>	<b>87.4 ± 0.2</b>	<u>69.1 ± 0.4</u>	<u>42.0 ± 0.1</u>	<u>69.3</u>