

# On Multi-Domain Long-Tailed Recognition, Imbalanced Domain Generalization and Beyond

Yuzhe Yang  
MIT

Hao Wang  
Rutgers University

Dina Katabi  
MIT

## Abstract

Real-world data often exhibit imbalanced label distributions. Existing studies on data imbalance focus on single-domain settings, i.e., samples are from the same data distribution. However, natural data can originate from distinct domains, where a minority class in one domain could have abundant instances from other domains. We formalize the task of Multi-Domain Long-Tailed Recognition (MDLT), which learns from multi-domain imbalanced data, addresses *label imbalance*, *domain shift*, and *divergent label distributions across domains*, and generalizes to all domain-class pairs. We first develop the *domain-class transferability graph*, and show that such transferability governs the success of learning in MDLT. We then propose **BoDA**, a theoretically grounded learning strategy that tracks the upper bound of transferability statistics, and ensures *balanced* alignment and calibration across imbalanced domain-class distributions. We curate five MDLT benchmarks based on widely-used multi-domain datasets, and compare **BoDA** to twenty algorithms that span different learning strategies. Extensive and rigorous experiments verify the superior performance of **BoDA**. Further, as a byproduct, **BoDA** establishes new state-of-the-art on Domain Generalization benchmarks, highlighting the importance of addressing data imbalance across domains, which can be crucial for improving generalization to unseen domains. Code and data are available at: <https://github.com/YyzHarry/multi-domain-imbalance>.

## 1 Introduction

Real-world data often exhibit label imbalance – i.e., instead of a uniform label distribution over classes, in reality, data are by their nature imbalanced: a few classes contain a large number of instances, whereas many others have only a few instances [5, 6, 52]. This phenomenon poses a challenge for deep recognition models, and has motivated several prior solutions [6, 10, 33, 39, 52, 53]. Such prior solutions focus on *single domain* scenarios, i.e., samples are from the same data distribution; they propose techniques for learning from imbalanced training data and generalizing to a balanced test set.

In contrast, this paper formulates the problem of *Multi-Domain Long-Tailed Recognition* (MDLT) as learning from multi-domain imbalanced data, with each domain having its own imbalanced label distribution, and generalizing to a test set that is balanced over all domain-class pairs. MDLT is a natural extension of the single domain case. It arises in real-world scenarios, where data targeted for one task can originate from different domains. For example, in visual recognition problems, minority classes from “photo” images could be complemented with potentially abundant samples from “sketch” images. Similarly, in autonomous driving, the minority accident class in “real” life could be enriched with accidents generated in “simulation”. Also, in medical diagnosis, data from distinct populations could enhance each other, where minority samples from one institution could be enriched with instances from others. In the above examples, different data types act as distinct

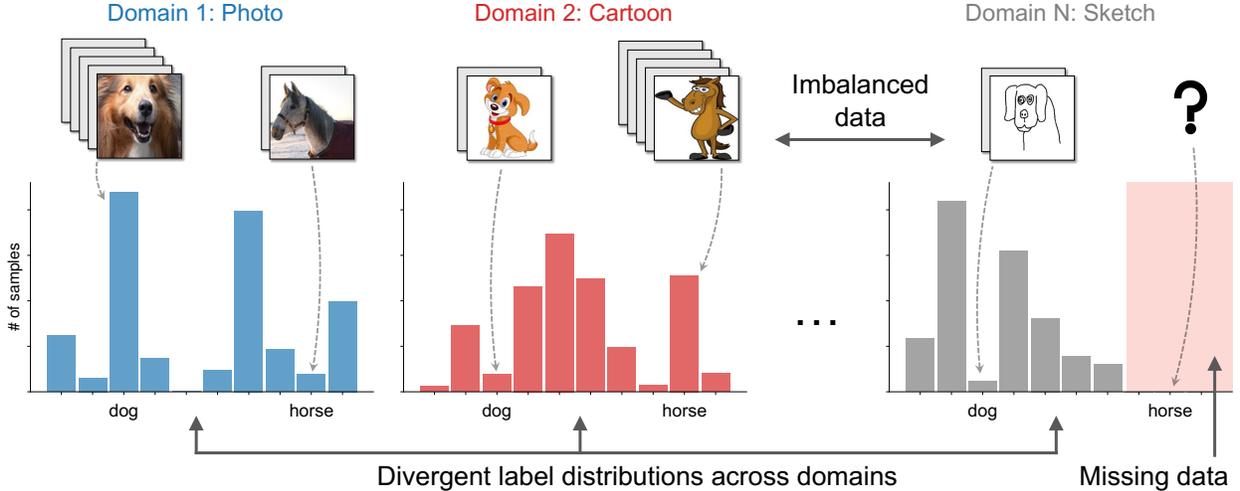


Figure 1: Multi-Domain Long-Tailed Recognition (MDLT) aims to learn from imbalanced data from multiple distinct domains, tackle label imbalance, domain shift, and divergent label distributions across domains, and generalize to the entire set of classes over all domains.

*domains*, and such multi-domain data could be leveraged to tackle the inherent data imbalance within each domain.

We note that MDLT has key differences from its single-domain counterpart:

- First, the label distribution for each domain is likely different from other domains. For example, in Fig. 1, both “Photo” and “Cartoon” domains exhibit imbalanced label distributions; Yet, the “horse” class in “Cartoon” has many more samples than in “Photo”. This creates challenges with *divergent label distributions across domains*, in addition to in-domain data imbalance.
- Second, multi-domain data inherently involves *domain shift*. Simply treating different domains as a whole and applying traditional data-imbalance methods is unlikely to yield the best results, as the domain gap can be arbitrarily large.
- Third, MDLT naturally motivates *zero-shot generalization within and across domains* – i.e., to generalize to both in-domain missing classes (Fig. 1 right part), as well as new domains with no training data, where the latter case is typically denoted as Domain Generalization (DG).

To deal with the above issues, we first develop the *domain-class transferability graph*, which quantifies the transferability between different domain-class pairs under data imbalance. In this graph, each node refers to a domain-class pair, and each edge refers to the distance between two domain-class pairs in the embedding space. We show that the transferability graph dictates the performance of imbalanced learning across domains. Inspired by this, we design BoDA (Balanced Domain-Class Distribution Alignment), a new loss function that encourages similarity between features of the same class in different domains, and penalizes similarity between features of different classes within and across domains. BoDA does so while accounting for that different classes have very different number of samples, and hence the statistics of their features are intrinsically imbalanced. Analytically, we prove that minimizing the BoDA loss optimizes an upper bound of the *balanced* transferability statistics, corroborating the effectiveness of BoDA for learning multi-domain imbalanced data.

For MDLT evaluation, we curate five MDLT benchmarks based on datasets widely used for domain generalization (DG). These datasets naturally exhibit heavy class imbalance within each domain and data shift across domains, highlighting that the MDLT problem is widely present in current benchmarks. We compare BoDA against twenty algorithms that span different learning strategies. Extensive experiments across benchmarks and algorithms verify that BoDA consistently outperforms all these baselines on all datasets.

Additionally, we examine how BoDA performs in the DG setting. We show that combining BoDA with the DG state-of-the-art (SOTA) consistently brings further gains, yielding a new SOTA for DG. These results shed light on how label imbalance can affect out-of-distribution generalization and highlight the importance of integrating label imbalance into practical DG algorithm design.

Our contributions are as follows:

- We formulate the MDLT problem as learning from multi-domain imbalanced data and generalizing across all domain-class pairs.
- We introduce the domain-class transferability graph, a unified model for investigating MDLT. We further show that the transferability statistics induced from such graph are crucial and govern the success of MDLT algorithms.
- We design BoDA, a simple, effective, and interpretable loss function for MDLT. We prove theoretically that minimizing the BoDA loss is equivalent to optimizing an upper bound of balanced transferability statistics.
- Extensive experiments on benchmark datasets verify the superior and consistent performance of BoDA. Further, combined with DG algorithms, BoDA establishes a new SOTA on DG benchmarks, highlighting the importance of tackling cross-domain data imbalance for domain generalization.

## 2 Related Work

**Long-Tailed Recognition.** The literature is rich with research on long-tailed recognition [33, 57]. Proposed solutions include re-balancing the data by either over-sampling the minority classes or under-sampling the majority classes [9, 20], re-weighting or adjusting the loss functions [6, 10, 12, 22], as well as leveraging relevant learning paradigms such as transfer learning [33], metric learning [55], meta-learning [43], two-stage training [23], ensemble learning [48, 56], and self-supervised learning [30, 52]. Recent studies have also explored imbalanced regression [53]. In contrast to these past works, we extend long-tailed recognition to the multi-domain setting, and introduce new techniques suitable for learning from multi-domain imbalanced data.

**Multi-Domain Learning.** Multi-domain learning (MDL) aims to learn a model of minimal risk from datasets drawn from different underlying distributions [13], and is a specific case of transfer learning [37]. In contrast to domain adaptation (DA) [3, 37], which aims to minimize the risk over a single “target” domain, MDL minimizes the risk over all “source” domains, and considers both average and worst risks over all distributions [41]. Past solutions for MDL include designing shared and domain-specific models [13, 49], leveraging multi-task learning [51], and learning domain-invariant features [15, 31, 41, 45]. Our work falls under the MDL framework, but considers the practical and realistic setting where the label distribution is imbalanced within each domain and across domains.

**Domain Generalization.** Unlike MDL which focuses on in-domain generalization, domain gener-

alization (DG) aims to learn from multiple training domains and generalize to unseen domains [59]. Previous approaches include learning domain-invariant features [15, 31, 34], learning transferable model parameters using meta-learning [28, 54], data augmentation [7, 60], and capturing causal relationships [1, 25]. Past work on DG has not investigated label imbalance within a domain and across domains. This paper shows that label imbalance plays a crucial role in DG, and that by combating data imbalance, we substantially boost DG performance on standard benchmarks.

### 3 Domain-Class Transferability Graph

When learning from MDLT, a natural question arises:

*How do we model MDLT in the presence of both **domain shift** and **class imbalance** within and across domains?*

We argue that in contrast to single-domain imbalanced learning where the basic unit one cares about is a *class* (i.e., minority *vs.* majority classes), in MDLT, the basic unit naturally translates to a **domain-class pair**.

**Problem Setup.** Given a multi-domain classification task with a discrete label space  $\mathcal{C} = \{1, \dots, C\}$  and a domain space  $\mathcal{D} = \{1, \dots, D\}$ , let  $\mathcal{S} = \{(\mathbf{x}_i, c_i, d_i)\}_{i=1}^N$  be the training set, where  $\mathbf{x}_i \in \mathbb{R}^l$  denotes the input,  $c_i \in \mathcal{C}$  is the class label, and  $d_i \in \mathcal{D}$  is the domain label. We denote as  $\mathbf{z} = f(\mathbf{x}; \theta)$  the representation of  $\mathbf{x}$ , where  $f : \mathcal{X} \rightarrow \mathcal{Z}$  maps the input into a representation space  $\mathcal{Z} \subseteq \mathbb{R}^h$ . The final prediction  $\hat{c} = g(\mathbf{z})$  is given by a classification function  $g : \mathcal{Z} \rightarrow \mathcal{C}$ . We denote the set of samples belonging to domain  $d$  and class  $c$  (i.e., the domain-class pair  $(d, c)$ ) as  $\mathcal{S}_{d,c} \subseteq \mathcal{S}$ , with  $N_{d,c} \triangleq |\mathcal{S}_{d,c}|$  as the number of samples. Similarly,  $\mathcal{Z}_{d,c} \subseteq \mathcal{Z}$  denotes the representation set for  $(d, c)$ . We use  $\mathcal{M} = \mathcal{D} \times \mathcal{C} := \{(d, c) : d \in \mathcal{D}, c \in \mathcal{C}\}$  to denote the set of all domain-class pairs.

**Definition 1** (Transferability). *Given a learned model and a distance function  $\mathbf{d} : \mathbb{R}^h \times \mathbb{R}^h \rightarrow \mathbb{R}$  in the feature space, the transferability from domain-class pair  $(d, c)$  to  $(d', c')$  is:*

$$\text{trans}((d, c), (d', c')) \triangleq \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c'})],$$

where  $\boldsymbol{\mu}_{d',c'} \triangleq \mathbb{E}_{\mathbf{z}' \in \mathcal{Z}_{d',c'}}[\mathbf{z}']$  is the first order statistics (i.e., mean) of  $(d', c')$ .

Intuitively, the transferability between two domain-class pairs is the average distance between their learned representations, characterizing how close they are in the feature space. By default,  $\mathbf{d}$  is chosen as the Euclidean distance, but it can also represent the higher order statistics of  $(d, c)$ . For example, the Mahalanobis distance [11] uses the covariance  $\boldsymbol{\Sigma}_{d,c} \triangleq \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [(\mathbf{z} - \boldsymbol{\mu}_{d,c})(\mathbf{z} - \boldsymbol{\mu}_{d,c})^\top]$ . In the remainder of the paper, with a slight abuse of the notation, we allow  $\boldsymbol{\mu}_{d,c}$  to represent both the first and higher order statistics for  $(d, c)$ .

**Definition 2** (Transferability Graph). *The transferability graph for a learned model is defined as  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where the vertices,  $\mathcal{V} \subseteq \{\boldsymbol{\mu}_{d,c}\}$ , represents the domain-class pairs, and the edges,  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ , are assigned weights equal to  $\text{trans}((d, c), (d', c'))$ .*

**Transferability Graph Visualization.** It is convenient to visualize the transferability graph of a learned model in a 2D Cartesian space. To do so, we use the average of  $\text{trans}((d, c), (d', c'))$  and

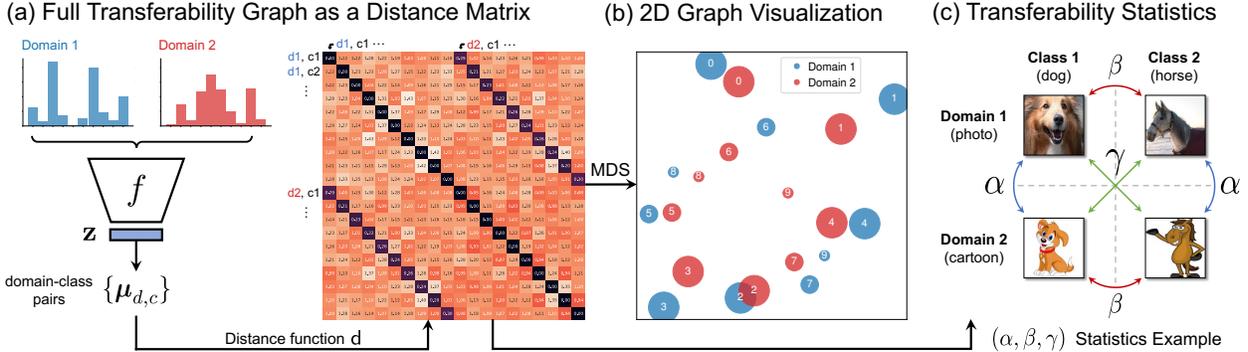


Figure 2: Overall framework of transferability graph. (a) Distribution statistics  $\{\mu_{d,c}\}$  is computed for all domain-class pairs, by which we generate a full transferability matrix. (b) MDS is used to project the graph into a 2D space for visualization. (c) We define  $(\alpha, \beta, \gamma)$  transferability statistics to further describe the whole transferability graph.

$\text{trans}((d', c'), (d, c))$  as a similarity measure between them. We can then visualize this similarity and the underlying transferability graph using multidimensional scaling (MDS) [8]. Figs. 2a and 2b show this process, where for each  $(d, c)$  pair, we estimate its distribution statistics  $\{\mu_{d,c}\}$  from the learned model and compute the transferability graph as a distance matrix. We then use MDS to project it into a 2D space, where each dot refers to one  $(d, c)$ , and the distance represents transferability.

**Definition 3** ( $(\alpha, \beta, \gamma)$  Transferability Statistics). *The transferability graph can be summarized by the following transferability statistics:*

$$\begin{aligned}
 \text{Different domains, same class: } & \alpha = \mathbb{E}_c \mathbb{E}_d \mathbb{E}_{d' \neq d} [\text{trans}((d, c), (d', c))] . \\
 \text{Same domain, different classes: } & \beta = \mathbb{E}_d \mathbb{E}_c \mathbb{E}_{c' \neq c} [\text{trans}((d, c), (d, c'))] . \\
 \text{Different domains, different classes: } & \gamma = \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_c \mathbb{E}_{c' \neq c} [\text{trans}((d, c), (d', c'))] .
 \end{aligned}$$

As illustrated in Fig. 2c,  $(\alpha, \beta, \gamma)$  captures the similarity between features of the same class across domains and different classes within and across domains.

## 4 What Makes for Good Representations in MDLT?

### 4.1 Divergent Label Distributions Hamper Transferable Features

MDLT has to deal with differences between the label distributions across domains. To understand the implications of this issue we start with an example.

**Motivating Example.** We construct **Digits-MLT**, a two-domain toy MDLT dataset that combines two digit datasets: MNIST-M [15] and SVHN [36]. The task is 10-class digit classification. Details of the datasets are in Appendix D. We manually vary the number of samples for each domain-class pair to simulate different label distributions, and train a plain ResNet-18 [21] using empirical risk minimization (ERM) for each case. We keep all test sets balanced and identical.

The results in Fig. 3 reveal interesting observations. When the per-domain label distributions are balanced and *identical* across domains, although a domain gap exists, it does not prohibit the model

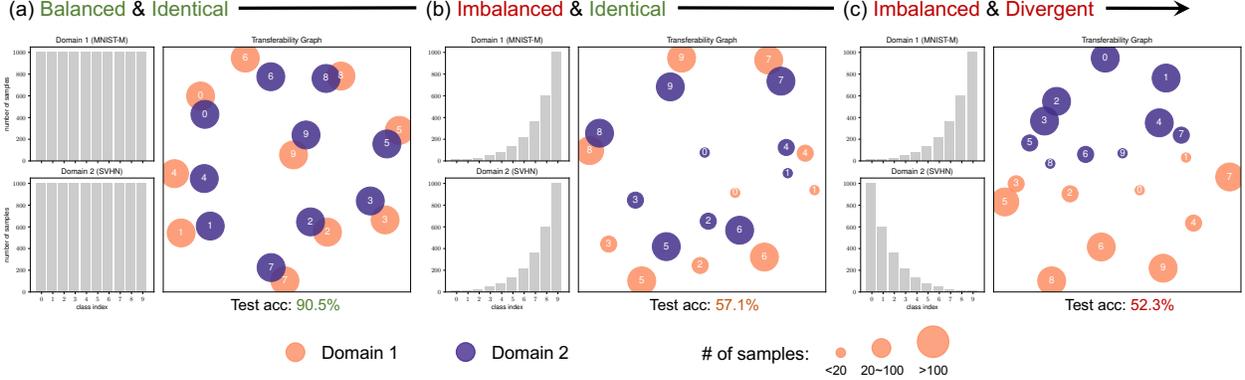


Figure 3: The evolving pattern of transferability graph when varying label proportions of Digits-MLT. (a) Label distributions for two domains are balanced and identical. (b) Label distributions for two domains are imbalanced but identical. (c) Label distributions for two domains are imbalanced and *divergent*.

from learning discriminative features of high accuracy (90.5%), as shown in Fig. 3a. If the label distributions are imbalanced but *identical*, as in Fig. 3b, ERM is still able to align similar classes in the two domains, where majority classes (e.g., class 9) are closer in terms of transferability than minority classes (e.g., class 0). In contrast, when the labels are both imbalanced and *mismatched* across domains, as in Fig. 3c, the learned features are no longer transferable, resulting in a clear gap across domains and the worst accuracy. This is because *divergent label distributions* across domains produce an undesirable shortcut; the model can minimize the classification loss simply by separating the two domains.

**Transferable Features are Desirable.** As the results indicate, *transferable* features across  $(d, c)$  pairs are needed, especially when imbalance occurs. In particular, the transferability link between the same class across domains should be greater than that between different classes within or across domains. This can be captured via the  $(\alpha, \beta, \gamma)$  transferability statistics, as we show next.

## 4.2 Transferability Statistics Characterize Generalization

**Motivating Example.** Again, we use Digits-MLT with varying label distributions. We consider three imbalance types to compose different label configurations: (1) **Uniform** (i.e., balanced labels), (2) **Forward-LT**, where the labels exhibit a long tail over class ids, and (3) **Backward-LT**, where labels are inversely long-tailed with respect to the class ids. For each configuration, we train 20 ERM models with varying hyperparameters. We then calculate the  $(\alpha, \beta, \gamma)$  statistics for each model, and plot its classification accuracy against  $(\beta + \gamma) - \alpha$ .

Fig. 4 reveals the following findings: (1) *The  $(\alpha, \beta, \gamma)$  statistics characterize a model’s performance in MDLT.* In particular, the  $(\beta + \gamma) - \alpha$  quantity displays a very strong correlation with test performance across the entire range and every label configuration. (2) *Data imbalance increases the risk of learning less transferable features.* When the label distributions are similar across domains (Fig. 4a), the models are robust to varying parameters, clustering in the upper-right region. However, as the labels become imbalanced (Figs. 4b, 4c) and further divergent (Figs. 4d, 4e), chances that the model learns non-transferable features (i.e., lower  $(\beta + \gamma) - \alpha$ ) increase, leading to a large

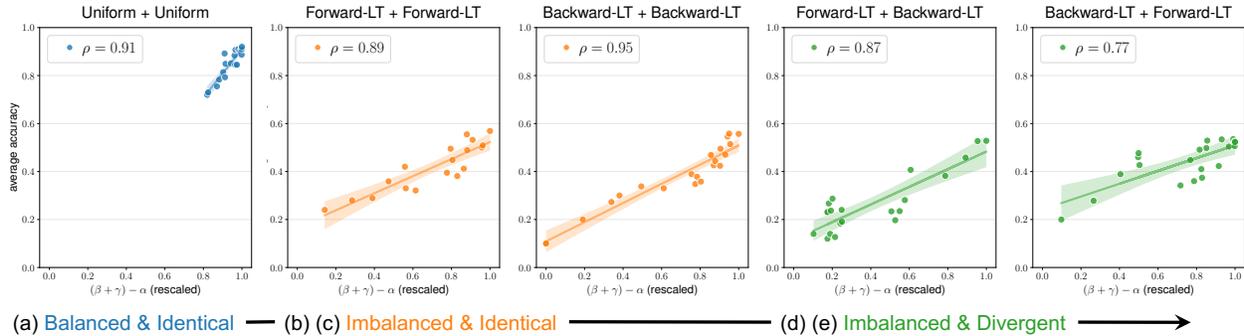


Figure 4: Correspondence between  $(\beta + \gamma) - \alpha$  quantity and test accuracy across different label configurations of `Digits-MLT`. Each plot refers to specific label distributions for two domains (e.g., (a) employs “Uniform” for domain 1 and “Uniform” for domain 2). Each point corresponds to a model trained with ERM using different hyperparameters.

drop in performance. We provide further evidence in Appendix H.4 showing that these observations hold regardless of datasets and training regimes.

### 4.3 A Loss that Bounds the Transferability Statistics

We use the above findings to design a new loss function particularly suitable for MDLT. We will first introduce the loss function then prove that it minimizes an upper bound of the  $(\alpha, \beta, \gamma)$  statistics. We start from a simple loss inspired by the metric learning objective [17, 44]. We call this loss  $\mathcal{L}_{\text{DA}}$  since it aims for Domain-Class Distribution Alignment, i.e., aligning the features of the same class across domains. Let  $(\mathbf{x}_i, c_i, d_i)$  denote a sample with feature  $\mathbf{z}_i$ . Given a set of training samples with feature set  $\mathcal{Z}$ , we have

$$\mathcal{L}_{\text{DA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) = \sum_{\mathbf{z}_i \in \mathcal{Z}} \frac{-1}{|\mathcal{D}| - 1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp(-\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i}))}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp(-\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))}. \quad (1)$$

Intuitively,  $\mathcal{L}_{\text{DA}}$  tackles label *divergence*, as  $(d, c)$  pairs that share same class would be pulled closer, and vice versa. It is also related to  $(\alpha, \beta, \gamma)$  statistics, as the numerator represents *positive* cross-domain pairs  $(\alpha)$ , and the denominator represents *negative* cross-class pairs  $(\beta, \gamma)$ . A detailed probabilistic interpretation of  $\mathcal{L}_{\text{DA}}$  is provided in Appendix B.2.

But,  $\mathcal{L}_{\text{DA}}$  does not address label *imbalance*. Note that  $(\alpha, \beta, \gamma)$  is defined in a *balanced* way, independent of the number of samples of each  $(d, c)$ . However, given an imbalanced dataset, most samples will come from majority domain-class pairs, which would dominate  $\mathcal{L}_{\text{DA}}$  and cause minority pairs to be overlooked.

**Balanced Domain-Class Distribution Alignment (BoDA).** To tackle data imbalance across  $(d, c)$  pairs, we modify the loss in Eqn. (1) to the BoDA loss:

$$\mathcal{L}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) = \sum_{\mathbf{z}_i \in \mathcal{Z}} \frac{-1}{|\mathcal{D}| - 1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i}))}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))}, \quad \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c}) = \frac{\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c})}{N_{d_i, c_i}}. \quad (2)$$

BoDA scales the original  $\mathbf{d}$  by a factor of  $1/N_{d_i, c_i}$ , i.e., it counters the effect of imbalanced domain-class pairs by introducing a *balanced* distance measure  $\tilde{\mathbf{d}}$ .

**Theorem 1** ( $\mathcal{L}_{\text{BoDA}}$  as an Upper Bound). *Given a multi-domain long-tailed dataset  $\mathcal{S}$  with domain label space  $\mathcal{D}$  and class label space  $\mathcal{C}$  satisfying  $|\mathcal{D}| > 1$  and  $|\mathcal{C}| > 1$ , let  $\mathcal{Z}$  be the representation set of all training samples, and  $(\alpha, \beta, \gamma)$  be the transferability statistics for  $\mathcal{S}$  defined in Definition 3. It holds that*

$$\mathcal{L}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) \geq N \log \left( |\mathcal{D}| - 1 + |\mathcal{D}|(|\mathcal{C}| - 1) \exp \left( \frac{|\mathcal{C}||\mathcal{D}|}{N} \cdot \alpha - \frac{|\mathcal{C}|}{N} \cdot \beta - \frac{|\mathcal{C}|(|\mathcal{D}| - 1)}{N} \cdot \gamma \right) \right). \quad (3)$$

The proof of Theorem 1 is in Appendix A.2. Theorem 1 has the following interesting implications: (1)  $\mathcal{L}_{\text{BoDA}}$  upper-bounds  $(\alpha, \beta, \gamma)$  statistics in a desired form that naturally translates to better performance. By minimizing  $\mathcal{L}_{\text{BoDA}}$ , we ensure a low  $\alpha$  (attract same classes) and high  $\beta, \gamma$  (separate different classes), which are essential conditions for generalization in MDLT. (2) The constant factors correspond to how much each component contributes to the transferability graph. Zooming on the arguments of  $\exp(\cdot)$ , we observe that the objective is proportional to  $\alpha - (\frac{1}{|\mathcal{D}|}\beta + \frac{|\mathcal{D}|-1}{|\mathcal{D}|\gamma})$ . According to Definition 3, we note that  $\alpha$  summarizes data similarity for the same class, while  $(\frac{1}{|\mathcal{D}|}\beta + \frac{|\mathcal{D}|-1}{|\mathcal{D}|\gamma})$  summarizes data similarity across different classes, using the weighted average of  $\beta$  and  $\gamma$ , where their weights are proportional to the number of associated domains (i.e., 1 for  $\beta$ ,  $(|\mathcal{D}| - 1)$  for  $\gamma$ ).

#### 4.4 Calibration for Data Imbalance Leads to Better Transfer

BoDA works by encouraging feature transfer for similar classes across domains, i.e., if  $(d, c)$  and  $(d', c)$  refer to the same class in different domains, then we want to transfer their features to each other. But, minority domain-class pairs naturally have worse  $\boldsymbol{\mu}_{d,c}$  estimates due to data scarcity, and forcing other pairs to transfer to them hurts learning. Thus, when bringing two domain-class pairs closer in the embedding space, we want the minority  $(d, c)$  to transfer to majority ones, not the inverse. The following example further clarifies this point.

**Motivating Example.** We use Digits-MLT with divergent labels (Fig. 5). We focus on *feature discrepancy*, i.e., the distance between training and test features for the same class. For each class in domain 1, we compute the distance in the feature space between the means of the training set and test set (solid line). We also compute the distance between the training data of domain 2 and test data of domain 1 (dashed line), for the same class.

As shown by the solid orange line in Fig. 5b, for minority domain-class pairs such as class “8” and “9” in domain 1, the distance in the feature space between training and testing is large. In fact, the test set of these minority domain-class pairs is closer to the training data for “8” and “9” in domain 2 than in their own domain, as shown by the dashed purple line. This example indicates that a better training would try to transfer the features of minority domain-class pairs to majority pairs with which they share the same class, as shown by the grey arrow in Fig. 5b. Such transfer will improve generalization to the test set.

**BoDA with Calibrated Distance.** The above discussion motivates a modification to BoDA to favor transfer to majority domain-class pairs:

$$\tilde{\mathcal{L}}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) = \sum_{\mathbf{z}_i \in \mathcal{Z}} \frac{-1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp(-\lambda_{d_i, c_i}^{d, c_i} \tilde{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i}))}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp(-\lambda_{d_i, c_i}^{d', c'} \tilde{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))}, \quad \lambda_{d, c}^{d', c'} = \left( \frac{N_{d', c'}}{N_{d, c}} \right)^\nu, \quad (4)$$

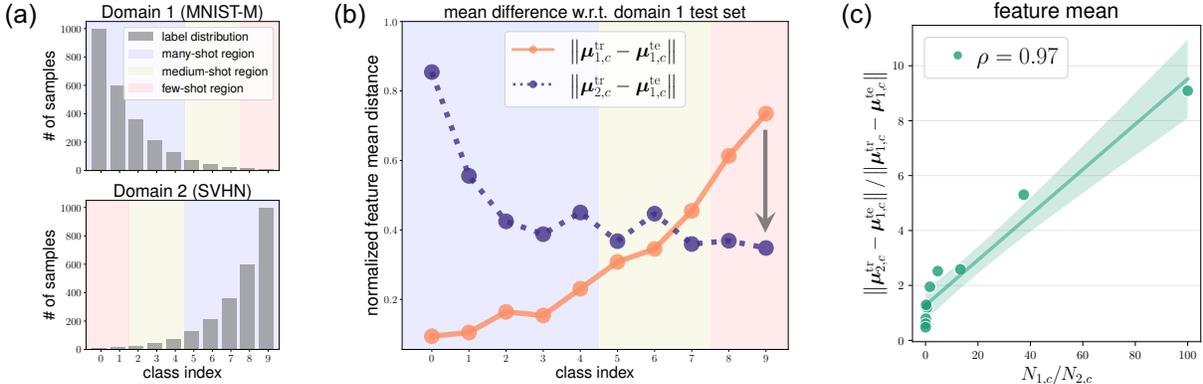


Figure 5: The need for *calibration*. (a) Per-domain label distribution of Digits-MLT. (b) Distance between training and test data. **Solid line** plots the distance between training and test data from the same domain-class pairs. **Dashed line** plots the distance between test data from a particular domain-class pair and the training data with which it shares the same class but differs in the domain. The **blue** and **red** background colors refer to majority and minority domain-class pairs, respectively. (c) Correspondence between the ratio of the sample size and their feature distances between testing and training across different domain-class pairs.

where  $\nu$  is a constant that allows for a sublinear relation (default  $\nu = 1$ ).  $\lambda_{d,c}^{d',c'}$  indicates how much we would like to transfer ( $d, c$ ) to ( $d', c'$ ), based on their relative sample size. Fig. 5c verifies that the ratio of the sample size is highly correlated with the ratio of the distance between testing and training. Further, Theorem 2 in Appendix A shows that  $\tilde{\mathcal{L}}_{\text{BoDA}}$  is an upper bound of the calibrated transferability statistics.

**Variants of BoDA: Matching Higher Order Statistics.** The distance  $d$  can be set to the Euclidean distance  $d(\mathbf{z}, \boldsymbol{\mu}_{d,c}) = \sqrt{(\mathbf{z} - \boldsymbol{\mu}_{d,c})^\top (\mathbf{z} - \boldsymbol{\mu}_{d,c})}$ , which captures the first order statistics. To match higher order statistics such as covariance,  $d(\mathbf{z}, \{\boldsymbol{\mu}_{d,c}, \boldsymbol{\Sigma}_{d,c}\}) = \sqrt{(\mathbf{z} - \boldsymbol{\mu}_{d,c})^\top \boldsymbol{\Sigma}_{d,c}^{-1} (\mathbf{z} - \boldsymbol{\mu}_{d,c})}$  is used, resembling the Mahalanobis distance [11]. We refer to these variants as  $\tilde{\mathcal{L}}_{\text{BoDA}}$  and  $\tilde{\mathcal{L}}_{\text{BoDA-M}}$ .

**Joint Loss.** BoDA serves as a representation learning scheme for MDLT, which operates over  $\mathcal{Z}$ . For classification, we train deep networks by combining  $\tilde{\mathcal{L}}_{\text{BoDA}}$  and the standard cross-entropy (CE) loss in an end-to-end fashion, where CE is applied to the output layer, and BoDA is applied to the latent features. We combine the losses as  $\mathcal{L}_{\text{CE}} + \omega \tilde{\mathcal{L}}_{\text{BoDA}}$ , with  $\omega$  as a trade-off hyperparameter.

## 5 What Makes for Good Classifiers in MDLT?

In the long-tailed recognition literature, an important finding is that decoupling *representation learning* and *classifier learning* leads to better results [23, 58]. In particular, instance-balanced sampling is used during the first stage of learning, while class-balanced sampling is used for re-training the classifier (with the representation fixed) in the second stage [23]. Motivated by this, we explore whether a similar decoupling benefits MDLT. We use three learning algorithms, ERM [46], DANN [31], and CORAL [45]. We train each algorithm with and without the second stage classifier learning, and report the average accuracy over all MDLT datasets (presented later).

As Table 1 shows, similar to what has been observed in the single domain case [23, 58], regardless of algorithm, decoupling the classifier learning consistently improves performance. Since BoDA can support both coupled and decoupled classifier learning, we use  $\text{BoDA}_r$  to refer to models that couple representation and classifier learning, and  $\text{BoDA}_{r,c}$  for models that decouple representation from classifier learning. In the classifier learning stage, we simply use class-balanced sampling.

Table 1: The benefits of decoupling the classifier.

Algorithm	w/o decouple	w/ decouple
ERM [46]	77.6 $\pm$ 0.2	<b>79.2</b> $\pm$ 0.3
DANN [15]	77.7 $\pm$ 0.6	<b>79.0</b> $\pm$ 0.1
CORAL [45]	78.0 $\pm$ 0.1	<b>79.6</b> $\pm$ 0.2

## 6 Benchmarking MDLT

**Datasets.** We curate five multi-domain datasets typically used in DG and adapt them for MDLT evaluation. To do so, for each dataset, we create two balanced datasets one for validation and the other for testing, and leave the rest for training. The size of the validation and test data sets is roughly 5% and 10% of original data, respectively. Table 10 in Appendix D provides the statistics of each MDLT dataset. Fig. 6 shows the label distributions across domains in the five datasets.

1. **VLCS-MLT.** We construct VLCS-MLT using the VLCS dataset [14], which is an object recognition dataset with 10,729 images from 4 domains and 5 classes.
2. **PACS-MLT.** PACS-MLT is constructed from the PACS dataset [27], an object recognition dataset with 9,991 images from 4 domains and 7 classes.
3. **OfficeHome-MLT.** We set up OfficeHome-MLT using the OfficeHome dataset [47] which contains 15,588 images from 4 domains and 65 classes.
4. **TerraInc-MLT.** TerraInc-MLT is created from TerraIncognita [2], a species classification dataset including 24,788 images from 4 domains and 10 classes.
5. **DomainNet-MLT.** We construct DomainNet-MLT using DomainNet [38], a large-scale multi-domain dataset for object recognition. It contains 586,575 images from 345 classes and 6 domains.

**Network Architectures.** For experiments on the synthetic Digits-MLT dataset, we use a simple CNN architecture as in [19]. For the MDLT datasets, we follow [19], and use ResNet-50 [21] for all algorithms.

**Competing Algorithms.** We compare BoDA to a large number of algorithms that span different learning strategies and categories, including (1) *vanilla*: **ERM** [46], (2) *distributionally robust optimization*: **GroupDRO** [40], (3) *data augmentation*: **Mixup** [50], **SagNet** [35], (4) *meta-learning*: **MLDG** [28], (5) *domain-invariant feature learning*: **IRM** [1], **DANN** [15], **CDANN** [31], **CORAL** [45], **MMD** [29], (6) *transfer learning*: **MTL** [4], (7) *multi-task learning*: **Fish** [42], and (8) *imbalanced learning*: **Focal** [32], **CBLoss** [10], **LDAM** [6], **BSoftmax** [39], **SSP** [52], **CRT** [23]. We provide detailed descriptions in Appendix E.2.

**Implementation and Evaluation Metrics.** For a fair evaluation, following [19], for each algorithm we conduct a random search of 20 trials over a joint distribution of all hyperparameters (see Appendix E.3 for details). We then use the validation set to select the best hyperparameters for each algorithm, fix them and rerun the experiments under three different random seeds to report the final average accuracy with standard deviation. Such process ensures the comparison is best-versus-best, and the hyperparameters are optimized for all algorithms. In addition to the average

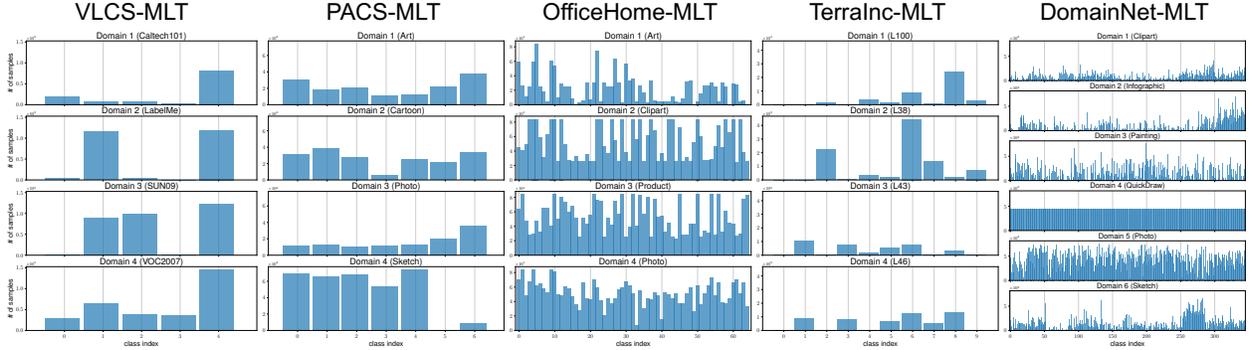


Figure 6: Overview of training set label distribution for five MDLT datasets. We set up MDLT benchmarks from datasets traditionally used for DG, and make validation/test sets balanced across all domain-class pairs. More details are provided in Appendix D.

Table 2: Results on VLCS-MLT.

Algorithm	Accuracy (by domain)		Accuracy (by shot)			
	Average	Worst	Many	Medium	Few	Zero
ERM [46]	76.3 ±0.4	53.6 ±1.1	84.6 ±0.5	76.6 ±0.4	—	32.9 ±0.4
IRM [1]	76.5 ±0.2	52.3 ±0.7	85.3 ±0.6	75.5 ±1.0	—	33.5 ±1.0
GroupDRO [40]	76.7 ±0.4	54.1 ±1.3	85.3 ±0.9	76.2 ±1.0	—	34.5 ±2.0
Mixup [50]	75.9 ±0.1	52.7 ±1.3	84.4 ±0.2	77.1 ±0.6	—	29.2 ±1.4
MLDG [28]	76.9 ±0.2	53.6 ±0.5	84.9 ±0.3	77.5 ±1.0	—	34.4 ±0.9
CORAL [45]	75.9 ±0.5	51.6 ±0.7	84.3 ±0.6	75.5 ±0.5	—	34.5 ±0.8
MMD [29]	76.3 ±0.6	53.4 ±0.3	84.5 ±0.8	77.1 ±0.5	—	32.7 ±0.3
DANN [15]	77.5 ±0.1	54.1 ±0.3	85.9 ±0.5	76.0 ±0.4	—	38.0 ±2.3
CDANN [31]	76.6 ±0.4	53.6 ±0.4	84.4 ±0.7	77.3 ±0.8	—	35.0 ±0.8
MTL [4]	76.3 ±0.3	52.9 ±0.5	84.8 ±0.9	76.2 ±0.6	—	33.3 ±1.4
SagNet [35]	76.3 ±0.2	52.3 ±0.2	85.3 ±0.3	75.1 ±0.2	—	32.9 ±0.3
Fish [42]	77.5 ±0.3	54.3 ±0.4	86.2 ±0.5	76.0 ±0.4	—	35.6 ±2.2
Focal [32]	75.6 ±0.4	52.3 ±0.2	84.0 ±0.2	75.5 ±0.6	—	32.7 ±0.9
CBLoss [10]	76.8 ±0.3	52.5 ±0.5	84.8 ±0.7	77.5 ±1.4	—	33.2 ±1.6
LDAM [6]	77.5 ±0.1	52.9 ±0.2	<b>86.5</b> ±0.4	75.5 ±0.5	—	35.2 ±0.6
BSoftmax [39]	76.7 ±0.5	52.9 ±0.9	84.4 ±0.9	78.2 ±0.6	—	34.3 ±0.9
SSP [52]	76.1 ±0.3	52.3 ±1.0	83.8 ±0.3	76.0 ±1.2	—	37.1 ±0.7
CRT [23]	76.3 ±0.2	51.4 ±0.3	84.5 ±0.1	77.3 ±0.0	—	31.7 ±1.0
BoDA <sub>r</sub>	76.9 ±0.5	51.4 ±0.3	85.3 ±0.3	77.3 ±0.2	—	33.3 ±0.5
BoDA-M <sub>r</sub>	77.5 ±0.3	53.4 ±0.3	85.8 ±0.2	77.3 ±0.2	—	35.7 ±0.7
BoDA <sub>r,c</sub>	77.3 ±0.2	53.4 ±0.3	85.3 ±0.3	78.0 ±0.2	—	38.6 ±0.7
BoDA-M <sub>r,c</sub>	<b>78.2</b> ±0.4	<b>55.4</b> ±0.5	85.3 ±0.3	<b>79.3</b> ±0.6	—	<b>43.3</b> ±1.1
BoDA vs. ERM	<b>+1.9</b>	<b>+1.8</b>	<b>+0.7</b>	<b>+2.7</b>	—	<b>+10.4</b>

Table 3: Results on PACS-MLT.

Algorithm	Accuracy (by domain)		Accuracy (by shot)			
	Average	Worst	Many	Medium	Few	Zero
ERM [46]	97.1 ±0.1	95.8 ±0.2	97.1 ±0.0	97.0 ±0.0	98.0 ±0.9	—
IRM [1]	96.7 ±0.2	95.2 ±0.4	96.8 ±0.2	96.7 ±0.7	94.7 ±1.4	—
GroupDRO [40]	97.0 ±0.1	95.3 ±0.4	97.3 ±0.1	95.3 ±1.2	94.7 ±3.6	—
Mixup [50]	96.7 ±0.2	95.1 ±0.2	97.0 ±0.1	96.7 ±0.3	91.3 ±2.7	—
MLDG [28]	96.6 ±0.1	94.1 ±0.3	96.8 ±0.1	96.3 ±0.7	92.7 ±0.5	—
CORAL [45]	96.6 ±0.5	94.3 ±0.7	96.6 ±0.5	97.0 ±0.8	94.7 ±0.5	—
MMD [29]	96.9 ±0.1	96.2 ±0.2	96.9 ±0.2	97.0 ±0.0	96.7 ±0.5	—
DANN [15]	96.5 ±0.0	94.3 ±0.1	96.5 ±0.1	98.0 ±0.0	94.7 ±2.4	—
CDANN [31]	96.1 ±0.1	94.5 ±0.2	96.1 ±0.1	96.3 ±0.5	94.0 ±0.9	—
MTL [4]	96.7 ±0.2	94.5 ±0.6	96.8 ±0.1	95.3 ±1.7	97.3 ±1.1	—
SagNet [35]	<b>97.2</b> ±0.1	95.2 ±0.3	<b>97.4</b> ±0.1	96.7 ±0.5	95.3 ±0.5	—
Fish [42]	96.9 ±0.2	95.2 ±0.2	97.0 ±0.1	97.0 ±0.5	94.7 ±1.1	—
Focal [32]	96.5 ±0.2	94.6 ±0.7	96.6 ±0.1	95.0 ±1.7	96.7 ±0.5	—
CBLoss [10]	96.9 ±0.1	95.1 ±0.4	96.8 ±0.2	97.0 ±1.2	<b>100.0</b> ±0.0	—
LDAM [6]	96.5 ±0.2	94.7 ±0.2	96.6 ±0.1	95.7 ±1.4	96.0 ±0.0	—
BSoftmax [39]	96.9 ±0.3	95.6 ±0.3	96.6 ±0.4	<b>98.7</b> ±0.7	99.3 ±0.5	—
SSP [52]	96.9 ±0.2	95.4 ±0.4	96.7 ±0.2	98.3 ±0.5	98.0 ±0.9	—
CRT [23]	96.3 ±0.1	94.9 ±0.1	96.3 ±0.1	97.3 ±0.3	94.0 ±0.9	—
BoDA <sub>r</sub>	97.0 ±0.1	95.1 ±0.4	97.0 ±0.1	96.3 ±0.5	98.0 ±0.9	—
BoDA-M <sub>r</sub>	97.1 ±0.1	94.9 ±0.1	97.3 ±0.1	96.3 ±0.5	96.0 ±0.0	—
BoDA <sub>r,c</sub>	<b>97.2</b> ±0.1	95.7 ±0.3	<b>97.4</b> ±0.1	97.0 ±0.0	94.7 ±1.1	—
BoDA-M <sub>r,c</sub>	97.1 ±0.2	<b>96.3</b> ±0.1	97.1 ±0.0	97.0 ±0.8	96.0 ±0.0	—
BoDA vs. ERM	<b>+0.1</b>	<b>+0.5</b>	<b>+0.3</b>	<b>+0.0</b>	<b>-2.0</b>	—

accuracy across domains, we also report the worst accuracy over domains, and further divide all domain-class pairs into *many-shot* (pairs with over 100 training samples), *medium-shot* (pairs with 20~100 training samples), *few-shot* (pairs with under 20 training samples), and *zero-shot* (pairs with no training data), and report the results for these subsets.

## 6.1 Main Results

We report the main results in this section for all MDLT datasets. The complete results and all additional experiments are provided in Appendix F and H.

**Benchmark Results on MDLT Datasets.** The performance of all methods on VLCS-MLT, PACS-MLT, OfficeHome-MLT, TerraInc-MLT and DomainNet-MLT are in Table 2, 3, 4, 5 and 6, respectively. We highlight rows in gray for BoDA and its variants, and bolden the best result in each

Table 4: Results on OfficeHome-MLT.

Algorithm	Accuracy (by domain)		Accuracy (by shot)			
	Average	Worst	Many	Medium	Few	Zero
ERM [46]	80.7 ±0.0	71.3 ±0.1	87.8 ±0.2	81.0 ±0.2	63.1 ±0.1	63.3 ±7.2
IRM [1]	80.6 ±0.4	70.7 ±0.2	87.6 ±0.4	81.5 ±0.4	61.1 ±0.9	56.7 ±1.4
GroupDRO [40]	80.1 ±0.3	68.7 ±0.9	88.1 ±0.2	80.8 ±0.4	59.8 ±1.2	51.7 ±3.6
Mixup [50]	81.2 ±0.2	72.3 ±0.6	87.9 ±0.4	81.8 ±0.1	64.1 ±0.4	60.0 ±4.1
MLDG [28]	80.4 ±0.2	70.2 ±0.6	87.1 ±0.1	81.3 ±0.3	61.3 ±1.0	61.7 ±1.4
CORAL [45]	81.9 ±0.1	<b>72.7</b> ±0.6	87.9 ±0.1	83.0 ±0.1	63.5 ±0.7	65.0 ±2.4
MMD [29]	78.4 ±0.4	67.7 ±0.8	85.2 ±0.2	79.4 ±0.7	58.8 ±0.4	56.7 ±3.6
DANN [15]	79.2 ±0.2	70.2 ±0.9	86.2 ±0.1	80.0 ±0.1	60.3 ±1.1	61.7 ±5.9
CDANN [31]	79.0 ±0.2	69.4 ±0.3	86.4 ±0.6	79.8 ±0.1	58.9 ±0.8	50.0 ±4.7
MTL [4]	79.5 ±0.2	69.8 ±0.6	87.3 ±0.3	79.8 ±0.2	61.1 ±0.2	51.7 ±2.7
SagNet [35]	80.9 ±0.1	70.5 ±0.5	87.8 ±0.4	81.9 ±0.1	61.2 ±0.9	56.7 ±3.6
Fish [42]	81.3 ±0.3	71.3 ±0.7	<b>88.2</b> ±0.2	81.9 ±0.3	63.2 ±0.8	61.7 ±1.4
Focal [32]	77.9 ±0.0	67.6 ±0.4	86.5 ±0.3	78.3 ±0.1	57.4 ±0.3	46.7 ±3.6
CBLoss [10]	79.8 ±0.2	69.5 ±0.7	86.6 ±0.4	80.6 ±0.2	61.1 ±1.4	65.0 ±2.4
LDAM [6]	80.3 ±0.2	69.9 ±0.5	87.1 ±0.2	81.3 ±0.3	61.1 ±0.2	51.7 ±2.7
BSoftmax [39]	80.4 ±0.2	70.9 ±0.5	86.7 ±0.5	81.3 ±0.3	62.4 ±1.0	60.0 ±4.1
SSP [52]	81.1 ±0.3	71.1 ±0.3	87.3 ±0.6	82.3 ±0.3	61.6 ±0.7	63.3 ±1.4
CRT [23]	81.2 ±0.0	72.5 ±0.2	87.7 ±0.1	81.8 ±0.1	64.0 ±0.1	65.0 ±2.4
BoDA <sub>r</sub>	81.5 ±0.1	71.8 ±0.1	87.7 ±0.2	82.3 ±0.1	<b>64.2</b> ±0.3	63.3 ±1.4
BoDA-M <sub>r</sub>	81.9 ±0.2	71.6 ±0.2	87.3 ±0.3	83.4 ±0.2	62.3 ±0.3	65.0 ±2.4
BoDA <sub>r,c</sub>	82.3 ±0.1	72.3 ±0.3	87.1 ±0.2	<b>83.9</b> ±0.3	63.2 ±0.2	65.0 ±2.4
BoDA-M <sub>r,c</sub>	<b>82.4</b> ±0.2	72.3 ±0.3	87.7 ±0.1	<b>83.9</b> ±0.6	<b>64.2</b> ±0.3	<b>66.7</b> ±2.7
BoDA vs. ERM	<b>+1.7</b>	<b>+1.0</b>	<b>-0.1</b>	<b>+2.9</b>	<b>+1.1</b>	<b>+3.4</b>

Table 5: Results on TerraInc-MLT.

Algorithm	Accuracy (by domain)		Accuracy (by shot)			
	Average	Worst	Many	Medium	Few	Zero
ERM [46]	75.3 ±0.3	67.4 ±0.3	85.6 ±0.8	69.6 ±3.2	66.1 ±2.4	14.4 ±2.8
IRM [1]	73.3 ±0.7	64.3 ±1.3	83.5 ±0.6	70.0 ±1.8	58.3 ±3.4	20.1 ±1.4
GroupDRO [40]	72.0 ±0.4	66.6 ±0.2	84.7 ±1.1	64.6 ±4.7	38.9 ±1.2	13.5 ±1.1
Mixup [50]	71.1 ±0.7	60.4 ±1.1	83.2 ±0.7	60.0 ±0.6	56.1 ±3.0	12.2 ±2.1
MLDG [28]	76.6 ±0.2	66.9 ±0.5	86.1 ±0.6	73.8 ±3.9	70.6 ±3.7	18.8 ±2.4
CORAL [45]	76.4 ±0.5	67.8 ±0.9	86.3 ±0.3	77.5 ±3.1	66.1 ±2.0	11.0 ±1.4
MMD [29]	73.3 ±0.4	63.7 ±1.1	84.0 ±0.4	67.9 ±2.7	60.6 ±1.6	13.6 ±2.6
DANN [15]	68.7 ±0.9	61.1 ±1.0	79.6 ±1.2	62.5 ±8.1	48.9 ±2.8	13.3 ±1.1
CDANN [31]	70.3 ±0.5	63.9 ±1.0	83.5 ±0.8	50.0 ±4.2	43.9 ±4.7	20.4 ±3.1
MTL [4]	75.0 ±0.7	67.7 ±1.4	85.2 ±0.7	73.8 ±1.6	61.1 ±2.8	12.4 ±4.0
SagNet [35]	75.1 ±1.6	66.5 ±2.1	85.5 ±0.9	77.1 ±5.0	57.8 ±4.3	13.0 ±3.4
Fish [42]	75.3 ±0.5	66.3 ±0.5	85.8 ±0.2	73.3 ±3.9	61.1 ±3.0	13.7 ±3.3
Focal [32]	75.7 ±0.4	65.3 ±1.1	85.7 ±0.3	76.2 ±3.9	68.9 ±3.2	12.6 ±1.9
CBLoss [10]	78.0 ±0.4	68.3 ±2.0	85.0 ±0.1	89.2 ±1.2	83.9 ±2.5	9.3 ±3.9
LDAM [6]	74.7 ±0.9	64.1 ±1.4	85.1 ±0.6	70.8 ±3.5	67.8 ±1.2	11.1 ±2.4
BSoftmax [39]	76.7 ±1.0	65.6 ±1.3	83.4 ±0.8	90.8 ±0.9	78.3 ±3.9	12.6 ±2.4
SSP [52]	78.5 ±0.7	67.3 ±0.4	85.5 ±1.0	87.8 ±0.9	82.6 ±1.2	13.2 ±2.8
CRT [23]	81.6 ±0.1	70.0 ±0.4	<b>89.7</b> ±0.2	90.4 ±0.3	83.9 ±0.5	12.9 ±0.0
BoDA <sub>r</sub>	78.6 ±0.4	68.5 ±0.3	86.4 ±0.1	85.0 ±1.0	80.0 ±0.9	13.7 ±2.1
BoDA-M <sub>r</sub>	79.4 ±0.6	71.3 ±0.4	88.4 ±0.3	76.2 ±2.7	88.3 ±1.6	14.4 ±1.4
BoDA <sub>r,c</sub>	82.3 ±0.3	68.5 ±0.6	89.2 ±0.2	<b>92.5</b> ±0.9	88.3 ±1.2	21.3 ±0.7
BoDA-M <sub>r,c</sub>	<b>83.0</b> ±0.4	<b>74.6</b> ±0.7	89.2 ±0.2	91.2 ±0.6	<b>91.7</b> ±2.0	<b>21.7</b> ±1.4
BoDA vs. ERM	<b>+7.7</b>	<b>+7.2</b>	<b>+3.6</b>	<b>+22.9</b>	<b>+25.6</b>	<b>+7.3</b>

Table 6: Results on DomainNet-MLT.

Algorithm	Accuracy (by domain)		Accuracy (by shot)			
	Average	Worst	Many	Medium	Few	Zero
ERM [46]	58.6 ±0.2	29.4 ±0.3	66.0 ±0.1	56.1 ±0.1	35.9 ±0.5	27.6 ±0.3
IRM [1]	57.1 ±0.1	27.6 ±0.1	64.7 ±0.1	54.3 ±0.3	33.5 ±0.3	25.8 ±0.3
GroupDRO [40]	53.6 ±0.1	25.9 ±0.2	61.8 ±0.1	49.1 ±0.3	30.7 ±0.7	22.0 ±0.1
Mixup [50]	57.6 ±0.1	28.7 ±0.0	64.9 ±0.2	54.5 ±0.1	35.6 ±0.2	27.3 ±0.3
MLDG [28]	58.5 ±0.0	28.7 ±0.1	66.0 ±0.1	55.7 ±0.1	35.3 ±0.2	26.9 ±0.3
CORAL [45]	59.4 ±0.1	30.1 ±0.4	66.4 ±0.1	57.1 ±0.0	37.7 ±0.6	29.9 ±0.2
MMD [29]	56.7 ±0.0	27.2 ±0.2	64.2 ±0.1	54.0 ±0.0	33.9 ±0.2	25.4 ±0.2
DANN [15]	55.8 ±0.1	26.9 ±0.4	63.0 ±0.1	52.7 ±0.1	34.2 ±0.4	26.8 ±0.4
CDANN [31]	56.0 ±0.1	27.7 ±0.1	63.2 ±0.0	52.7 ±0.2	34.3 ±0.5	27.6 ±0.1
MTL [4]	58.6 ±0.1	29.3 ±0.2	65.9 ±0.1	56.0 ±0.4	35.4 ±0.1	28.2 ±0.3
SagNet [35]	58.9 ±0.0	29.4 ±0.2	66.3 ±0.1	56.4 ±0.0	36.2 ±0.3	27.2 ±0.4
Fish [42]	59.6 ±0.1	29.1 ±0.1	67.1 ±0.1	57.2 ±0.1	36.8 ±0.4	27.8 ±0.3
Focal [32]	57.8 ±0.2	27.5 ±0.1	65.2 ±0.2	55.1 ±0.2	35.8 ±0.1	26.3 ±0.1
CBLoss [10]	58.9 ±0.1	30.1 ±0.1	64.3 ±0.0	61.0 ±0.3	42.5 ±0.4	28.1 ±0.2
LDAM [6]	59.2 ±0.0	29.2 ±0.2	66.6 ±0.0	57.0 ±0.0	37.1 ±0.2	27.8 ±0.3
BSoftmax [39]	58.9 ±0.1	29.9 ±0.1	64.3 ±0.1	60.9 ±0.3	42.4 ±0.6	28.2 ±0.1
SSP [52]	59.7 ±0.0	31.6 ±0.2	64.3 ±0.1	62.6 ±0.1	45.0 ±0.3	30.5 ±0.0
CRT [23]	60.4 ±0.2	31.6 ±0.1	66.8 ±0.0	61.6 ±0.1	45.7 ±0.1	29.7 ±0.1
BoDA <sub>r</sub>	60.1 ±0.2	32.6 ±0.1	65.7 ±0.2	60.6 ±0.1	42.6 ±0.3	30.5 ±0.2
BoDA-M <sub>r</sub>	60.1 ±0.2	32.2 ±0.2	65.9 ±0.2	60.7 ±0.1	42.9 ±0.3	30.0 ±0.1
BoDA <sub>r,c</sub>	<b>61.7</b> ±0.1	<b>33.4</b> ±0.1	<b>67.0</b> ±0.1	62.7 ±0.1	46.0 ±0.2	<b>32.2</b> ±0.3
BoDA-M <sub>r,c</sub>	<b>61.7</b> ±0.2	33.3 ±0.1	<b>67.0</b> ±0.1	<b>63.0</b> ±0.3	<b>46.6</b> ±0.4	31.8 ±0.2
BoDA vs. ERM	<b>+3.1</b>	<b>+4.0</b>	<b>+1.0</b>	<b>+6.9</b>	<b>+10.7</b>	<b>+4.6</b>

Table 7: Results over all MDLT benchmarks.

Algorithm	VLCS-MLT	PACS-MLT	OfficeHome-MLT	TerraInc-MLT	DomainNet-MLT	Avg
ERM [46]	76.3 ±0.4	97.1 ±0.1	80.7 ±0.0	75.3 ±0.3	58.6 ±0.2	77.6
IRM [1]	76.5 ±0.2	96.7 ±0.2	80.6 ±0.4	73.3 ±0.7	57.1 ±0.1	76.8
GroupDRO [40]	76.7 ±0.4	97.0 ±0.1	80.1 ±0.3	72.0 ±0.4	53.6 ±0.1	75.9
Mixup [50]	75.9 ±0.1	96.7 ±0.2	81.2 ±0.2	71.1 ±0.7	57.6 ±0.1	76.5
MLDG [28]	76.9 ±0.2	96.6 ±0.1	80.4 ±0.2	76.6 ±0.2	58.5 ±0.0	77.8
CORAL [45]	75.9 ±0.5	96.6 ±0.5	81.9 ±0.1	76.4 ±0.5	59.4 ±0.1	78.0
MMD [29]	76.3 ±0.6	96.9 ±0.1	78.4 ±0.4	73.3 ±0.4	56.7 ±0.0	76.3
DANN [15]	77.5 ±0.1	96.5 ±0.0	79.2 ±0.2	68.7 ±0.9	55.8 ±0.1	75.5
CDANN [31]	76.6 ±0.4	96.1 ±0.1	79.0 ±0.2	70.3 ±0.5	56.0 ±0.1	75.6
MTL [4]	76.3 ±0.3	96.7 ±0.2	79.5 ±0.2	75.0 ±0.7	58.6 ±0.1	77.2
SagNet [35]	76.3 ±0.2	<b>97.2</b> ±0.1	80.9 ±0.1	75.1 ±1.6	58.9 ±0.0	77.7
Fish [42]	77.5 ±0.3	96.9 ±0.2	81.3 ±0.3	75.3 ±0.5	59.6 ±0.1	78.1
Focal [32]	75.6 ±0.4	96.5 ±0.2	77.9 ±0.0	75.7 ±0.4	57.8 ±0.2	76.7
CBLoss [10]	76.8 ±0.3	96.9 ±0.1	79.8 ±0.2	78.0 ±0.2	58.9 ±0.1	78.1
LDAM [6]	77.5 ±0.1	96.5 ±0.2	80.3 ±0.2	74.7 ±0.9	59.2 ±0.0	77.7
BSoftmax [39]	76.7 ±0.5	96.9 ±0.3	80.4 ±0.2	76.7 ±1.0	58.9 ±0.1	77.9
SSP [52]	76.1 ±0.3	96.9 ±0.2	81.1 ±0.3	78.5 ±0.7	59.7 ±0.0	78.5
CRT [23]	76.3 ±0.2	96.3 ±0.1	81.2 ±0.0	81.6 ±0.1	60.4 ±0.2	79.2
BoDA <sub>r</sub>	76.9 ±0.5	97.0 ±0.1	81.5 ±0.1	78.6 ±0.4	60.1 ±0.2	78.8
BoDA-M <sub>r</sub>	77.5 ±0.3	97.1 ±0.1	81.9 ±0.2	79.4 ±0.6	60.1 ±0.2	79.2
BoDA <sub>r,c</sub>	77.3 ±0.2	<b>97.2</b> ±0.1	82.3 ±0.1	82.3 ±0.3	<b>61.7</b> ±0.1	80.2
BoDA-M <sub>r,c</sub>	<b>78.2</b> ±0.4	97.1 ±0.2	<b>82.4</b> ±0.2	<b>83.0</b> ±0.4	<b>61.7</b> ±0.2	<b>80.5</b>
BoDA vs. ERM	<b>+1.9</b>	<b>+0.1</b>	<b>+1.7</b>	<b>+7.7</b>	<b>+3.1</b>	<b>+2.9</b>

column. First, as all tables indicate, BoDA consistently achieves the best average accuracy across all datasets. It also achieves the best worst-case accuracy most of the time. Moreover, on certain datasets (e.g., OfficeHome-MLT), MDL methods perform better (e.g., CORAL), while on others (e.g., TerraInc-MLT), imbalanced methods achieve higher gains (e.g., CRT); Nevertheless, regardless of dataset, BoDA outperforms all methods, highlighting its effectiveness for the MDLT task. Finally, compared to ERM, BoDA slightly improves the average and many-shot performance, while substantially boosting the performance for the medium-shot, few-shot, and zero-shot pairs. Table 7 summarizes the averaged accuracy across all datasets, where BoDA brings large overall improvements of  $\sim 3\%$ .

**A Closer Look at Accuracy Gains.** We further explore how BoDA performs across *all* domain-class pairs. Fig. 7 shows the absolute accuracy gains of BoDA over ERM on OfficeHome-MLT, where

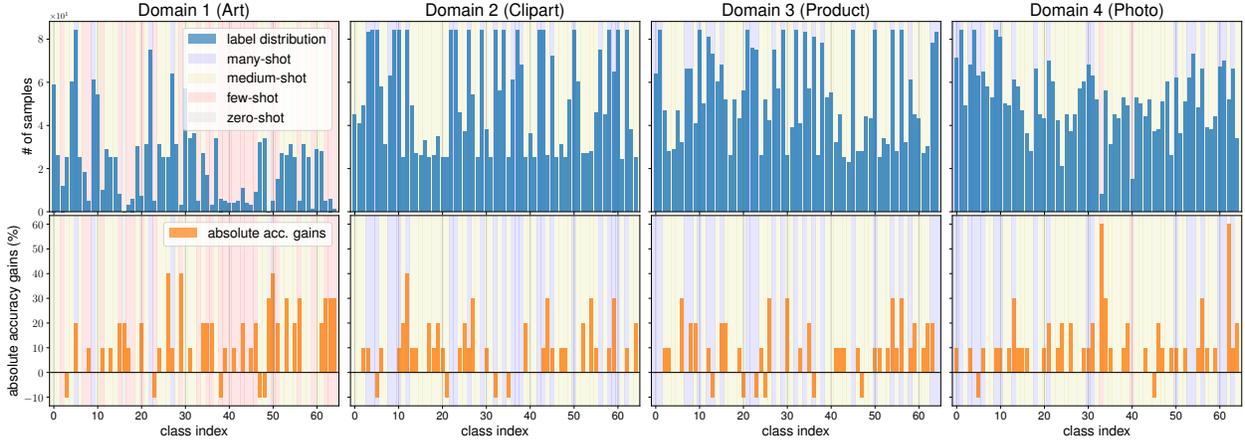


Figure 7: The absolute accuracy improvements of BoDA *vs.* ERM over all domain-class pairs on OfficeHome-MLT. BoDA establishes large improvements w.r.t. all regions, especially for the few-shot and zero-shot ones. Results for other datasets are in Appendix H.2.

BoDA consistently improves the performance over all domains. The improvements are especially large for domain “Art”, where most of the classes lie in the *few-shot* region. For certain classes, BoDA can improve up to 50% accuracy, indicating its effectiveness on tackling MDLT.

**Ablation Studies on BoDA Components (Appendix H.1).** We study the effects of (1) adding balanced distance (i.e., BoDA *vs.* vanilla DA), and (2) different choices of distance calibration coefficient  $\lambda_{d,c}^{d',c'}$  in BoDA. We observe that BoDA improves over DA by a large margin (2.3% on average over all MDLT datasets), highlighting the importance of using *balanced* distance. Interestingly, as for  $\lambda_{d,c}^{d',c'}$ , we find that BoDA is pretty robust to different choices within a given range, and obtain similar gains (1.9% to 2.9% over ERM).

## 6.2 Understanding the Behavior of BoDA on MDLT

To better understand how the design of BoDA contributes to its ability to outperform other algorithms, we go back to the Digits-MLT dataset, but this time we run BoDA as opposed to ERM.

**Better Learned Representations for Minority Data.** Similar to Fig. 5, we plot in Fig. 8b the feature mean distance between training and test data for BoDA on Digits-MLT. The plot shows that BoDA learns better representations with smaller feature discrepancy, especially for minority classes.

**Improved Transferability against Severe Imbalance.** Fig. 8c plots the transferability graph induced by BoDA. It shows that even in the presence of severe and divergent label imbalance (Fig. 8a), BoDA still learns transferable features. Further, BoDA learns a *balanced* feature space that separates different classes away. The better learned features translate to better accuracy (9.5% absolute accuracy gains *vs.* ERM in Fig. 3c). We provide more related results in Appendix H.3 and H.5.

**Tightness of the Bound.** We study whether the BoDA bound derived in Theorem 1 is tight. We train a ResNet-18 on Digits-MLT for 5,000 steps to ensure convergence. We compute the loss over all samples, and combine the results over 3 random seeds. Table 8 confirms the bound is empirically tight.

Table 8: BoDA bound.

	$\mathcal{L}_{\text{BoDA}}$
Empirical	$2.92947 \pm 7.3\text{e-}3$
Theoretical	$2.92513 \pm 7.8\text{e-}3$

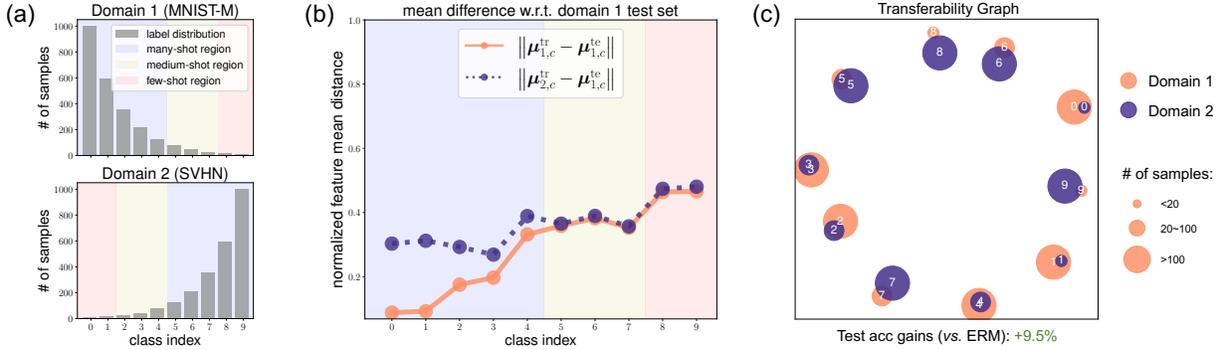


Figure 8: BoDA analysis. (a) Label distribution setup. (b) Distance of feature mean between train and test data. BoDA enables better learned tail ( $d, c$ ) with smaller feature discrepancy. (c) BoDA learns features that are more aligned across domains even in the presence of divergent labels, and significantly improves upon ERM by 9.5%.

Table 9: BoDA strengthens performance on Domain Generalization (DG) benchmarks. Full tables including detailed results for each DG dataset are provided in Appendix G.

Algorithm	VLCS	PACS	OfficeHome	TerraInc	DomainNet	Avg
ERM	77.5 $\pm$ 0.4	85.5 $\pm$ 0.2	66.5 $\pm$ 0.3	46.1 $\pm$ 1.8	40.9 $\pm$ 0.1	63.3
Current SOTA [45]	<b>78.8</b> $\pm$ 0.6	86.2 $\pm$ 0.3	68.7 $\pm$ 0.3	47.6 $\pm$ 1.0	41.5 $\pm$ 0.1	64.5
BoDA <sub>r,c</sub>	78.5 $\pm$ 0.3	<b>86.9</b> $\pm$ 0.4	<b>69.3</b> $\pm$ 0.1	<b>50.2</b> $\pm$ 0.4	<b>42.7</b> $\pm$ 0.1	<b>65.5</b>
BoDA <sub>r,c</sub> + Current SOTA [45]	79.1 $\pm$ 0.1	87.9 $\pm$ 0.5	69.9 $\pm$ 0.2	50.7 $\pm$ 0.6	43.5 $\pm$ 0.3	66.2
BoDA vs. ERM	<b>+1.6</b>	<b>+2.4</b>	<b>+3.4</b>	<b>+4.6</b>	<b>+2.6</b>	<b>+2.9</b>

## 7 Beyond MDLT: (Imbalanced) Domain Generalization

Domain Generalization (DG) refers to learning from multiple domains and generalizing to unseen domains. Since naturally the learning domains differ in their label distributions and may even have class imbalance within each domain, we investigate whether tackling cross-domain data imbalance can further strengthen the performance for DG. Note that all datasets we adapted for MDLT are standard benchmarks for DG, which confirms that data imbalance is an intrinsic problem in DG, but has been overlooked by past works.

We study whether BoDA can improve performance for DG. To test BoDA, we follow standard DG evaluation protocol [19], and compare to the current SOTA [45]. Table 9 reveals the following findings: First, BoDA alone can improve upon the current SOTA on four out of the five datasets, and achieves notable average performance gains. Moreover, combined with the current SOTA, BoDA further boosts the result by a notable margin across all datasets, suggesting that label imbalance is orthogonal to existing DG-specific algorithms. Finally, similar to MDLT, the gains depend on how severe the imbalance is within a dataset – e.g., TerraInc exhibits the most severe label imbalance across domains, on which BoDA achieves the highest gains. Detailed results for each DG dataset are provided in Appendix G. These intriguing results shed light on how label imbalance can affect out-of-distribution generalization, and highlight the importance of integrating label imbalance for practical DG algorithm design.

## 8 Conclusion

We formalize the MDLT task as learning from multi-domain imbalanced data, and generalizing to all domain-class pairs. We introduce the domain-class transferability graph, and propose BoDA, a theoretically grounded loss that tackles MDLT. Extensive results on five curated real-world MDLT benchmarks verify its superiority. Furthermore, incorporating BoDA into DG algorithms establishes a new SOTA on DG benchmarks. Our work opens up new avenues for realistic multi-domain learning and generalization in the presence of data imbalance.

## Acknowledgments

This work is supported by the GIST-MIT Research Collaboration grant funded by GIST. Yuzhe Yang is supported by the MathWorks Fellowship.

## References

- [1] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- [2] S. Beery, G. Van Horn, and P. Perona. Recognition in terra incognita. In *ECCV*, 2018.
- [3] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.
- [4] G. Blanchard, A. A. Deshmukh, U. Dogan, G. Lee, and C. Scott. Domain generalization by marginal transfer learning. *Journal of Machine Learning Research*, 22(2):1–55, 2021.
- [5] M. Buda, A. Maki, and M. A. Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106:249–259, 2018.
- [6] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *NeurIPS*, 2019.
- [7] F. M. Carlucci, A. D’Innocente, S. Bucci, B. Caputo, and T. Tommasi. Domain generalization by solving jigsaw puzzles. In *CVPR*, 2019.
- [8] J. D. Carroll and P. Arabie. Multidimensional scaling. *Measurement, judgment and decision making*, pages 179–250, 1998.
- [9] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- [10] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, 2019.
- [11] R. De Maesschalck, D. Jouan-Rimbaud, and D. L. Massart. The mahalanobis distance. *Chemo-metrics and intelligent laboratory systems*, 50(1):1–18, 2000.

- [12] Q. Dong, S. Gong, and X. Zhu. Imbalanced deep learning by minority class incremental rectification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(6):1367–1381, 2019.
- [13] M. Dredze, A. Kulesza, and K. Crammer. Multi-domain learning by confidence-weighted parameter combination. *Machine Learning*, 79(1):123–149, 2010.
- [14] C. Fang, Y. Xu, and D. N. Rockmore. Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias. In *ICCV*, 2013.
- [15] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *Journal of machine learning research*, 17(1):2096–2030, 2016.
- [16] A. Globerson, G. Chechik, F. Pereira, and N. Tishby. Euclidean embedding of co-occurrence data. In *NeurIPS*, 2004.
- [17] J. Goldberger, G. E. Hinton, S. Roweis, and R. R. Salakhutdinov. Neighbourhood components analysis. In *NeurIPS*, 2004.
- [18] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773, 2012.
- [19] I. Gulrajani and D. Lopez-Paz. In search of lost domain generalization. In *ICLR*, 2021.
- [20] H. He, Y. Bai, E. A. Garcia, and S. Li. Adasyn: Adaptive synthetic sampling approach for imbalanced learning. In *IEEE international joint conference on neural networks*, pages 1322–1328, 2008.
- [21] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [22] C. Huang, Y. Li, C. L. Chen, and X. Tang. Deep imbalanced learning for face recognition and attribute prediction. *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [23] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis. Decoupling representation and classifier for long-tailed recognition. *ICLR*, 2020.
- [24] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [25] D. Krueger, E. Caballero, J.-H. Jacobsen, A. Zhang, J. Binas, R. L. Priol, and A. Courville. Out-of-distribution generalization via risk extrapolation (rex). *arXiv preprint arXiv:2003.00688*, 2020.
- [26] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [27] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales. Deeper, broader and artier domain generalization. In *ICCV*, 2017.
- [28] D. Li, Y. Yang, Y.-Z. Song, and T. Hospedales. Learning to generalize: Meta-learning for domain generalization. In *AAAI*, 2018.

- [29] H. Li, S. J. Pan, S. Wang, and A. C. Kot. Domain generalization with adversarial feature learning. In *CVPR*, 2018.
- [30] T. Li, P. Cao, Y. Yuan, L. Fan, Y. Yang, R. Feris, P. Indyk, and D. Katabi. Targeted supervised contrastive learning for long-tailed recognition. *arXiv preprint arXiv:2111.13998*, 2021.
- [31] Y. Li, M. Gong, X. Tian, T. Liu, and D. Tao. Domain generalization via conditional invariant representations. In *AAAI*, 2018.
- [32] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *ICCV*, pages 2980–2988, 2017.
- [33] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu. Large-scale long-tailed recognition in an open world. In *CVPR*, 2019.
- [34] K. Muandet, D. Balduzzi, and B. Schölkopf. Domain generalization via invariant feature representation. In *ICML*, 2013.
- [35] H. Nam, H. Lee, J. Park, W. Yoon, and D. Yoo. Reducing domain gap by reducing style bias. In *CVPR*, 2021.
- [36] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [37] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009.
- [38] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang. Moment matching for multi-source domain adaptation. In *ICCV*, 2019.
- [39] J. Ren, C. Yu, X. Ma, H. Zhao, S. Yi, et al. Balanced meta-softmax for long-tailed visual recognition. In *NeurIPS*, 2020.
- [40] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *ICLR*, 2020.
- [41] A. Schoenauer-Sebag, L. Heinrich, M. Schoenauer, M. Sebag, L. F. Wu, and S. J. Altschuler. Multi-domain adversarial learning. In *ICLR*, 2019.
- [42] Y. Shi, J. Seely, P. H. Torr, N. Siddharth, A. Hannun, N. Usunier, and G. Synnaeve. Gradient matching for domain generalization. *arXiv preprint arXiv:2104.09937*, 2021.
- [43] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng. Meta-weight-net: Learning an explicit mapping for sample weighting. *arXiv preprint arXiv:1902.07379*, 2019.
- [44] K. Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NeurIPS*, 2016.
- [45] B. Sun and K. Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *ECCV*, 2016.

- [46] V. N. Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5):988–999, 1999.
- [47] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan. Deep hashing network for unsupervised domain adaptation. In *CVPR*, 2017.
- [48] X. Wang, L. Lian, Z. Miao, Z. Liu, and S. Yu. Long-tailed recognition by routing diverse distribution-aware experts. In *ICLR*, 2021.
- [49] T. Xiao, H. Li, W. Ouyang, and X. Wang. Learning deep feature representations with domain guided dropout for person re-identification. In *CVPR*, 2016.
- [50] M. Xu, J. Zhang, B. Ni, T. Li, C. Wang, Q. Tian, and W. Zhang. Adversarial domain adaptation with domain mixup. In *AAAI*, 2020.
- [51] Y. Yang and T. M. Hospedales. A unified perspective on multi-domain and multi-task learning. In *ICLR*, 2015.
- [52] Y. Yang and Z. Xu. Rethinking the value of labels for improving class-imbalanced learning. In *NeurIPS*, 2020.
- [53] Y. Yang, K. Zha, Y.-C. Chen, H. Wang, and D. Katabi. Delving into deep imbalanced regression. In *ICML*, 2021.
- [54] M. Zhang, H. Marklund, A. Gupta, S. Levine, and C. Finn. Adaptive risk minimization: A meta-learning approach for tackling group shift. *arXiv preprint arXiv:2007.02931*, 2020.
- [55] X. Zhang, Z. Fang, Y. Wen, Z. Li, and Y. Qiao. Range loss for deep face recognition with long-tailed training data. In *ICCV*, 2017.
- [56] Y. Zhang, B. Hooi, L. Hong, and J. Feng. Test-agnostic long-tailed recognition by test-time aggregating diverse experts with self-supervision. *arXiv preprint arXiv:2107.09249*, 2021.
- [57] Y. Zhang, B. Kang, B. Hooi, S. Yan, and J. Feng. Deep long-tailed learning: A survey. *arXiv preprint arXiv:2110.04596*, 2021.
- [58] B. Zhou, Q. Cui, X.-S. Wei, and Z.-M. Chen. Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition. *CVPR*, 2020.
- [59] K. Zhou, Z. Liu, Y. Qiao, T. Xiang, and C. C. Loy. Domain generalization in vision: A survey. *arXiv preprint arXiv:2103.02503*, 2021.
- [60] K. Zhou, Y. Yang, Y. Qiao, and T. Xiang. Domain generalization with mixstyle. In *ICLR*, 2021.

## A Theoretical Analysis and Complete Proofs

In this section, we explain the details of Theorem 1 in the main paper, and also formally describe Theorem 2. We start with giving additional definitions and providing a useful lemma and its proof, which invoked through the proof of the theorems. We then formally prove the arguments in Theorem 1 and 2.

### A.1 Additional Definition, Lemma, and Theorem

**Definition 4** ( $(\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma})$  Calibrated Transferability Statistics). *The transferability graph can be further described by the following three components:*

$$\begin{aligned}\tilde{\alpha} &= \mathbb{E}_c \mathbb{E}_d \mathbb{E}_{d' \neq d} \left[ \lambda_{d,c}^{d',c} \cdot \text{trans}((d, c), (d', c)) \right], \\ \tilde{\beta} &= \mathbb{E}_d \mathbb{E}_c \mathbb{E}_{c' \neq c} \left[ \lambda_{d,c}^{d,c'} \cdot \text{trans}((d, c), (d, c')) \right], \\ \tilde{\gamma} &= \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_c \mathbb{E}_{c' \neq c} \left[ \lambda_{d,c}^{d',c'} \cdot \text{trans}((d, c), (d', c')) \right],\end{aligned}$$

where  $\lambda_{d,c}^{d',c'} = \left( \frac{N_{d',c'}}{N_{d,c}} \right)^\nu$  denotes the distance calibration coefficient.

**Lemma 1.** *Let  $\eta, \pi > 0$  and  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ ,  $\varphi(x) = \log(\eta + \pi \exp(x))$ . Given a finite sequence  $x_1, x_2, \dots, x_M \in \mathbb{R}$ , it holds that*

$$\frac{1}{M} \sum_{i=1}^M \varphi(x_i) \geq \varphi \left( \frac{1}{M} \sum_{i=1}^M x_i \right).$$

*Proof.* Note that  $\varphi$  is smooth and thus twice differentiable for all  $x \in \mathbb{R}$ . We obtain the second derivative of  $\varphi$  as

$$\varphi''(x) = \frac{\eta \pi \exp(x)}{(\eta + \pi \exp(x))^2} > 0, \quad \forall x \in \mathbb{R}.$$

Therefore,  $\varphi$  is convex. Thus, by Jensen's inequality, we obtain that  $\frac{1}{M} \sum_{i=1}^M \varphi(x_i) \geq \varphi \left( \frac{1}{M} \sum_{i=1}^M x_i \right)$ , which completes the proof.  $\square$

**Theorem 2** ( $\tilde{\mathcal{L}}_{\text{BoDA}}$  as an Upper Bound). *Given a multi-domain long-tailed dataset  $\mathcal{S}$  with domain label space  $\mathcal{D}$  and class label space  $\mathcal{C}$  satisfying  $|\mathcal{D}| > 1$  and  $|\mathcal{C}| > 1$ , let  $\mathcal{Z}$  be the representation set of all training samples. It holds that*

$$\tilde{\mathcal{L}}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) \geq N \log \left( |\mathcal{D}| - 1 + |\mathcal{D}|(|\mathcal{C}| - 1) \exp \left( \frac{|\mathcal{C}||\mathcal{D}|}{N} \cdot \tilde{\alpha} - \frac{|\mathcal{C}|}{N} \cdot \tilde{\beta} - \frac{|\mathcal{C}|(|\mathcal{D}| - 1)}{N} \cdot \tilde{\gamma} \right) \right), \quad (5)$$

where  $(\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma})$  are the calibrated transferability statistics for  $\mathcal{S}$  defined in Definition 4.

## A.2 Proof of Theorem 1

Recall that  $\mathcal{M} = \mathcal{D} \times \mathcal{C} := \{(d, c) : d \in \mathcal{D}, c \in \mathcal{C}\}$  is the set of all domain-class pairs.  $\mathcal{L}_{\text{BoDA}}$  is given by

$$\begin{aligned} \mathcal{L}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) &= \sum_{\mathbf{z}_i \in \mathcal{Z}} \frac{-1}{|\mathcal{D}| - 1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i}))}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))} \\ &= \sum_{\mathbf{z}_i \in \mathcal{Z}} \ell_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\}), \end{aligned}$$

where  $\ell_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\})$  is the *sample-wise* BoDA loss. We rewrite  $\ell_{\text{BoDA}}$  in the following format

$$\begin{aligned} \ell_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\}) &= -\frac{1}{|\mathcal{D}| - 1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i}))}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))} \\ &= \log \left( \frac{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))}{\prod_{d \in \mathcal{D} \setminus \{d_i\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i}))^{|\mathcal{D}| - 1}} \right) \\ &= \log \left( \frac{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))}{\exp\left(-\frac{1}{|\mathcal{D}| - 1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i})\right)} \right). \end{aligned} \quad (6)$$

We will first focus on the term in the numerator of Eqn. (6). We can rewrite the sum into two terms

$$\begin{aligned} &\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})) \\ &= \underbrace{\sum_{d' \in \mathcal{D} \setminus \{d_i\}} \sum_{c' \in \{c_i\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))}_{T_1} + \underbrace{\sum_{d' \in \mathcal{D}} \sum_{c' \in \mathcal{C} \setminus \{c_i\}} \exp(-\tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}))}_{T_2}. \end{aligned}$$

Since the exponential function  $\exp(\cdot)$  is convex, we apply Jensen's inequality on both  $T_1$  and  $T_2$ :

$$\begin{aligned} T_1 &\geq (|\mathcal{D}| - 1) \exp \left( -\frac{1}{|\mathcal{D}| - 1} \sum_{d' \in \mathcal{D} \setminus \{d_i\}} \sum_{c' \in \{c_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}) \right) \\ &= (|\mathcal{D}| - 1) \exp \left( -\frac{1}{|\mathcal{D}| - 1} \sum_{d' \in \mathcal{D} \setminus \{d_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c_i}) \right), \\ T_2 &\geq |\mathcal{D}|(|\mathcal{C}| - 1) \exp \left( -\frac{1}{|\mathcal{D}|(|\mathcal{C}| - 1)} \sum_{d' \in \mathcal{D}} \sum_{c' \in \mathcal{C} \setminus \{c_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}) \right). \end{aligned}$$

Thus, by using  $\exp(x)/\exp(y) = \exp(x - y)$  and rearranging terms, we bound  $\ell_{\text{BoDA}}$  by

$$\begin{aligned} & \ell_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\}) \\ & \geq \log \left( |\mathcal{D}| - 1 + |\mathcal{D}|(|\mathcal{C}| - 1) \exp \left( \underbrace{\frac{1}{|\mathcal{D}| - 1} \sum_{d' \in \mathcal{D} \setminus \{d_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c_i}) - \frac{1}{|\mathcal{D}|(|\mathcal{C}| - 1)} \sum_{d' \in \mathcal{D}} \sum_{c' \in \mathcal{C} \setminus \{c_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})}_{T(\mathbf{z}_i, \{\boldsymbol{\mu}\})} \right) \right). \end{aligned}$$

Leveraging Lemma 1, by setting  $\eta = |\mathcal{D}| - 1$ ,  $\pi = |\mathcal{D}|(|\mathcal{C}| - 1)$ , and  $x_i = T(\mathbf{z}_i, \{\boldsymbol{\mu}\})$ , we further bound  $\mathcal{L}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\})$  by

$$\begin{aligned} \mathcal{L}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) &= \sum_{\mathbf{z}_i \in \mathcal{Z}} \ell_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\}) \\ &\geq \sum_{\mathbf{z}_i \in \mathcal{Z}} \log \left( (|\mathcal{D}| - 1 + |\mathcal{D}|(|\mathcal{C}| - 1) \exp(T(\mathbf{z}_i, \{\boldsymbol{\mu}\}))) \right) \\ &\geq |\mathcal{Z}| \log \left( |\mathcal{D}| - 1 + |\mathcal{D}|(|\mathcal{C}| - 1) \exp \left( \frac{1}{|\mathcal{Z}|} \sum_{\mathbf{z}_i \in \mathcal{Z}} T(\mathbf{z}_i, \{\boldsymbol{\mu}\}) \right) \right). \end{aligned} \quad (7)$$

Note that the argument of the  $\exp(\cdot)$  in Eqn. (7) can be expanded and further rearranged as

$$\begin{aligned} \frac{1}{|\mathcal{Z}|} \sum_{\mathbf{z}_i \in \mathcal{Z}} T(\mathbf{z}_i, \{\boldsymbol{\mu}\}) &= \frac{1}{|\mathcal{Z}|} \sum_{\mathbf{z}_i \in \mathcal{Z}} \frac{1}{|\mathcal{D}| - 1} \sum_{d' \in \mathcal{D} \setminus \{d_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c_i}) - \\ & \quad \frac{1}{|\mathcal{Z}|} \sum_{\mathbf{z}_i \in \mathcal{Z}} \frac{1}{|\mathcal{D}|(|\mathcal{C}| - 1)} \sum_{d' \in \mathcal{D}} \sum_{c' \in \mathcal{C} \setminus \{c_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'}) \\ &= \frac{1}{|\mathcal{Z}|} \underbrace{\frac{1}{|\mathcal{D}| - 1} \sum_{\mathbf{z}_i \in \mathcal{Z}} \sum_{d' \in \mathcal{D} \setminus \{d_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c_i})}_{T_\alpha} - \\ & \quad \frac{1}{|\mathcal{Z}|} \underbrace{\frac{1}{|\mathcal{D}|(|\mathcal{C}| - 1)} \sum_{\mathbf{z}_i \in \mathcal{Z}} \sum_{c' \in \mathcal{C} \setminus \{c_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d_i, c'})}_{T_\beta} - \\ & \quad \underbrace{\frac{1}{|\mathcal{Z}|} \frac{1}{|\mathcal{D}|(|\mathcal{C}| - 1)} \sum_{\mathbf{z}_i \in \mathcal{Z}} \sum_{d' \in \mathcal{D} \setminus \{d_i\}} \sum_{c' \in \mathcal{C} \setminus \{c_i\}} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})}_{T_\gamma}. \end{aligned} \quad (8)$$

Recall that each  $\mathbf{z}_i \in \mathcal{Z}$  belongs to a domain-class pair  $(d_i, c_i)$ , and  $\mathcal{Z}_{d,c}$  denotes the representation set of  $\mathcal{S}_{d,c}$  with size  $N_{d,c}$ . For simplicity, we remove the subscript  $i$  in the following derivation. We

can further rewrite  $T_\alpha, T_\beta, T_\gamma$  as

$$\begin{aligned}
T_\alpha &= \frac{1}{|\mathcal{Z}|} \frac{1}{|\mathcal{D}| - 1} \sum_{c \in \mathcal{C}} \sum_{d \in \mathcal{D}} \sum_{d' \in \mathcal{D} \setminus \{d\}} \sum_{\mathbf{z} \in \mathcal{Z}_{d,c}} \tilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c}) \\
&= \frac{1}{|\mathcal{Z}|} \frac{1}{|\mathcal{D}| - 1} |\mathcal{C}| |\mathcal{D}| (|\mathcal{D}| - 1) \mathbb{E}_c \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} \underbrace{[N_{d,c} \cdot \tilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c})]}_{\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c})} \\
&= \frac{|\mathcal{C}| |\mathcal{D}|}{|\mathcal{Z}|} \underbrace{\mathbb{E}_c \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c})]}_{\alpha}, \tag{9}
\end{aligned}$$

$$\begin{aligned}
T_\beta &= \frac{1}{|\mathcal{Z}|} \frac{1}{|\mathcal{D}| (|\mathcal{C}| - 1)} \sum_{c \in \mathcal{C}} \sum_{d \in \mathcal{D}} \sum_{c' \in \mathcal{C} \setminus \{c\}} \sum_{\mathbf{z} \in \mathcal{Z}_{d,c}} \tilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d,c'}) \\
&= \frac{1}{|\mathcal{Z}|} \frac{1}{|\mathcal{D}| (|\mathcal{C}| - 1)} |\mathcal{C}| |\mathcal{D}| (|\mathcal{C}| - 1) \mathbb{E}_d \mathbb{E}_c \mathbb{E}_{c' \neq c} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} \underbrace{[N_{d,c} \cdot \tilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d,c'})]}_{\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d,c'})} \\
&= \frac{|\mathcal{C}|}{|\mathcal{Z}|} \underbrace{\mathbb{E}_d \mathbb{E}_c \mathbb{E}_{c' \neq c} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d,c'})]}_{\beta}, \tag{10}
\end{aligned}$$

$$\begin{aligned}
T_\gamma &= \frac{1}{|\mathcal{Z}|} \frac{1}{|\mathcal{D}| (|\mathcal{C}| - 1)} \sum_{c \in \mathcal{C}} \sum_{d \in \mathcal{D}} \sum_{d' \in \mathcal{D} \setminus \{d\}} \sum_{c' \in \mathcal{C} \setminus \{c\}} \sum_{\mathbf{z} \in \mathcal{Z}_{d,c}} \tilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c'}) \\
&= \frac{1}{|\mathcal{Z}|} \frac{|\mathcal{C}| |\mathcal{D}| (|\mathcal{D}| - 1) (|\mathcal{C}| - 1)}{|\mathcal{D}| (|\mathcal{C}| - 1)} \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_c \mathbb{E}_{c' \neq c} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} \underbrace{[N_{d,c} \cdot \tilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c'})]}_{\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c'})} \\
&= \frac{|\mathcal{C}| (|\mathcal{D}| - 1)}{|\mathcal{Z}|} \underbrace{\mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_c \mathbb{E}_{c' \neq c} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c'})]}_{\gamma}, \tag{11}
\end{aligned}$$

where  $(\alpha, \beta, \gamma)$  are the transferability statistics for  $\mathcal{S}$  as in Definition 3. Finally, replace  $|\mathcal{Z}| = N$  and combine Eqn. (7), (8), (9), (10), and (11), we have

$$\mathcal{L}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) \geq N \log \left( |\mathcal{D}| - 1 + |\mathcal{D}| (|\mathcal{C}| - 1) \exp \left( \frac{|\mathcal{C}| |\mathcal{D}|}{N} \cdot \alpha - \frac{|\mathcal{C}|}{N} \cdot \beta - \frac{|\mathcal{C}| (|\mathcal{D}| - 1)}{N} \cdot \gamma \right) \right).$$

This completes the proof.

### A.3 Proof of Theorem 2

We first define a notion of *calibrated distance*  $\hat{\mathbf{d}}$ . Let  $\mathbf{z} \in \mathcal{Z}_{d,c}$ , we have

$$\hat{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c'}) \triangleq \lambda_{d,c}^{d',c'} \cdot \tilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c'}) = \left( \frac{N_{d',c'}}{N_{d,c}} \right)^\nu \cdot \tilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c'}).$$

From Theorem 1, by substituting  $\tilde{\mathbf{d}}$  with  $\hat{\mathbf{d}}$ , it holds that

$$\begin{aligned}
\tilde{\mathcal{L}}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) &= \mathcal{L}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) \Big|_{\tilde{\mathbf{d}} \rightarrow \hat{\mathbf{d}}} \\
&\geq N \log \left( |\mathcal{D}| - 1 + |\mathcal{D}| (|\mathcal{C}| - 1) \exp (T'_\alpha - T'_\beta - T'_\gamma) \right), \tag{12}
\end{aligned}$$

where  $T'_\alpha$ ,  $T'_\beta$ , and  $T'_\gamma$  can be expressed as

$$\begin{aligned}
T'_\alpha &= \frac{|\mathcal{C}||\mathcal{D}|}{N} \mathbb{E}_c \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [N_{d,c} \cdot \widehat{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c})] \\
&= \frac{|\mathcal{C}||\mathcal{D}|}{N} \mathbb{E}_c \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\lambda_{d,c}^{d',c} \cdot \underbrace{N_{d,c} \cdot \widetilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c})}_{\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c})}] \\
&= \frac{|\mathcal{C}||\mathcal{D}|}{N} \underbrace{\mathbb{E}_c \mathbb{E}_d \mathbb{E}_{d' \neq d} \left[ \lambda_{d,c}^{d',c} \cdot \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c})] \right]}_{\widetilde{\alpha}}, \tag{13}
\end{aligned}$$

$$\begin{aligned}
T'_\beta &= \frac{|\mathcal{C}|}{N} \mathbb{E}_d \mathbb{E}_c \mathbb{E}_{c' \neq c} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [N_{d,c} \cdot \widehat{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d,c'})] \\
&= \frac{|\mathcal{C}|}{N} \mathbb{E}_d \mathbb{E}_c \mathbb{E}_{c' \neq c} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\lambda_{d,c}^{d,c'} \cdot \underbrace{N_{d,c} \cdot \widetilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d,c'})}_{\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d,c'})}] \\
&= \frac{|\mathcal{C}|}{N} \underbrace{\mathbb{E}_d \mathbb{E}_c \mathbb{E}_{c' \neq c} \left[ \lambda_{d,c}^{d,c'} \cdot \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d,c'})] \right]}_{\widetilde{\beta}}, \tag{14}
\end{aligned}$$

$$\begin{aligned}
T'_\gamma &= \frac{|\mathcal{C}|(|\mathcal{D}| - 1)}{N} \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_c \mathbb{E}_{c' \neq c} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [N_{d,c} \cdot \widehat{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c'})] \\
&= \frac{|\mathcal{C}|(|\mathcal{D}| - 1)}{N} \mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_c \mathbb{E}_{c' \neq c} \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\lambda_{d,c}^{d',c'} \cdot \underbrace{N_{d,c} \cdot \widetilde{\mathbf{d}}(\mathbf{z}, \boldsymbol{\mu}_{d',c'})}_{\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c'})}] \\
&= \frac{|\mathcal{C}|(|\mathcal{D}| - 1)}{N} \underbrace{\mathbb{E}_d \mathbb{E}_{d' \neq d} \mathbb{E}_c \mathbb{E}_{c' \neq c} \left[ \lambda_{d,c}^{d',c'} \cdot \mathbb{E}_{\mathbf{z} \in \mathcal{Z}_{d,c}} [\mathbf{d}(\mathbf{z}, \boldsymbol{\mu}_{d',c'})] \right]}_{\widetilde{\gamma}}, \tag{15}
\end{aligned}$$

where  $(\widetilde{\alpha}, \widetilde{\beta}, \widetilde{\gamma})$  are formally defined in Definition 4. Combine Eqn. (12), (13), (14), and (15), we have

$$\widetilde{\mathcal{L}}_{\text{BoDA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) \geq N \log \left( |\mathcal{D}| - 1 + |\mathcal{D}|(|\mathcal{C}| - 1) \exp \left( \frac{|\mathcal{C}||\mathcal{D}|}{N} \cdot \widetilde{\alpha} - \frac{|\mathcal{C}|}{N} \cdot \widetilde{\beta} - \frac{|\mathcal{C}|(|\mathcal{D}| - 1)}{N} \cdot \widetilde{\gamma} \right) \right),$$

which completes the proof.

## B Additional Discussions, Properties, and Interpretations

### B.1 Unified Interpretation for Single- and Multi-Domain Imbalance

In the main paper we show that, in the multi-domain setting, label imbalance implicitly brings *label divergence* across domains, which brings additional challenges and potentially harms MDLT performance. Here we provide a unified viewpoint from the *label divergence* perspective to explain single- and multi-domain data imbalance.

To elaborate, in single domain imbalanced learning, we essentially cope with the divergence between the imbalanced training label distribution and the uniform test label distribution:

$$\text{div}(p(y) \parallel \mathcal{U}),$$

where  $\text{div}(\cdot\|\cdot)$  indicates certain divergence measure. In contrast, when extending to the multi-domain scenario, given  $|\mathcal{D}|$  domains with (different) imbalanced label distributions, the target divergence becomes

$$\underbrace{\sum_d \text{div}(p_d(y) \|\mathcal{U})}_{\text{imbalanced training}} + \text{const} \cdot \underbrace{\sum_{d \neq d'} \text{div}(p_d(y) \|\ p_{d'}(y))}_{\text{divergence across domains}},$$

where one not only needs to tackle the imbalanced training data for each domain  $d \in \mathcal{D}$  in order to generalize to the balanced test set, but also takes into consideration the *label divergence* across domains.

Such interpretation echoes our BoDA objective: We design the DA loss for cross-domain distribution alignment to tackle the latter term, and further adapt it to BoDA via balanced distance to address the former term.

## B.2 A Probabilistic Perspective of $\mathcal{L}_{\text{DA}}$ Derivation

Recall  $\mathcal{M} = \mathcal{D} \times \mathcal{C}$  the set of all  $(d, c)$  pairs. Let  $(\mathbf{x}_i, c_i, d_i)$  denote a sample with feature  $\mathbf{z}_i$ . Following the metric learning setting [17], we model the likelihood of  $\boldsymbol{\mu}_{d,c}$  given  $\mathbf{z}_i$  to decay exponentially with respect to their distance in the representation space. Such modeling can be viewed as performing a random walk with transition probability inversely related to distance [16]. For domain-class pairs that share the same class label but different domain labels with  $\mathbf{x}_i$  (i.e.,  $(d, c_i), d \neq d_i$ ), the normalized likelihood of  $\boldsymbol{\mu}_{d,c_i}$  given  $\mathbf{z}_i$  can be written as

$$\mathbb{P}((d, c_i) | \mathbf{z}_i) = \frac{\exp(-\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d,c_i}))}{\sum_{(d',c') \in \mathcal{M} \setminus \{(d_i,c_i)\}} \exp(-\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d',c'}))},$$

where the denominator is a sum over all domain-class pairs except  $(d_i, c_i)$ . As motivated, we want to concentrate all  $\mathbf{z}_i$  from the same class across different domains (i.e., smaller  $\alpha$ ), while separating  $\mathbf{z}_i$  from different classes within and across domains (i.e., larger  $\beta, \gamma$ ). Therefore, the positive domain-class pairs with  $\mathbf{x}_i$  are those share the same class labels but different domain labels. As a result, we define the per-sample loss as the average negative log-likelihood over all positive domain-class pairs:

$$\ell_{\text{DA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\}) = -\frac{1}{|\mathcal{D}| - 1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp(-\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d,c_i}))}{\sum_{(d',c') \in \mathcal{M} \setminus \{(d_i,c_i)\}} \exp(-\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d',c'}))}.$$

Given a set of all training samples with representation set as  $\mathcal{Z}$ , the total loss can then be derived as

$$\mathcal{L}_{\text{DA}}(\mathcal{Z}, \{\boldsymbol{\mu}\}) = \sum_{\mathbf{z}_i \in \mathcal{Z}} \frac{-1}{|\mathcal{D}| - 1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp(-\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d,c_i}))}{\sum_{(d',c') \in \mathcal{M} \setminus \{(d_i,c_i)\}} \exp(-\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d',c'}))}.$$

## B.3 Intrinsic Hardness-Aware Property of BoDA

Below, we demonstrate an additional property of BoDA: the intrinsic *hardness-aware* property. Specifically, we analyze the gradients of BoDA loss with respect to positive  $(d, c)$  pairs and different negative  $(d, c)$  pairs. We observe that the gradient contributions from *hard* positives/negatives are

larger than that from the *easy* ones, indicating that BoDA automatically concentrates on the *hard*  $(d, c)$  pairs, where penalties are given according to their hardness.

Recall that the sample-wise calibrated BoDA loss  $\tilde{\ell}_{\text{BoDA}}$  can be written as

$$\begin{aligned}
& \tilde{\ell}_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\}) \\
&= -\frac{1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp\left(-\lambda_{d_i, c_i}^{d, c_i} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i})\right)}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp\left(-\lambda_{d_i, c_i}^{d', c'} \tilde{\mathbf{d}}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})\right)} \\
&= -\frac{1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \log \frac{\exp\left(-\frac{\lambda_{d_i, c_i}^{d, c_i}}{N_{d_i, c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i})\right)}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp\left(-\frac{\lambda_{d_i, c_i}^{d', c'}}{N_{d_i, c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})\right)}, \tag{16}
\end{aligned}$$

where  $\mathbf{z}_i \in \mathcal{Z}_{d_i, c_i}$ . For convenience, we further define the probability of  $\mathbf{z}_i$  being recognized as belonging to  $\boldsymbol{\mu}_{d, c}$  as

$$P_{d, c}^i \triangleq \frac{\exp\left(-\frac{\lambda_{d_i, c_i}^{d, c}}{N_{d_i, c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c})\right)}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp\left(-\frac{\lambda_{d_i, c_i}^{d', c'}}{N_{d_i, c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})\right)}, \quad (d, c) \in \mathcal{M} \setminus \{(d_i, c_i)\}.$$

Note that the essential goal of Eqn. (16) is to align (minimize) *positive* distances  $\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i})$  and to separate (maximize) *negative* distances  $\mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})$ . Therefore, we analyze the gradients with respect to positive distance and different negative distances to explore the properties of  $\tilde{\ell}_{\text{BoDA}}$ . Specifically, we have

$$\begin{aligned}
& \frac{\partial \tilde{\ell}_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\})}{\partial \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i})} \\
&= \frac{-1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \frac{\partial}{\partial \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i})} \left\{ -\frac{\lambda_{d_i, c_i}^{d, c_i}}{N_{d_i, c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i}) - \log \sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp\left(-\frac{\lambda_{d_i, c_i}^{d', c'}}{N_{d_i, c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})\right) \right\} \\
&= \frac{1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \frac{\lambda_{d_i, c_i}^{d, c_i}}{N_{d_i, c_i}} \left( 1 - \frac{\exp\left(-\frac{\lambda_{d_i, c_i}^{d, c_i}}{N_{d_i, c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d, c_i})\right)}{\sum_{(d', c') \in \mathcal{M} \setminus \{(d_i, c_i)\}} \exp\left(-\frac{\lambda_{d_i, c_i}^{d', c'}}{N_{d_i, c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d', c'})\right)} \right) \\
&= \frac{1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \frac{N_{d, c_i}^\nu}{N_{d_i, c_i}^{(1+\nu)}} (1 - P_{d, c_i}^i) \\
&\propto \sum_{d \in \mathcal{D} \setminus \{d_i\}} N_{d, c_i}^\nu (1 - P_{d, c_i}^i),
\end{aligned}$$

$$\begin{aligned}
& \frac{\partial \tilde{\ell}_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\})}{\partial \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d',c'})} \\
&= \frac{-1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \frac{\partial}{\partial \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d',c'})} \left\{ -\frac{\lambda_{d_i,c_i}^{d,c_i}}{N_{d_i,c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d,c_i}) - \log \sum_{(d',c') \in \mathcal{M} \setminus \{(d_i,c_i)\}} \exp \left( -\frac{\lambda_{d_i,c_i}^{d',c'}}{N_{d_i,c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d',c'}) \right) \right\} \\
&= -\frac{1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \frac{\lambda_{d_i,c_i}^{d',c'}}{N_{d_i,c_i}} \frac{\exp \left( -\frac{\lambda_{d_i,c_i}^{d',c'}}{N_{d_i,c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d,c_i}) \right)}{\sum_{(d',c') \in \mathcal{M} \setminus \{(d_i,c_i)\}} \exp \left( -\frac{\lambda_{d_i,c_i}^{d',c'}}{N_{d_i,c_i}} \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d',c'}) \right)} \\
&= -\frac{1}{|\mathcal{D}|-1} \sum_{d \in \mathcal{D} \setminus \{d_i\}} \frac{N_{d',c'}^\nu}{N_{d_i,c_i}^{(1+\nu)}} P_{d',c'}^i \\
&\propto -N_{d',c'}^\nu P_{d',c'}^i.
\end{aligned}$$

Combine the above results, we have

$$\text{positive: } \frac{\partial \tilde{\ell}_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\})}{\partial \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d,c_i})} \propto \sum_{d \in \mathcal{D} \setminus \{d_i\}} N_{d,c_i}^\nu (1 - P_{d,c_i}^i), \quad (17)$$

$$\text{negative: } \frac{\partial \tilde{\ell}_{\text{BoDA}}(\mathbf{z}_i, \{\boldsymbol{\mu}\})}{\partial \mathbf{d}(\mathbf{z}_i, \boldsymbol{\mu}_{d',c'})} \propto -N_{d',c'}^\nu P_{d',c'}^i. \quad (18)$$

**Interpretation.** Eqn. (17) and (18) illustrate several interesting and important properties of BoDA:

1. *Intrinsic hard positive and negative mining.* For positive pairs, we observe that the gradient magnitudes are proportional to  $(1 - P_{d,c_i}^i)$ , where for an easy  $(d, c_i)$  pair,  $P_{d,c_i}^i \approx 1$  and  $(1 - P_{d,c_i}^i) \approx 0$ , and for a hard  $(d, c_i)$  pair,  $P_{d,c_i}^i \approx 0$  and  $(1 - P_{d,c_i}^i) \approx 1$ , indicating that the gradient contributions from *hard* positives are larger than *easy* ones. Similarly, for negative pairs, the gradient magnitudes are proportional to  $P_{d',c'}^i$ , where an easy  $(d', c')$  pair has  $P_{d',c'}^i \approx 0$  and a hard  $(d, c_i)$  pair induces  $P_{d',c'}^i \approx 1$ , showing that the gradient contribution is large for hard negatives and small for easy negatives. Therefore, BoDA is a hardness-aware loss with intrinsic hard positive/negative mining property.
2. *Scaling gradients according to the number of samples of each  $(d, c)$ .* Furthermore, as we have shown in Fig. 5, when data are imbalanced across different  $(d, c)$  pairs, minority pairs with smaller number of samples would induce worse  $\boldsymbol{\mu}_{d,c}$  estimates. We further observe that the gradients for both positive and negative pairs are proportional to their number of samples (i.e.,  $N_{d,c_i}^\nu$  and  $N_{d',c'}^\nu$ ). This suggests that BoDA automatically adjusts the gradient scale for each  $(d, c)$  according to how accurate the estimation of  $\boldsymbol{\mu}_{d,c}$  is. The appealing property highlights that BoDA also implicitly calibrates the gradient scale, emphasizing gradients from majority pairs (which are more reliable) while suppressing gradients from minority pairs (which are less reliable). Such behavior is essential for better statistics transfer as we demonstrated in the main paper.

## C Pseudo Code for BoDA

We provide the pseudo code of BoDA in Algorithm 1.

---

**Algorithm 1** Balanced Domain-Class Distribution Alignment (BoDA)

---

**Input:** Training set  $\mathcal{D} = \{(\mathbf{x}_i, c_i, d_i)\}_{i=1}^N$ , all domain-class pairs  $\mathcal{M} = \{(d, c)\}$ , encoder  $f$ , classifier  $g$ , total training epochs  $E$ , calibration parameter  $\nu$ , loss weight  $\omega$ , momentum  $\alpha$

**for all**  $(d, c) \in \mathcal{M}$  **do**

    Initialize the feature statistics  $\{\boldsymbol{\mu}_{d,c}^{(0)}, \boldsymbol{\Sigma}_{d,c}^{(0)}\}$

**end for**

**for**  $e = 0$  **to**  $E$  **do**

**repeat**

        Sample a mini-batch  $\{(\mathbf{x}_i, c_i, d_i)\}_{i=1}^m$  from  $\mathcal{D}$

**for**  $i = 1$  **to**  $m$  (in parallel) **do**

$\mathbf{z}_i = f(\mathbf{x}_i)$

$\hat{c}_i = g(\mathbf{z}_i)$

**end for**

        Calculate  $\tilde{\mathcal{L}}_{\text{BoDA}}$  using  $\{\mathbf{z}_i\}$  based on Eqn. (4)

        Calculate  $\mathcal{L}_{\text{CE}}$  using  $\frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{c}_i, c_i)$

        Do one training step with loss  $\mathcal{L}_{\text{CE}} + \omega \tilde{\mathcal{L}}_{\text{BoDA}}$

**until** iterate over all training samples at current epoch  $e$

    /\* Update feature statistics with momentum updating \*/

**for all**  $(d, c) \in \mathcal{M}$  **do**

        Estimate current feature statistics  $\{\boldsymbol{\mu}_{d,c}, \boldsymbol{\Sigma}_{d,c}\}$

$\boldsymbol{\mu}_{d,c}^{(e+1)} \leftarrow \alpha \times \boldsymbol{\mu}_{d,c}^{(e)} + (1 - \alpha) \times \boldsymbol{\mu}_{d,c}$

$\boldsymbol{\Sigma}_{d,c}^{(e+1)} \leftarrow \alpha \times \boldsymbol{\Sigma}_{d,c}^{(e)} + (1 - \alpha) \times \boldsymbol{\Sigma}_{d,c}$

**end for**

**end for**

---

## D Details of MDLT Datasets

In this section, we provide the detailed information of the curated MDLT datasets we used in our experiments. Table 10 provides an overview of the datasets. Table 11 provides the image examples across domains for each MDLT dataset.

**Digits-MLT.** We construct Digits-MLT by combining two digit datasets: (1) MNIST-M [15], a variant of the original MNIST handwritten digit classification dataset [26] with colorful background, and (2) SVHN [36]. The original MNIST-M dataset contains 60,000 training samples and 10,000 testing examples, and the original SVHN dataset contains 73,257 images for training and 26,032 images for testing. Both datasets have examples of dimension  $(3, 32, 32)$  and 10 classes. We create Digits-MLT with controllable degrees of data imbalance, where we keep the maximum number of samples each  $(d, c)$  to be 1,000, and manually vary the imbalance degree to adjust the number of samples for minority  $(d, c)$ . For validation and test set, we use the original test set of the two datasets, but keep the number of samples each  $(d, c)$  to be 800.

**VLCS-MLT.** The original VLCS dataset [14] is an object recognition dataset that comprises photographic domains  $d \in \{\text{Caltech101, LabelMe, SUN09, VOC2007}\}$ , with scenes captured from urban to rural. The dataset contains 5 classes with 10,729 examples of dimension  $(3, 224, 224)$ . To construct VLCS-MLT, for each  $(d, c)$  we split out a validation set of size 15 and a test set of size 30, and

Table 10: Detailed statistics of the curated MDLT datasets used in our experiments. For the synthetic **Digits-MLT** dataset, we manually vary the minimum ( $d, c$ ) size to simulate different degrees of imbalance.

Dataset	# Domains	# Classes	Max ( $d, c$ ) size	Min ( $d, c$ ) size	# Training set	# Val. set	# Test set
Digits-MLT	2	10	1,000	10 ~ 1,000	20,000 ~ 4,956	16,000	16,000
VLCS-MLT	4	5	1,454	0	9,872	285	572
PACS-MLT	4	7	741	5	7,891	700	1,400
OfficeHome-MLT	4	65	84	0	11,688	1,300	2,600
TerraInc-MLT	4	10	4,455	0	23,269	353	708
DomainNet-MLT	6	345	778	0	468,574	39,240	78,761

Table 11: Overview of images from different domains in all MDLT datasets. For each dataset, we pick a single class and show illustrative images from each domain.

Dataset	Domains					
Digits-MLT	MNIST-M	SVHN				
						
VLCS-MLT	Caltech101	LabelMe	SUN09	VOC2007		
						
PACS-MLT	Art	Cartoon	Photo	Sketch		
						
OfficeHome-MLT	Art	Clipart	Product	Photo		
						
TerraInc-MLT	L100	L38	L43	L46		
						
DomainNet-MLT	Clipart	Infographic	Painting	QuickDraw	Photo	Sketch
						

leave the rest for training.

**PACS-MLT.** The original PACS dataset [27] is an object recognition dataset that comprises four domains  $d \in \{ \text{art, cartoons, photos, sketches} \}$  with image style changes. It contains 7 classes with 9,991 examples of dimension (3, 224, 224). We construct PACS-MLT in a similar manner as VLCS-MLT,

where we split out a validation set of size 25 and a test set of size 50 for each  $(d, c)$ , and leave the rest for training.

**OfficeHome-MLT.** The original OfficeHome dataset [47] includes domains  $d \in \{ \text{art, clipart, product, real} \}$ , containing 15,588 examples of dimension  $(3, 224, 224)$  and 65 classes. We make OfficeHome-MLT by splitting out a validation set of size 5 and a test set of size 10 for each  $(d, c)$ , leaving the rest for training.

**TerraInc-MLT.** TerraInc-MLT is constructed from TerraIncognita dataset [2], a species classification dataset that contains photographs of wild animals taken by camera traps at locations  $d \in \{ \text{L100, L38, L43, L46} \}$ . The dataset contains 10 classes with 24,788 examples of dimension  $(3, 224, 224)$ . For each  $(d, c)$ , we split out a validation set of size 10 and a test set of size 20, and use all remaining samples for training.

**DomainNet-MLT.** We construct DomainNet-MLT using DomainNet dataset [38], a large-scale multi-domain dataset for object recognition that consists of six domains  $d \in \{ \text{clipart, infograph, painting, quickdraw, real, sketch} \}$ , 345 classes, and 586,575 examples of size  $(3, 224, 224)$ . To construct DomainNet-MLT, for each  $(d, c)$  we split out a validation set of size 20 and a test set of size 40, and leave the rest for training.

## E Experimental Settings

### E.1 Implementation Details

For the synthetic Digits-MLT dataset, we fix the network architecture as a small MNIST CNN [19] for all algorithms, and use no data augmentation. For all other MDLT datasets, following [19], we use the pretrained ResNet-50 model [21] as the backbone network for all algorithms, and use the same data augmentation protocol as [19]: random crop and resize to  $224 \times 224$  pixels, random horizontal flips, random color jitter, grayscaling the image with 10% probability, and normalization using the ImageNet channel statistics. We train all models using the Adam optimizer [24] for 5,000 steps on all MDLT datasets except DomainNet-MLT, on which we train longer for 15,000 steps to ensure convergence. We fix a batch size of 64 per domain for Digits-MLT experiments, a batch size of 32 per domain for DomainNet-MLT experiments, and a batch size of 24 per domain for experiments on all other datasets.

For all MDLT datasets except OfficeHome-MLT and TerraInc-MLT, we define *many-shot*  $(d, c)$  pairs as with over 100 training samples, *medium-shot* as with 20~100 training samples, and *few-shot* as with under 20 training samples. For OfficeHome-MLT, we define *many-shot* as  $(d, c)$  pairs with over 60 training samples, *medium-shot* as with 20~60 training samples, and *few-shot* as with under 20 training samples. For TerraInc-MLT, we define *many-shot* as  $(d, c)$  pairs with over 100 training samples, *medium-shot* as with 25~100 training samples, and *few-shot* as with under 25 training samples.

### E.2 Competing Algorithms

We compare BoDA to a large number of algorithms that span different learning strategies. We group them according to their categories, and provide detailed descriptions for each algorithm below.

- *Vanilla*: The empirical risk minimization (**ERM**) [46] minimizes the sum of errors across all domains and samples.
- *Distributionally robust optimization*: Group distributionally robust optimization (**GroupDRO**) [40] performs ERM while increasing the importance of domains with larger errors.
- *Cross-domain data augmentation*: Inter-domain mixup (**Mixup**) [50] performs ERM on linear interpolations of examples from random pairs of domains and their labels. Style-agnostic network (**SagNet**) [35] disentangles style encodings from image content by randomizing and augmenting styles.
- *Meta-learning*: Meta-learning for domain generalization (**MLDG**) [28] leverages meta-learning to learn how to generalize across domains.
- *Domain-invariant representation learning*: Invariant risk minimization (**IRM**) [1] learns a feature representation such that the optimal linear classifier on top of that representation matches across domains. Domain adversarial neural networks (**DANN**) [15] employ an adversarial network to match feature distributions. Class-conditional DANN (**CDANN**) [31] builds upon DANN but further matches the conditional distributions across domains for all labels. Deep correlation alignment (**CORAL**) [45] matches the mean and covariance of feature distributions. Maximum mean discrepancy (**MMD**) [29] matches the MMD [18] of feature distributions.
- *Transfer learning*: Marginal transfer learning (**MTL**) [4] estimates a mean embedding per domain, passed as a second argument to the classifier.
- *Multi-task learning*: Gradient matching for domain generalization (**Fish**) [42] maximizes the inner product between gradients from different domains through a multi-task objective.
- *Imbalanced learning*: Focal loss (**Focal**) [32] reduces the relative loss for well-classified samples and focuses on difficult samples. Class-balanced loss (**CBLoss**) [10] proposes re-weighting by the inverse effective number of samples. The LDAM loss (**LDAM**) [6] employs a modified marginal loss that favors minority samples more. Balanced-Softmax (**BSoftmax**) [39] extends Softmax to an unbiased estimation that considers the number of samples of each class. Self-supervised pre-training (**SSP**) [52] uses self-supervised learning as a first-stage pre-training to alleviate the network dependence on imbalanced labels. Classifier re-training (**CRT**) [23] decomposes the representation and classifier learning into two stages, where it fine-tunes the classifier using class-balanced sampling with representation fixed in the second stage.

### E.3 Hyperparameters Search Protocol

For a fair evaluation across different algorithms, following the training protocol in [19], for each algorithm we conduct a random search of 20 trials over a joint distribution of its all hyperparameters. We then use the validation set to select the best hyperparameters for each algorithm, fix them and rerun the experiments under 3 different random seeds to report the final average accuracy (and standard deviation). Such process ensures the comparison is best-versus-best, and the hyperparameters are optimized for all algorithms.

We detail the hyperparameter choices for each algorithm in Table 12.

### E.4 Settings for DG Experiments

For DG experiments, we strictly follow the training protocols described in [19]. Across all benchmark DG datasets, we keep the same hyperparameter search space for BoDA as in Table 12. We fix all other training parameters unchanged so that the results of BoDA are directly comparable to the

Table 12: Hyperparameters search space for all experiments.

Condition	Parameter	Default value	Random distribution
<b>General:</b>			
ResNet	learning rate	0.00005	$10^{\text{Uniform}(-5, -3.5)}$
	dropout	0	RandomChoice([0, 0.1, 0.5])
	generator learning rate	0.00005	$10^{\text{Uniform}(-5, -3.5)}$
	discriminator learning rate	0.00005	$10^{\text{Uniform}(-5, -3.5)}$
not ResNet	learning rate	0.001	$10^{\text{Uniform}(-4.5, -3.5)}$
	generator learning rate	0.001	$10^{\text{Uniform}(-4.5, -2.5)}$
	discriminator learning rate	0.001	$10^{\text{Uniform}(-4.5, -2.5)}$
Digits-MLT	weight decay	0	0
	generator weight decay	0	0
not Digits-MLT	weight decay	0	$10^{\text{Uniform}(-6, -2)}$
	generator weight decay	0	$10^{\text{Uniform}(-6, -2)}$
<b>Algorithm-specific:</b>			
IRM	lambda	100	$10^{\text{Uniform}(-1, 5)}$
	iterations of penalty annealing	500	$10^{\text{Uniform}(0, 4)}$
GroupDRO	eta	0.01	$10^{\text{Uniform}(-3, -1)}$
Mixup	alpha	0.2	$10^{\text{Uniform}(0, 4)}$
MLDG	beta	1	$10^{\text{Uniform}(-1, 1)}$
CORAL, MMD	gamma	1	$10^{\text{Uniform}(-1, 1)}$
DANN, CDANN	lambda	1.0	$10^{\text{Uniform}(-2, 2)}$
	discriminator weight decay	0	$10^{\text{Uniform}(-6, -2)}$
	discriminator steps	1	$2^{\text{Uniform}(0, 3)}$
	gradient penalty	0	$10^{\text{Uniform}(-2, 1)}$
	adam $\beta_1$	0.5	RandomChoice([0, 0.5])
MTL	ema	0.99	RandomChoice([.5, .9, .99, 1])
SagNet	adversary weight	0.1	$10^{\text{Uniform}(-2, 1)}$
Fish	meta learning rate	0.5	RandomChoice([.05, .1, .5])
Focal	gamma	1	$0.5 * 10^{\text{Uniform}(0, 1)}$
CBLoss	beta	0.9999	$1 - 10^{\text{Uniform}(-5, -2)}$
LDAM	max_m	0.5	$10^{\text{Uniform}(-1, -0.1)}$
	scale	30	RandomChoice([10, 30])
BoDA	nu	1	$10^{\text{Uniform}(-0.5, 0)}$
	BoDA loss weight	0.1	$10^{\text{Uniform}(-2, -0.5)}$

results in [19].

For model selection, we use the *training-domain validation set* protocol in [19] with 80% – 20% training-validation split, and the average out-domain test performance is reported across all runs for each domain.

## F Complete Results for MDLT

We provide complete evaluation results on the five MDLT datasets. In addition to the reported results in the main paper, for each dataset we also include the accuracy on each domain together with the averaged and the worst accuracy.

### F.1 VLCS-MLT

Table 13: Complete evaluation results on VLCS-MLT.

Algorithm	Accuracy (by domain)						Accuracy (by shot)			
	C	L	S	V	Average	Worst	Many	Medium	Few	Zero
ERM	99.3 ±0.3	53.6 ±1.1	65.9 ±1.2	86.4 ±0.7	76.3 ±0.4	53.6 ±1.1	84.6 ±0.5	76.6 ±0.4	—	32.9 ±0.4
IRM	99.1 ±0.4	52.3 ±0.7	68.8 ±1.4	86.0 ±0.3	76.5 ±0.2	52.3 ±0.7	85.3 ±0.6	75.5 ±1.0	—	33.5 ±1.0
GroupDRO	98.7 ±0.3	54.1 ±1.3	67.5 ±1.5	86.7 ±0.3	76.7 ±0.4	54.1 ±1.3	85.3 ±0.9	76.2 ±1.0	—	34.5 ±2.0
Mixup	99.3 ±0.3	52.7 ±1.3	66.1 ±0.0	85.3 ±1.1	75.9 ±0.1	52.7 ±1.3	84.4 ±0.2	77.1 ±0.6	—	29.2 ±1.4
MLDG	99.3 ±0.3	53.6 ±0.5	68.3 ±0.4	86.4 ±0.5	76.9 ±0.2	53.6 ±0.5	84.9 ±0.3	77.5 ±1.0	—	34.4 ±0.9
CORAL	99.3 ±0.3	51.6 ±0.7	67.5 ±1.8	85.3 ±0.9	75.9 ±0.5	51.6 ±0.7	84.3 ±0.6	75.5 ±0.5	—	34.5 ±0.8
MMD	99.6 ±0.2	53.4 ±0.3	65.6 ±0.8	86.7 ±1.1	76.3 ±0.6	53.4 ±0.3	84.5 ±0.8	77.1 ±0.5	—	32.7 ±0.3
DANN	99.6 ±0.2	54.1 ±0.3	69.9 ±0.2	86.7 ±0.0	77.5 ±0.1	54.1 ±0.3	85.9 ±0.5	76.0 ±0.4	—	38.0 ±2.3
CDANN	99.6 ±0.4	53.6 ±0.4	67.5 ±0.6	85.8 ±0.8	76.6 ±0.4	53.6 ±0.4	84.4 ±0.7	77.3 ±0.8	—	35.0 ±0.8
MTL	99.1 ±0.2	52.9 ±0.5	66.7 ±0.4	86.7 ±0.6	76.3 ±0.3	52.9 ±0.5	84.8 ±0.9	76.2 ±0.6	—	33.3 ±1.4
SagNet	99.6 ±0.4	52.3 ±0.2	67.2 ±0.2	86.2 ±1.0	76.3 ±0.2	52.3 ±0.2	85.3 ±0.3	75.1 ±0.2	—	32.9 ±0.3
Fish	98.7 ±0.3	54.3 ±0.4	69.4 ±0.8	87.6 ±0.4	77.5 ±0.3	54.3 ±0.4	86.2 ±0.5	76.0 ±0.4	—	35.6 ±2.2
Focal	99.1 ±0.4	52.3 ±0.2	66.1 ±0.8	84.9 ±0.2	75.6 ±0.4	52.3 ±0.2	84.0 ±0.2	75.5 ±0.6	—	32.7 ±0.9
CBLoss	99.1 ±0.2	52.5 ±0.5	68.5 ±1.0	87.1 ±1.0	76.8 ±0.3	52.5 ±0.5	84.8 ±0.7	77.5 ±1.4	—	33.2 ±1.6
LDAM	98.9 ±0.2	52.9 ±0.2	69.4 ±1.4	<b>88.0</b> ±1.3	77.5 ±0.1	52.9 ±0.2	<b>86.5</b> ±0.4	75.5 ±0.5	—	35.2 ±0.6
BSoftmax	99.3 ±0.3	52.9 ±0.9	68.0 ±0.2	86.7 ±0.8	76.7 ±0.5	52.9 ±0.9	84.4 ±0.9	78.2 ±0.6	—	34.3 ±0.9
SSP	99.1 ±0.2	52.3 ±1.0	68.0 ±0.2	85.1 ±0.4	76.1 ±0.3	52.3 ±1.0	83.8 ±0.3	76.0 ±1.2	—	37.1 ±0.7
CRT	99.6 ±0.3	51.4 ±0.3	66.9 ±0.8	86.9 ±0.4	76.3 ±0.2	51.4 ±0.3	84.5 ±0.1	77.3 ±0.0	—	31.7 ±1.0
BoDA <sub>r</sub>	99.3 ±0.3	51.4 ±0.3	70.2 ±0.4	86.7 ±0.3	76.9 ±0.5	51.4 ±0.3	85.3 ±0.3	77.3 ±0.2	—	33.3 ±0.5
BoDA-M <sub>r</sub>	<b>100.0</b> ±0.0	53.4 ±0.3	68.5 ±0.4	<b>88.0</b> ±0.8	77.5 ±0.3	53.4 ±0.3	85.8 ±0.2	77.3 ±0.2	—	35.7 ±0.7
BoDA <sub>r,c</sub>	99.3 ±0.3	53.4 ±0.3	68.5 ±0.4	<b>88.0</b> ±0.4	77.3 ±0.2	53.4 ±0.3	85.3 ±0.3	78.0 ±0.2	—	38.6 ±0.7
BoDA-M <sub>r,c</sub>	<b>100.0</b> ±0.0	<b>55.4</b> ±0.5	<b>72.6</b> ±0.3	84.7 ±0.5	<b>78.2</b> ±0.4	<b>55.4</b> ±0.5	85.3 ±0.3	<b>79.3</b> ±0.6	—	<b>43.3</b> ±1.1
BoDA vs. ERM	<b>+0.7</b>	<b>+1.8</b>	<b>+6.7</b>	<b>+1.6</b>	<b>+1.9</b>	<b>+1.8</b>	<b>+0.7</b>	<b>+2.7</b>	—	<b>+10.4</b>

## F.2 PACS-MLT

Table 14: Complete evaluation results on PACS-MLT.

Algorithm	Accuracy (by domain)						Accuracy (by shot)			
	A	C	P	S	Average	Worst	Many	Medium	Few	Zero
ERM	96.8 ±0.1	97.0 ±0.3	98.9 ±0.3	95.8 ±0.2	97.1 ±0.1	95.8 ±0.2	97.1 ±0.0	97.0 ±0.0	98.0 ±0.9	—
IRM	96.8 ±0.1	96.3 ±0.7	98.7 ±0.2	95.2 ±0.4	96.7 ±0.2	95.2 ±0.4	96.8 ±0.2	96.7 ±0.7	94.7 ±1.4	—
GroupDRO	96.9 ±0.2	97.0 ±0.4	99.0 ±0.1	95.3 ±0.4	97.0 ±0.1	95.3 ±0.4	97.3 ±0.1	95.3 ±1.2	94.7 ±3.6	—
Mixup	96.5 ±0.3	96.9 ±0.7	98.5 ±0.2	95.1 ±0.2	96.7 ±0.2	95.1 ±0.2	97.0 ±0.1	96.7 ±0.3	91.3 ±2.7	—
MLDG	96.6 ±0.2	97.2 ±0.3	98.5 ±0.1	94.1 ±0.3	96.6 ±0.1	94.1 ±0.3	96.8 ±0.1	96.3 ±0.7	92.7 ±0.5	—
CORAL	96.9 ±0.4	97.0 ±0.5	98.3 ±0.3	94.3 ±0.7	96.6 ±0.5	94.3 ±0.7	96.6 ±0.5	97.0 ±0.8	94.7 ±0.5	—
MMD	96.8 ±0.2	97.1 ±0.4	97.4 ±0.3	96.3 ±0.3	96.9 ±0.1	96.2 ±0.2	96.9 ±0.2	97.0 ±0.0	96.7 ±0.5	—
DANN	95.7 ±0.3	97.2 ±0.4	98.9 ±0.1	94.3 ±0.1	96.5 ±0.0	94.3 ±0.1	96.5 ±0.1	98.0 ±0.0	94.7 ±2.4	—
CDANN	95.5 ±0.5	96.7 ±0.2	97.2 ±0.3	94.9 ±0.5	96.1 ±0.1	94.5 ±0.2	96.1 ±0.1	96.3 ±0.5	94.0 ±0.9	—
MTL	96.3 ±0.4	97.9 ±0.3	98.2 ±0.3	94.6 ±0.7	96.7 ±0.2	94.5 ±0.6	96.8 ±0.1	95.3 ±1.7	97.3 ±1.1	—
SagNet	97.0 ±0.2	97.8 ±0.4	98.9 ±0.1	95.2 ±0.3	<b>97.2</b> ±0.1	95.2 ±0.3	<b>97.4</b> ±0.1	96.7 ±0.5	95.3 ±0.5	—
Fish	95.5 ±0.2	97.9 ±0.4	98.2 ±0.3	95.9 ±0.5	96.9 ±0.2	95.2 ±0.2	97.0 ±0.1	97.0 ±0.5	94.7 ±1.1	—
Focal	96.6 ±0.4	96.6 ±0.8	98.1 ±0.2	94.6 ±0.7	96.5 ±0.2	94.6 ±0.7	96.6 ±0.1	95.0 ±1.7	96.7 ±0.5	—
CBLoss	<b>97.3</b> ±0.1	97.4 ±0.5	97.8 ±0.6	95.1 ±0.4	96.9 ±0.1	95.1 ±0.4	96.8 ±0.2	97.0 ±1.2	<b>100.0</b> ±0.0	—
LDAM	96.9 ±0.1	96.6 ±0.6	97.9 ±0.1	94.7 ±0.2	96.5 ±0.2	94.7 ±0.2	96.6 ±0.1	95.7 ±1.4	96.0 ±0.0	—
BSoftmax	96.0 ±0.5	96.9 ±0.6	98.8 ±0.6	95.9 ±0.1	96.9 ±0.3	95.6 ±0.3	96.6 ±0.4	<b>98.7</b> ±0.7	99.3 ±0.5	—
SSP	96.2 ±0.5	96.8 ±0.2	98.9 ±0.1	95.7 ±0.3	96.9 ±0.2	95.4 ±0.4	96.7 ±0.2	98.3 ±0.5	98.0 ±0.9	—
CRT	95.3 ±0.2	96.7 ±0.1	98.5 ±0.1	94.9 ±0.1	96.3 ±0.1	94.9 ±0.1	96.3 ±0.1	97.3 ±0.3	94.0 ±0.9	—
BoDA <sub>r</sub>	96.9 ±0.4	97.4 ±0.2	98.6 ±0.2	95.1 ±0.4	97.0 ±0.1	95.1 ±0.4	97.0 ±0.1	96.3 ±0.5	98.0 ±0.9	—
BoDA-M <sub>r</sub>	96.6 ±0.2	<b>98.0</b> ±0.2	99.1 ±0.2	94.9 ±0.1	97.1 ±0.1	94.9 ±0.1	97.3 ±0.1	96.3 ±0.5	96.0 ±0.0	—
BoDA <sub>r,c</sub>	96.3 ±0.1	97.4 ±0.5	<b>99.4</b> ±0.3	95.7 ±0.3	<b>97.2</b> ±0.1	95.7 ±0.3	<b>97.4</b> ±0.1	97.0 ±0.0	94.7 ±1.1	—
BoDA-M <sub>r,c</sub>	96.3 ±0.4	97.7 ±0.2	98.1 ±0.4	<b>96.4</b> ±0.2	97.1 ±0.2	<b>96.3</b> ±0.1	97.1 ±0.0	97.0 ±0.8	96.0 ±0.0	—
BoDA vs. ERM	<b>-0.5</b>	<b>+0.7</b>	<b>+0.5</b>	<b>+0.6</b>	<b>+0.1</b>	<b>+0.5</b>	<b>+0.3</b>	<b>+0.0</b>	<b>-2.0</b>	—

### F.3 OfficeHome-MLT

Table 15: Complete evaluation results on OfficeHome-MLT.

Algorithm	Accuracy (by domain)						Accuracy (by shot)			
	A	C	P	R	Average	Worst	Many	Medium	Few	Zero
ERM	71.3 ±0.1	78.4 ±0.2	89.6 ±0.3	83.3 ±0.2	80.7 ±0.0	71.3 ±0.1	87.8 ±0.2	81.0 ±0.2	63.1 ±0.1	63.3 ±7.2
IRM	70.7 ±0.2	78.5 ±0.8	89.4 ±0.5	83.8 ±0.6	80.6 ±0.4	70.7 ±0.2	87.6 ±0.4	81.5 ±0.4	61.1 ±0.9	56.7 ±1.4
GroupDRO	68.7 ±0.9	79.0 ±0.2	89.4 ±0.4	83.3 ±0.5	80.1 ±0.3	68.7 ±0.9	88.1 ±0.2	80.8 ±0.4	59.8 ±1.2	51.7 ±3.6
Mixup	72.3 ±0.6	79.1 ±0.4	89.7 ±0.1	83.9 ±0.2	81.2 ±0.2	72.3 ±0.6	87.9 ±0.4	81.8 ±0.1	64.1 ±0.4	60.0 ±4.1
MLDG	70.2 ±0.6	78.2 ±0.5	89.4 ±0.4	83.7 ±0.3	80.4 ±0.2	70.2 ±0.6	87.1 ±0.1	81.3 ±0.3	61.3 ±1.0	61.7 ±1.4
CORAL	<b>72.7</b> ±0.6	80.9 ±0.3	89.9 ±0.2	84.2 ±0.4	81.9 ±0.1	<b>72.7</b> ±0.6	87.9 ±0.1	83.0 ±0.1	63.5 ±0.7	65.0 ±2.4
MMD	67.7 ±0.8	77.8 ±0.2	87.4 ±0.5	80.6 ±0.4	78.4 ±0.4	67.7 ±0.8	85.2 ±0.2	79.4 ±0.7	58.8 ±0.4	56.7 ±3.6
DANN	70.2 ±0.9	77.3 ±0.3	87.3 ±0.5	82.1 ±0.4	79.2 ±0.2	70.2 ±0.9	86.2 ±0.1	80.0 ±0.1	60.3 ±1.1	61.7 ±5.9
CDANN	69.4 ±0.3	77.2 ±0.3	87.7 ±0.2	81.5 ±0.3	79.0 ±0.2	69.4 ±0.3	86.4 ±0.6	79.8 ±0.1	58.9 ±0.8	50.0 ±4.7
MTL	69.8 ±0.6	77.6 ±0.3	87.9 ±0.1	82.4 ±0.3	79.5 ±0.2	69.8 ±0.6	87.3 ±0.3	79.8 ±0.2	61.1 ±0.2	51.7 ±2.7
SagNet	70.5 ±0.5	79.6 ±0.5	89.3 ±0.4	83.9 ±0.1	80.9 ±0.1	70.5 ±0.5	87.8 ±0.4	81.9 ±0.1	61.2 ±0.9	56.7 ±3.6
Fish	71.3 ±0.7	79.1 ±0.1	<b>90.2</b> ±0.6	84.7 ±0.4	81.3 ±0.3	71.3 ±0.7	<b>88.2</b> ±0.2	81.9 ±0.3	63.2 ±0.8	61.7 ±1.4
Focal	67.6 ±0.4	76.6 ±0.8	87.1 ±0.5	80.2 ±0.3	77.9 ±0.0	67.6 ±0.4	86.5 ±0.3	78.3 ±0.1	57.4 ±0.3	46.7 ±3.6
CBLoss	69.5 ±0.7	78.7 ±0.3	88.9 ±0.4	82.2 ±0.1	79.8 ±0.2	69.5 ±0.7	86.6 ±0.4	80.6 ±0.2	61.1 ±1.4	65.0 ±2.4
LDAM	69.9 ±0.5	78.9 ±0.4	89.4 ±0.3	83.0 ±0.4	80.3 ±0.2	69.9 ±0.5	87.1 ±0.2	81.3 ±0.3	61.1 ±0.2	51.7 ±2.7
BSoftmax	70.9 ±0.5	78.7 ±0.2	89.0 ±0.8	83.0 ±0.3	80.4 ±0.2	70.9 ±0.5	86.7 ±0.5	81.3 ±0.3	62.4 ±1.0	60.0 ±4.1
SSP	71.1 ±0.3	79.6 ±0.8	89.4 ±0.3	84.2 ±0.2	81.1 ±0.3	71.1 ±0.3	87.3 ±0.6	82.3 ±0.3	61.6 ±0.7	63.3 ±1.4
CRT	72.5 ±0.2	79.6 ±0.2	88.9 ±0.1	83.6 ±0.2	81.2 ±0.0	72.5 ±0.2	87.7 ±0.1	81.8 ±0.1	64.0 ±0.1	65.0 ±2.4
BoDA <sub>r</sub>	71.8 ±0.1	80.3 ±0.3	89.1 ±0.4	84.6 ±0.2	81.5 ±0.1	71.8 ±0.1	87.7 ±0.2	82.3 ±0.1	<b>64.2</b> ±0.3	63.3 ±1.4
BoDA-M <sub>r</sub>	71.6 ±0.2	80.5 ±0.3	89.2 ±0.2	85.7 ±0.4	81.9 ±0.2	71.6 ±0.2	87.3 ±0.3	83.4 ±0.2	62.3 ±0.3	65.0 ±2.4
BoDA <sub>r,c</sub>	72.3 ±0.3	80.8 ±0.2	89.4 ±0.4	<b>86.3</b> ±0.3	82.3 ±0.1	72.3 ±0.3	87.1 ±0.2	<b>83.9</b> ±0.3	63.2 ±0.2	65.0 ±2.4
BoDA-M <sub>r,c</sub>	72.3 ±0.3	<b>81.5</b> ±0.4	89.5 ±0.3	85.8 ±0.2	<b>82.4</b> ±0.2	72.3 ±0.3	87.7 ±0.1	<b>83.9</b> ±0.6	<b>64.2</b> ±0.3	<b>66.7</b> ±2.7
BoDA <i>vs.</i> ERM	<b>+1.0</b>	<b>+3.1</b>	<b>-0.1</b>	<b>+3.0</b>	<b>+1.7</b>	<b>+1.0</b>	<b>-0.1</b>	<b>+2.9</b>	<b>+1.1</b>	<b>+3.4</b>

## F.4 TerraInc-MLT

Table 16: Complete evaluation results on TerraInc-MLT.

Algorithm	Accuracy (by domain)						Accuracy (by shot)			
	L100	L38	L43	L46	Average	Worst	Many	Medium	Few	Zero
ERM	80.3 ±1.3	71.2 ±0.7	82.2 ±0.3	67.4 ±0.3	75.3 ±0.3	67.4 ±0.3	85.6 ±0.8	69.6 ±3.2	66.1 ±2.4	14.4 ±2.8
IRM	78.2 ±0.9	69.6 ±2.0	81.1 ±0.7	64.3 ±1.3	73.3 ±0.7	64.3 ±1.3	83.5 ±0.6	70.0 ±1.8	58.3 ±3.4	20.1 ±1.4
GroupDRO	68.3 ±1.0	68.8 ±1.3	82.6 ±0.2	68.1 ±0.8	72.0 ±0.4	66.6 ±0.2	84.7 ±1.1	64.6 ±4.7	38.9 ±1.2	13.5 ±1.1
Mixup	75.4 ±1.4	70.2 ±1.3	78.3 ±0.6	60.4 ±1.1	71.1 ±0.7	60.4 ±1.1	83.2 ±0.7	60.0 ±0.6	56.1 ±3.0	12.2 ±2.1
MLDG	82.3 ±0.9	73.5 ±2.0	83.8 ±1.4	66.9 ±0.5	76.6 ±0.2	66.9 ±0.5	86.1 ±0.6	73.8 ±3.9	70.6 ±3.7	18.8 ±2.4
CORAL	81.6 ±1.0	72.0 ±0.6	84.2 ±0.2	67.8 ±0.9	76.4 ±0.5	67.8 ±0.9	86.3 ±0.3	77.5 ±3.1	66.1 ±2.0	11.0 ±1.4
MMD	78.9 ±0.6	68.8 ±1.0	81.9 ±0.9	63.7 ±1.1	73.3 ±0.4	63.7 ±1.1	84.0 ±0.4	67.9 ±2.7	60.6 ±1.6	13.6 ±2.6
DANN	74.1 ±0.8	63.1 ±1.9	75.9 ±0.2	61.5 ±0.9	68.7 ±0.9	61.1 ±1.0	79.6 ±1.2	62.5 ±8.1	48.9 ±2.8	13.3 ±1.1
CDANN	73.0 ±1.3	67.8 ±2.0	75.0 ±0.6	65.2 ±1.1	70.3 ±0.5	63.9 ±1.0	83.5 ±0.8	50.0 ±4.2	43.9 ±4.7	20.4 ±3.1
MTL	79.4 ±0.8	70.8 ±0.6	81.9 ±0.8	67.8 ±1.4	75.0 ±0.7	67.7 ±1.4	85.2 ±0.7	73.8 ±1.6	61.1 ±2.8	12.4 ±4.0
SagNet	79.4 ±1.8	71.2 ±0.7	83.4 ±2.4	66.5 ±2.1	75.1 ±1.6	66.5 ±2.1	85.5 ±0.9	77.1 ±5.0	57.8 ±4.3	13.0 ±3.4
Fish	80.1 ±1.9	70.2 ±0.2	84.4 ±0.9	66.3 ±0.5	75.3 ±0.5	66.3 ±0.5	85.8 ±0.2	73.3 ±3.9	61.1 ±3.0	13.7 ±3.3
Focal	80.9 ±0.7	71.6 ±1.6	84.4 ±1.3	66.1 ±1.7	75.7 ±0.4	65.3 ±1.1	85.7 ±0.3	76.2 ±3.9	68.9 ±3.2	12.6 ±1.9
CBLoss	84.9 ±0.6	78.0 ±1.2	80.7 ±0.3	68.3 ±2.0	78.0 ±0.4	68.3 ±2.0	85.0 ±0.1	89.2 ±1.2	83.9 ±2.5	9.3 ±3.9
LDAM	83.0 ±0.9	70.6 ±0.6	81.3 ±1.1	64.1 ±1.4	74.7 ±0.9	64.1 ±1.4	85.1 ±0.6	70.8 ±3.5	67.8 ±1.2	11.1 ±2.4
BSoftmax	83.5 ±2.1	75.5 ±0.4	82.1 ±0.7	65.6 ±1.3	76.7 ±1.0	65.6 ±1.3	83.4 ±0.8	90.8 ±0.9	78.3 ±3.9	12.6 ±2.4
SSP	82.6 ±1.3	80.7 ±1.8	83.2 ±0.6	67.3 ±0.4	78.5 ±0.7	67.3 ±0.4	85.5 ±1.0	87.8 ±0.9	82.6 ±1.2	13.2 ±2.8
CRT	89.0 ±0.1	81.8 ±0.3	85.8 ±0.3	70.0 ±0.4	81.6 ±0.1	70.0 ±0.4	<b>89.7</b> ±0.2	90.4 ±0.3	83.9 ±0.5	12.9 ±0.0
BoDA <sub>r</sub>	86.7 ±0.7	74.1 ±1.1	85.2 ±0.7	68.5 ±0.3	78.6 ±0.4	68.5 ±0.3	86.4 ±0.1	85.0 ±1.0	80.0 ±0.9	13.7 ±2.1
BoDA-M <sub>r</sub>	87.8 ±0.9	76.5 ±0.9	82.2 ±0.3	71.3 ±0.4	79.4 ±0.6	71.3 ±0.4	88.4 ±0.3	76.2 ±2.7	88.3 ±1.6	14.4 ±1.4
BoDA <sub>r,c</sub>	88.3 ±0.6	<b>82.9</b> ±0.5	<b>89.3</b> ±0.9	68.5 ±0.6	82.3 ±0.3	68.5 ±0.6	89.2 ±0.2	<b>92.5</b> ±0.9	88.3 ±1.2	21.3 ±0.7
BoDA-M <sub>r,c</sub>	<b>90.4</b> ±0.3	81.2 ±0.7	85.8 ±0.4	<b>74.6</b> ±0.7	<b>83.0</b> ±0.4	<b>74.6</b> ±0.7	89.2 ±0.2	91.2 ±0.6	<b>91.7</b> ±2.0	<b>21.7</b> ±1.4
BoDA <i>vs.</i> ERM	<b>+10.1</b>	<b>+11.7</b>	<b>+7.1</b>	<b>+7.2</b>	<b>+7.7</b>	<b>+7.2</b>	<b>+3.6</b>	<b>+22.9</b>	<b>+25.6</b>	<b>+7.3</b>

## F.5 DomainNet-MLT

Table 17: Complete evaluation results on DomainNet-MLT.

Algorithm	Accuracy (by domain)							Accuracy (by shot)				
	clip	info	paint	quick	real	sketch	Average	Worst	Many	Medium	Few	Zero
ERM	68.6 ±0.1	29.4 ±0.3	57.1 ±0.2	62.8 ±0.3	72.1 ±0.2	61.7 ±0.2	58.6 ±0.2	29.4 ±0.3	66.0 ±0.1	56.1 ±0.1	35.9 ±0.5	27.6 ±0.3
IRM	66.7 ±0.2	27.6 ±0.1	56.0 ±0.2	60.1 ±0.1	72.0 ±0.0	60.2 ±0.2	57.1 ±0.1	27.6 ±0.1	64.7 ±0.1	54.3 ±0.3	33.5 ±0.3	25.8 ±0.3
GroupDRO	60.1 ±0.2	25.9 ±0.2	50.3 ±0.1	63.9 ±0.2	64.9 ±0.2	56.7 ±0.3	53.6 ±0.1	25.9 ±0.2	61.8 ±0.1	49.1 ±0.3	30.7 ±0.7	22.0 ±0.1
Mixup	67.6 ±0.2	28.7 ±0.0	56.4 ±0.2	60.0 ±0.4	72.1 ±0.1	60.9 ±0.1	57.6 ±0.1	28.7 ±0.0	64.9 ±0.2	54.5 ±0.1	35.6 ±0.2	27.3 ±0.3
MLDG	68.0 ±0.2	28.7 ±0.1	57.2 ±0.1	61.6 ±0.2	73.3 ±0.1	61.9 ±0.2	58.5 ±0.0	28.7 ±0.1	66.0 ±0.1	55.7 ±0.1	35.3 ±0.2	26.9 ±0.3
CORAL	69.1 ±0.3	30.1 ±0.4	57.8 ±0.2	63.4 ±0.2	72.8 ±0.2	63.3 ±0.3	59.4 ±0.1	30.1 ±0.4	66.4 ±0.1	57.1 ±0.0	37.7 ±0.6	29.9 ±0.2
MMD	66.1 ±0.1	27.2 ±0.2	55.9 ±0.1	59.3 ±0.2	71.9 ±0.1	60.0 ±0.2	56.7 ±0.0	27.2 ±0.2	64.2 ±0.1	54.0 ±0.0	33.9 ±0.2	25.4 ±0.2
DANN	65.5 ±0.3	26.9 ±0.4	55.2 ±0.1	57.4 ±0.2	70.6 ±0.1	59.0 ±0.2	55.8 ±0.1	26.9 ±0.4	63.0 ±0.1	52.7 ±0.1	34.2 ±0.4	26.8 ±0.4
CDANN	65.9 ±0.1	27.7 ±0.1	55.3 ±0.1	57.6 ±0.2	70.9 ±0.2	58.7 ±0.1	56.0 ±0.1	27.7 ±0.1	63.2 ±0.0	52.7 ±0.2	34.3 ±0.5	27.6 ±0.1
MTL	68.2 ±0.2	29.3 ±0.2	57.3 ±0.1	62.1 ±0.1	72.9 ±0.1	61.8 ±0.2	58.6 ±0.1	29.3 ±0.2	65.9 ±0.1	56.0 ±0.4	35.4 ±0.1	28.2 ±0.3
SagNet	68.5 ±0.1	29.4 ±0.2	57.8 ±0.2	62.1 ±0.2	73.3 ±0.1	62.4 ±0.1	58.9 ±0.0	29.4 ±0.2	66.3 ±0.1	56.4 ±0.0	36.2 ±0.3	27.2 ±0.4
Fish	68.7 ±0.1	29.1 ±0.1	58.4 ±0.1	64.1 ±0.1	73.9 ±0.1	63.7 ±0.1	59.6 ±0.1	29.1 ±0.1	67.1 ±0.1	57.2 ±0.1	36.8 ±0.4	27.8 ±0.3
Focal	67.6 ±0.1	27.5 ±0.1	56.5 ±0.3	62.3 ±0.3	71.7 ±0.3	61.4 ±0.3	57.8 ±0.2	27.5 ±0.1	65.2 ±0.2	55.1 ±0.2	35.8 ±0.1	26.3 ±0.1
CBLoss	68.3 ±0.2	30.1 ±0.1	57.8 ±0.1	60.8 ±0.1	73.3 ±0.2	63.3 ±0.1	58.9 ±0.1	30.1 ±0.1	64.3 ±0.0	61.0 ±0.3	42.5 ±0.4	28.1 ±0.2
LDAM	68.8 ±0.2	29.2 ±0.2	57.1 ±0.1	65.0 ±0.0	72.3 ±0.1	63.1 ±0.1	59.2 ±0.0	29.2 ±0.2	66.6 ±0.0	57.0 ±0.0	37.1 ±0.2	27.8 ±0.3
BSoftmax	68.5 ±0.1	29.9 ±0.1	57.8 ±0.1	60.5 ±0.3	73.4 ±0.1	63.3 ±0.0	58.9 ±0.1	29.9 ±0.1	64.3 ±0.1	60.9 ±0.3	42.4 ±0.6	28.2 ±0.1
SSP	69.7 ±0.1	31.6 ±0.2	58.8 ±0.1	59.7 ±0.3	73.9 ±0.1	64.2 ±0.1	59.7 ±0.0	31.6 ±0.2	64.3 ±0.1	62.6 ±0.1	45.0 ±0.3	30.5 ±0.0
CRT	70.0 ±0.1	31.6 ±0.1	59.2 ±0.2	64.0 ±0.1	73.4 ±0.1	64.4 ±0.1	60.4 ±0.2	31.6 ±0.1	66.8 ±0.0	61.6 ±0.1	45.7 ±0.1	29.7 ±0.1
BoDA <sub>r</sub>	70.0 ±0.1	32.6 ±0.1	59.1 ±0.1	61.2 ±0.4	73.3 ±0.1	64.1 ±0.1	60.1 ±0.2	32.6 ±0.1	65.7 ±0.2	60.6 ±0.1	42.6 ±0.3	30.5 ±0.2
BoDA-M <sub>r</sub>	70.6 ±0.1	32.2 ±0.2	57.7 ±0.3	<b>65.5</b> ±0.3	70.2 ±0.1	64.5 ±0.1	60.1 ±0.2	32.2 ±0.2	65.9 ±0.2	60.7 ±0.1	42.9 ±0.3	30.0 ±0.1
BoDA <sub>r,c</sub>	<b>72.0</b> ±0.2	<b>33.4</b> ±0.1	60.7 ±0.2	63.6 ±0.2	<b>74.6</b> ±0.1	65.5 ±0.2	<b>61.7</b> ±0.1	<b>33.4</b> ±0.1	<b>67.0</b> ±0.1	62.7 ±0.1	46.0 ±0.2	<b>32.2</b> ±0.3
BoDA-M <sub>r,c</sub>	71.8 ±0.1	33.3 ±0.1	<b>60.8</b> ±0.1	63.7 ±0.3	<b>74.6</b> ±0.1	<b>65.8</b> ±0.2	<b>61.7</b> ±0.2	33.3 ±0.1	<b>67.0</b> ±0.1	<b>63.0</b> ±0.3	<b>46.6</b> ±0.4	31.8 ±0.2
BoDA vs. ERM	<b>+3.4</b>	<b>+4.0</b>	<b>+3.7</b>	<b>+0.9</b>	<b>+2.5</b>	<b>+4.1</b>	<b>+3.1</b>	<b>+4.0</b>	<b>+1.0</b>	<b>+6.9</b>	<b>+10.7</b>	<b>+4.6</b>

## G Complete Results for DG

We provide detailed results of Table 9 across five DG benchmarks [19]. Results for all algorithms except BoDA are directly copied from [19].

### G.1 VLCS

Table 18: Complete domain generalization results on VLCS.

Algorithm	C	L	S	V	Avg
ERM	97.7 $\pm$ 0.4	64.3 $\pm$ 0.9	73.4 $\pm$ 0.5	74.6 $\pm$ 1.3	77.5
IRM	98.6 $\pm$ 0.1	64.9 $\pm$ 0.9	73.4 $\pm$ 0.6	77.3 $\pm$ 0.9	78.5
GroupDRO	97.3 $\pm$ 0.3	63.4 $\pm$ 0.9	69.5 $\pm$ 0.8	76.7 $\pm$ 0.7	76.7
Mixup	98.3 $\pm$ 0.6	64.8 $\pm$ 1.0	72.1 $\pm$ 0.5	74.3 $\pm$ 0.8	77.4
MLDG	97.4 $\pm$ 0.2	65.2 $\pm$ 0.7	71.0 $\pm$ 1.4	75.3 $\pm$ 1.0	77.2
CORAL	98.3 $\pm$ 0.1	<b>66.1</b> $\pm$ 1.2	73.4 $\pm$ 0.3	77.5 $\pm$ 1.2	<b>78.8</b>
MMD	97.7 $\pm$ 0.1	64.0 $\pm$ 1.1	72.8 $\pm$ 0.2	75.3 $\pm$ 3.3	77.5
DANN	<b>99.0</b> $\pm$ 0.3	65.1 $\pm$ 1.4	73.1 $\pm$ 0.3	77.2 $\pm$ 0.6	78.6
CDANN	97.1 $\pm$ 0.3	65.1 $\pm$ 1.2	70.7 $\pm$ 0.8	77.1 $\pm$ 1.5	77.5
MTL	97.8 $\pm$ 0.4	64.3 $\pm$ 0.3	71.5 $\pm$ 0.7	75.3 $\pm$ 1.7	77.2
SagNet	97.9 $\pm$ 0.4	64.5 $\pm$ 0.5	71.4 $\pm$ 1.3	77.5 $\pm$ 0.5	77.8
ARM	98.7 $\pm$ 0.2	63.6 $\pm$ 0.7	71.3 $\pm$ 1.2	76.7 $\pm$ 0.6	77.6
VREx	98.4 $\pm$ 0.3	64.4 $\pm$ 1.4	74.1 $\pm$ 0.4	76.2 $\pm$ 1.3	78.3
RSC	97.9 $\pm$ 0.1	62.5 $\pm$ 0.7	72.3 $\pm$ 1.2	75.6 $\pm$ 0.8	77.1
BoDA	98.1 $\pm$ 0.3	64.5 $\pm$ 0.4	<b>74.3</b> $\pm$ 0.3	<b>78.0</b> $\pm$ 0.6	78.5

### G.2 PACS

Table 19: Complete domain generalization results on PACS.

Algorithm	A	C	P	S	Avg
ERM	84.7 $\pm$ 0.4	80.8 $\pm$ 0.6	97.2 $\pm$ 0.3	79.3 $\pm$ 1.0	85.5
IRM	84.8 $\pm$ 1.3	76.4 $\pm$ 1.1	96.7 $\pm$ 0.6	76.1 $\pm$ 1.0	83.5
GroupDRO	83.5 $\pm$ 0.9	79.1 $\pm$ 0.6	96.7 $\pm$ 0.3	78.3 $\pm$ 2.0	84.4
Mixup	86.1 $\pm$ 0.5	78.9 $\pm$ 0.8	97.6 $\pm$ 0.1	75.8 $\pm$ 1.8	84.6
MLDG	85.5 $\pm$ 1.4	80.1 $\pm$ 1.7	97.4 $\pm$ 0.3	76.6 $\pm$ 1.1	84.9
CORAL	<b>88.3</b> $\pm$ 0.2	80.0 $\pm$ 0.5	97.5 $\pm$ 0.3	78.8 $\pm$ 1.3	86.2
MMD	86.1 $\pm$ 1.4	79.4 $\pm$ 0.9	96.6 $\pm$ 0.2	76.5 $\pm$ 0.5	84.6
DANN	86.4 $\pm$ 0.8	77.4 $\pm$ 0.8	97.3 $\pm$ 0.4	73.5 $\pm$ 2.3	83.7
CDANN	84.6 $\pm$ 1.8	75.5 $\pm$ 0.9	96.8 $\pm$ 0.3	73.5 $\pm$ 0.6	82.6
MTL	87.5 $\pm$ 0.8	77.1 $\pm$ 0.5	96.4 $\pm$ 0.8	77.3 $\pm$ 1.8	84.6
SagNet	87.4 $\pm$ 1.0	80.7 $\pm$ 0.6	97.1 $\pm$ 0.1	80.0 $\pm$ 0.4	86.3
ARM	86.8 $\pm$ 0.6	76.8 $\pm$ 0.5	97.4 $\pm$ 0.3	79.3 $\pm$ 1.2	85.1
VREx	86.0 $\pm$ 1.6	79.1 $\pm$ 0.6	96.9 $\pm$ 0.5	77.7 $\pm$ 1.7	84.9
RSC	85.4 $\pm$ 0.8	79.7 $\pm$ 1.8	97.6 $\pm$ 0.3	78.2 $\pm$ 1.2	85.2
BoDA	88.2 $\pm$ 0.2	<b>81.7</b> $\pm$ 0.3	<b>97.8</b> $\pm$ 0.2	<b>80.2</b> $\pm$ 0.3	<b>86.9</b>

### G.3 OfficeHome

Table 20: Complete domain generalization results on OfficeHome.

Algorithm	A	C	P	R	Avg
ERM	61.3 $\pm$ 0.7	52.4 $\pm$ 0.3	75.8 $\pm$ 0.1	76.6 $\pm$ 0.3	66.5
IRM	58.9 $\pm$ 2.3	52.2 $\pm$ 1.6	72.1 $\pm$ 2.9	74.0 $\pm$ 2.5	64.3
GroupDRO	60.4 $\pm$ 0.7	52.7 $\pm$ 1.0	75.0 $\pm$ 0.7	76.0 $\pm$ 0.7	66.0
Mixup	62.4 $\pm$ 0.8	54.8 $\pm$ 0.6	76.9 $\pm$ 0.3	78.3 $\pm$ 0.2	68.1
MLDG	61.5 $\pm$ 0.9	53.2 $\pm$ 0.6	75.0 $\pm$ 1.2	77.5 $\pm$ 0.4	66.8
CORAL	65.3 $\pm$ 0.4	54.4 $\pm$ 0.5	76.5 $\pm$ 0.1	78.4 $\pm$ 0.5	68.7
MMD	60.4 $\pm$ 0.2	53.3 $\pm$ 0.3	74.3 $\pm$ 0.1	77.4 $\pm$ 0.6	66.3
DANN	59.9 $\pm$ 1.3	53.0 $\pm$ 0.3	73.6 $\pm$ 0.7	76.9 $\pm$ 0.5	65.9
CDANN	61.5 $\pm$ 1.4	50.4 $\pm$ 2.4	74.4 $\pm$ 0.9	76.6 $\pm$ 0.8	65.8
MTL	61.5 $\pm$ 0.7	52.4 $\pm$ 0.6	74.9 $\pm$ 0.4	76.8 $\pm$ 0.4	66.4
SagNet	63.4 $\pm$ 0.2	54.8 $\pm$ 0.4	75.8 $\pm$ 0.4	78.3 $\pm$ 0.3	68.1
ARM	58.9 $\pm$ 0.8	51.0 $\pm$ 0.5	74.1 $\pm$ 0.1	75.2 $\pm$ 0.3	64.8
VREx	60.7 $\pm$ 0.9	53.0 $\pm$ 0.9	75.3 $\pm$ 0.1	76.6 $\pm$ 0.5	66.4
RSC	60.7 $\pm$ 1.4	51.4 $\pm$ 0.3	74.8 $\pm$ 1.1	75.1 $\pm$ 1.3	65.5
<b>BoDA</b>	<b>65.4 <math>\pm</math>0.1</b>	<b>55.4 <math>\pm</math>0.3</b>	<b>77.1 <math>\pm</math>0.1</b>	<b>79.5 <math>\pm</math>0.3</b>	<b>69.3</b>

### G.4 TerraInc

Table 21: Complete domain generalization results on TerraInc.

Algorithm	L100	L38	L43	L46	Avg
ERM	49.8 $\pm$ 4.4	42.1 $\pm$ 1.4	56.9 $\pm$ 1.8	35.7 $\pm$ 3.9	46.1
IRM	54.6 $\pm$ 1.3	39.8 $\pm$ 1.9	56.2 $\pm$ 1.8	39.6 $\pm$ 0.8	47.6
GroupDRO	41.2 $\pm$ 0.7	38.6 $\pm$ 2.1	56.7 $\pm$ 0.9	36.4 $\pm$ 2.1	43.2
Mixup	<b>59.6 <math>\pm</math>2.0</b>	42.2 $\pm$ 1.4	55.9 $\pm$ 0.8	33.9 $\pm$ 1.4	47.9
MLDG	54.2 $\pm$ 3.0	44.3 $\pm$ 1.1	55.6 $\pm$ 0.3	36.9 $\pm$ 2.2	47.7
CORAL	51.6 $\pm$ 2.4	42.2 $\pm$ 1.0	57.0 $\pm$ 1.0	39.8 $\pm$ 2.9	47.6
MMD	41.9 $\pm$ 3.0	34.8 $\pm$ 1.0	57.0 $\pm$ 1.9	35.2 $\pm$ 1.8	42.2
DANN	51.1 $\pm$ 3.5	40.6 $\pm$ 0.6	57.4 $\pm$ 0.5	37.7 $\pm$ 1.8	46.7
CDANN	47.0 $\pm$ 1.9	41.3 $\pm$ 4.8	54.9 $\pm$ 1.7	39.8 $\pm$ 2.3	45.8
MTL	49.3 $\pm$ 1.2	39.6 $\pm$ 6.3	55.6 $\pm$ 1.1	37.8 $\pm$ 0.8	45.6
SagNet	53.0 $\pm$ 2.9	43.0 $\pm$ 2.5	57.9 $\pm$ 0.6	40.4 $\pm$ 1.3	48.6
ARM	49.3 $\pm$ 0.7	38.3 $\pm$ 2.4	55.8 $\pm$ 0.8	38.7 $\pm$ 1.3	45.5
VREx	48.2 $\pm$ 4.3	41.7 $\pm$ 1.3	56.8 $\pm$ 0.8	38.7 $\pm$ 3.1	46.4
RSC	50.2 $\pm$ 2.2	39.2 $\pm$ 1.4	56.3 $\pm$ 1.4	40.8 $\pm$ 0.6	46.6
<b>BoDA</b>	<b>54.0 <math>\pm</math>0.3</b>	<b>46.5 <math>\pm</math>0.2</b>	<b>59.5 <math>\pm</math>0.3</b>	<b>41.0 <math>\pm</math>0.4</b>	<b>50.2</b>

## G.5 DomainNet

Table 22: Complete domain generalization results on DomainNet.

Algorithm	clip	info	paint	quick	real	sketch	Avg
ERM	58.1 $\pm$ 0.3	18.8 $\pm$ 0.3	46.7 $\pm$ 0.3	12.2 $\pm$ 0.4	59.6 $\pm$ 0.1	49.8 $\pm$ 0.4	40.9
IRM	48.5 $\pm$ 2.8	15.0 $\pm$ 1.5	38.3 $\pm$ 4.3	10.9 $\pm$ 0.5	48.2 $\pm$ 5.2	42.3 $\pm$ 3.1	33.9
GroupDRO	47.2 $\pm$ 0.5	17.5 $\pm$ 0.4	33.8 $\pm$ 0.5	9.3 $\pm$ 0.3	51.6 $\pm$ 0.4	40.1 $\pm$ 0.6	33.3
Mixup	55.7 $\pm$ 0.3	18.5 $\pm$ 0.5	44.3 $\pm$ 0.5	12.5 $\pm$ 0.4	55.8 $\pm$ 0.3	48.2 $\pm$ 0.5	39.2
MLDG	59.1 $\pm$ 0.2	19.1 $\pm$ 0.3	45.8 $\pm$ 0.7	13.4 $\pm$ 0.3	59.6 $\pm$ 0.2	50.2 $\pm$ 0.4	41.2
CORAL	59.2 $\pm$ 0.1	19.7 $\pm$ 0.2	46.6 $\pm$ 0.3	13.4 $\pm$ 0.4	59.8 $\pm$ 0.2	50.1 $\pm$ 0.6	41.5
MMD	32.1 $\pm$ 13.3	11.0 $\pm$ 4.6	26.8 $\pm$ 11.3	8.7 $\pm$ 2.1	32.7 $\pm$ 13.8	28.9 $\pm$ 11.9	23.4
DANN	53.1 $\pm$ 0.2	18.3 $\pm$ 0.1	44.2 $\pm$ 0.7	11.8 $\pm$ 0.1	55.5 $\pm$ 0.4	46.8 $\pm$ 0.6	38.3
CDANN	54.6 $\pm$ 0.4	17.3 $\pm$ 0.1	43.7 $\pm$ 0.9	12.1 $\pm$ 0.7	56.2 $\pm$ 0.4	45.9 $\pm$ 0.5	38.3
MTL	57.9 $\pm$ 0.5	18.5 $\pm$ 0.4	46.0 $\pm$ 0.1	12.5 $\pm$ 0.1	59.5 $\pm$ 0.3	49.2 $\pm$ 0.1	40.6
SagNet	57.7 $\pm$ 0.3	19.0 $\pm$ 0.2	45.3 $\pm$ 0.3	12.7 $\pm$ 0.5	58.1 $\pm$ 0.5	48.8 $\pm$ 0.2	40.3
ARM	49.7 $\pm$ 0.3	16.3 $\pm$ 0.5	40.9 $\pm$ 1.1	9.4 $\pm$ 0.1	53.4 $\pm$ 0.4	43.5 $\pm$ 0.4	35.5
VREx	47.3 $\pm$ 3.5	16.0 $\pm$ 1.5	35.8 $\pm$ 4.6	10.9 $\pm$ 0.3	49.6 $\pm$ 4.9	42.0 $\pm$ 3.0	33.6
RSC	55.0 $\pm$ 1.2	18.3 $\pm$ 0.5	44.4 $\pm$ 0.6	12.2 $\pm$ 0.2	55.7 $\pm$ 0.7	47.8 $\pm$ 0.9	38.9
BoDA	<b>62.1</b> $\pm$ 0.4	<b>20.5</b> $\pm$ 0.7	<b>48.0</b> $\pm$ 0.1	<b>13.8</b> $\pm$ 0.6	<b>60.6</b> $\pm$ 0.4	<b>51.4</b> $\pm$ 0.3	<b>42.7</b>

## G.6 Averages

Table 23: Complete domain generalization results over all DG benchmarks.

Algorithm	VLCS	PACS	OfficeHome	TerraInc	DomainNet	Avg
ERM	77.5 $\pm$ 0.4	85.5 $\pm$ 0.2	66.5 $\pm$ 0.3	46.1 $\pm$ 1.8	40.9 $\pm$ 0.1	63.3
IRM	78.5 $\pm$ 0.5	83.5 $\pm$ 0.8	64.3 $\pm$ 2.2	47.6 $\pm$ 0.8	33.9 $\pm$ 2.8	61.6
GroupDRO	76.7 $\pm$ 0.6	84.4 $\pm$ 0.8	66.0 $\pm$ 0.7	43.2 $\pm$ 1.1	33.3 $\pm$ 0.2	60.7
Mixup	77.4 $\pm$ 0.6	84.6 $\pm$ 0.6	68.1 $\pm$ 0.3	47.9 $\pm$ 0.8	39.2 $\pm$ 0.1	63.4
MLDG	77.2 $\pm$ 0.4	84.9 $\pm$ 1.0	66.8 $\pm$ 0.6	47.7 $\pm$ 0.9	41.2 $\pm$ 0.1	63.6
CORAL	<b>78.8</b> $\pm$ 0.6	86.2 $\pm$ 0.3	68.7 $\pm$ 0.3	47.6 $\pm$ 1.0	41.5 $\pm$ 0.1	64.5
MMD	77.5 $\pm$ 0.9	84.6 $\pm$ 0.5	66.3 $\pm$ 0.1	42.2 $\pm$ 1.6	23.4 $\pm$ 9.5	58.8
DANN	78.6 $\pm$ 0.4	83.6 $\pm$ 0.4	65.9 $\pm$ 0.6	46.7 $\pm$ 0.5	38.3 $\pm$ 0.1	62.6
CDANN	77.5 $\pm$ 0.1	82.6 $\pm$ 0.9	65.8 $\pm$ 1.3	45.8 $\pm$ 1.6	38.3 $\pm$ 0.3	62.0
MTL	77.2 $\pm$ 0.4	84.6 $\pm$ 0.5	66.4 $\pm$ 0.5	45.6 $\pm$ 1.2	40.6 $\pm$ 0.1	62.9
SagNet	77.8 $\pm$ 0.5	86.3 $\pm$ 0.2	68.1 $\pm$ 0.1	48.6 $\pm$ 1.0	40.3 $\pm$ 0.1	64.2
ARM	77.6 $\pm$ 0.3	85.1 $\pm$ 0.4	64.8 $\pm$ 0.3	45.5 $\pm$ 0.3	35.5 $\pm$ 0.2	61.7
VREx	78.3 $\pm$ 0.2	84.9 $\pm$ 0.6	66.4 $\pm$ 0.6	46.4 $\pm$ 0.6	33.6 $\pm$ 2.9	61.9
RSC	77.1 $\pm$ 0.5	85.2 $\pm$ 0.9	65.5 $\pm$ 0.9	46.6 $\pm$ 1.0	38.9 $\pm$ 0.5	62.7
BoDA	78.5 $\pm$ 0.3	<b>86.9</b> $\pm$ 0.4	<b>69.3</b> $\pm$ 0.1	<b>50.2</b> $\pm$ 0.4	<b>42.7</b> $\pm$ 0.1	<b>65.5</b>

## H Additional Analysis and Studies

### H.1 Ablation Studies for BoDA

**Effect of Balanced Distance.** We study the effect of adding balanced distance in BoDA compared to the vanilla DA loss. As Table 24 demonstrates, incorporating balanced distance in BoDA is essential for addressing MDLT: we observe that BoDA improves over DA by a large margin, resulting in an averaged improvements of 2.3% over all MDLT benchmarks. The improvements are especially large on datasets with severe data imbalance across domains (e.g., TerraInc-MLT).

Table 24: Ablation study on effect of adding balanced distance in BoDA.

	VLCS-MLT	PACS-MLT	OfficeHome-MLT	TerraInc-MLT	DomainNet-MLT	Avg
DA	76.6 $\pm$ 0.4	96.8 $\pm$ 0.2	80.7 $\pm$ 0.3	76.4 $\pm$ 0.5	58.9 $\pm$ 0.2	77.9
BoDA	<b>77.3</b> $\pm$ 0.2	<b>97.2</b> $\pm$ 0.1	<b>82.3</b> $\pm$ 0.1	<b>82.3</b> $\pm$ 0.3	<b>61.7</b> $\pm$ 0.1	<b>80.2</b>
Gains	<b>+0.7</b>	<b>+0.4</b>	<b>+1.6</b>	<b>+5.9</b>	<b>+2.8</b>	<b>+2.3</b>

**Effect of Different Distance Calibration Coefficient  $\lambda_{d,c}^{d',c'}$ .** We further investigate the effect of different distance calibration coefficients in BoDA. Recall that  $\lambda_{d,c}^{d',c'} = (N_{d',c'}/N_{d,c})^\nu$  indicates how much we would like to transfer  $(d, c)$  to  $(d', c')$ , based on their relative sample sizes. We vary the value of  $\nu$ , and study its effect on BoDA performance across all MDLT datasets. Table 25 reveals several interesting findings. First, when  $\nu = 0$  (i.e., no calibration is used as the coefficient is always equal to 1), BoDA performance is lower than those with a positive  $\nu$ , confirming the effectiveness of the calibrated distance. Moreover, when we vary  $\nu$  between 0.5 – 1.5, the overall performance gains are similar across different choices, where  $\nu$  around 0.9 seems to achieve the best results. Finally, when compared to ERM, we demonstrate that BoDA consistently obtains notable gains across different  $\nu$ .

Table 25: Ablation study on effect of distance calibration coefficient  $\lambda_{d,c}^{d',c'}$  in BoDA. We vary the value of  $\nu$  and report the averaged results over all five MDLT datasets.

$\nu$	0	0.5	0.7	0.9	1	1.1	1.2	1.5	ERM
BoDA	78.9	80.1	80.0	80.2	80.1	79.8	79.6	79.2	77.6

### H.2 Absolute Accuracy Gains on All MDLT Benchmarks

We provide additional results for understanding how BoDA performs across *all* domain-class pair when cross-domain imbalance occurs. Similar to Fig. 7 in the main text, we plot the absolute gains of BoDA over ERM on all five MDLT datasets, shown in Figs. 9, 10, 11, 12, and 13. Across all datasets, we observe that BoDA establishes large improvements w.r.t. all regions, especially for the few-shot and zero-shot ones.

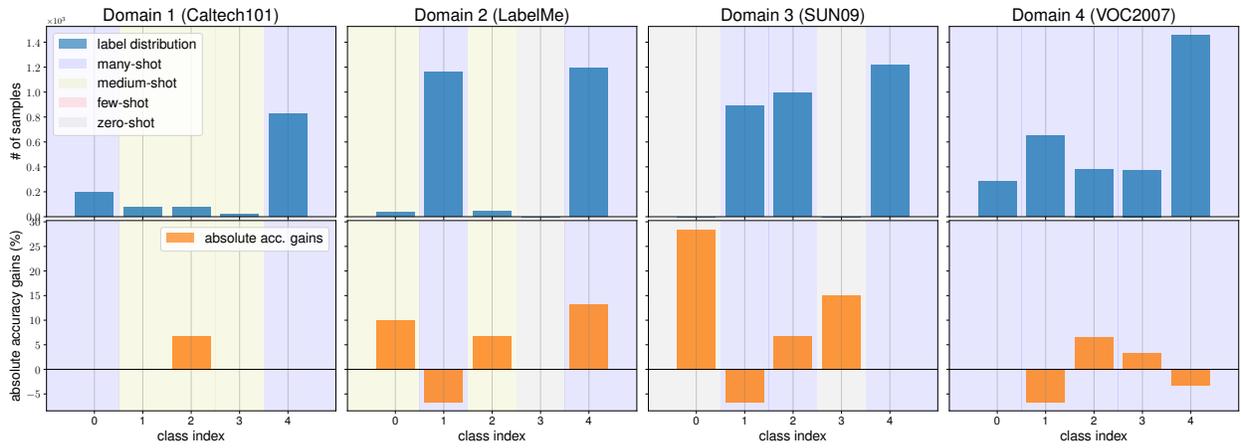


Figure 9: The absolute accuracy gains of BoDA *vs.* ERM over all domain-class pairs on VLCS-MLT.

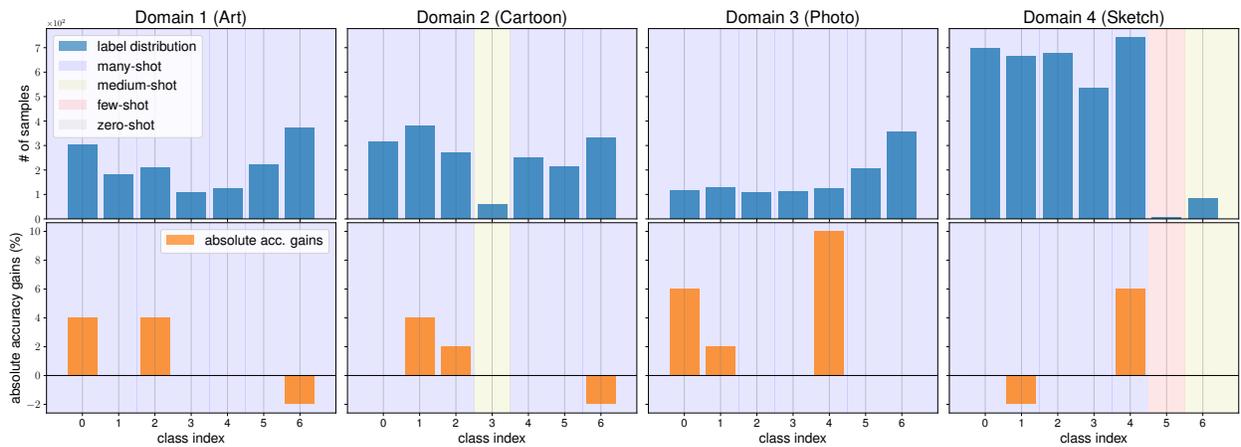


Figure 10: The absolute accuracy gains of BoDA *vs.* ERM over all domain-class pairs on PACS-MLT.

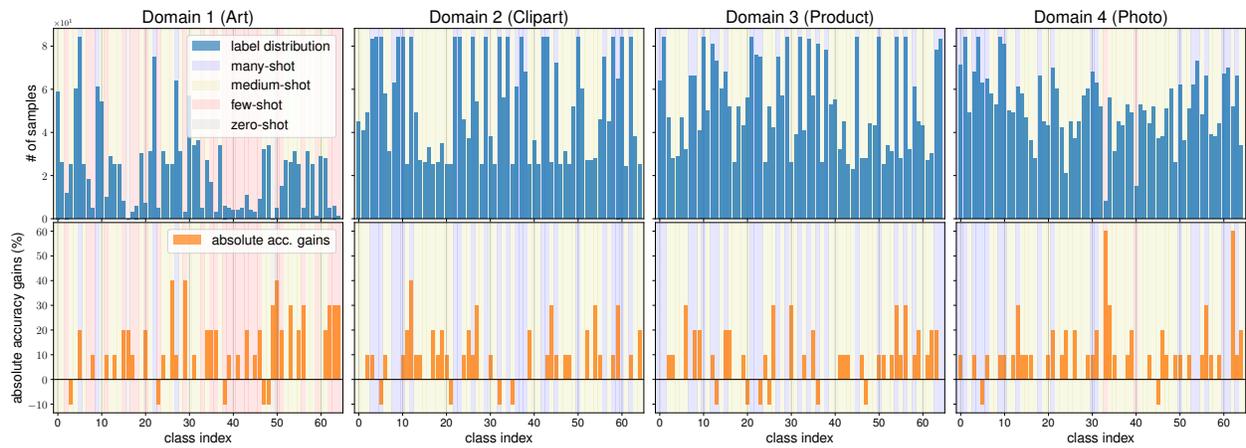


Figure 11: The absolute accuracy gains of BoDA *vs.* ERM over all domain-class pairs on OfficeHome-MLT.

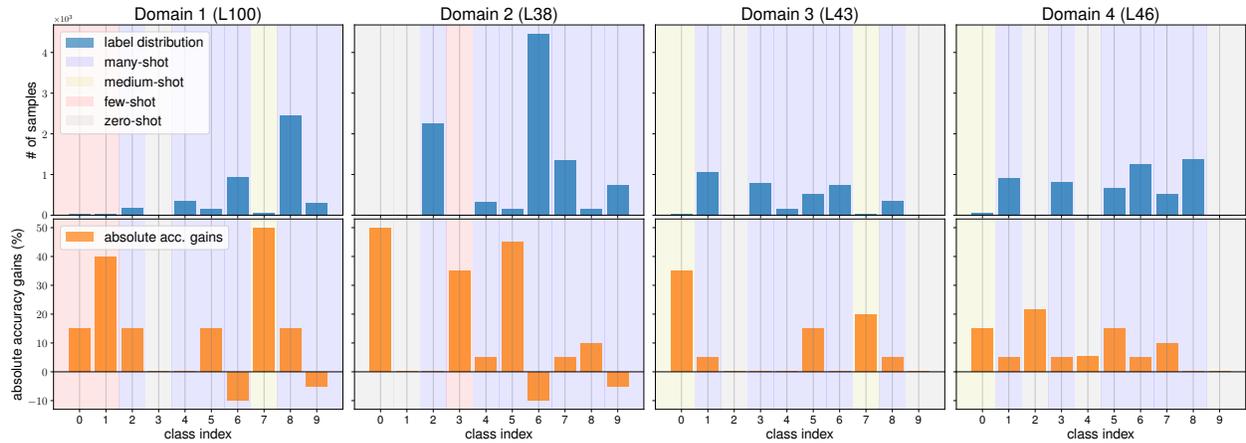


Figure 12: The absolute accuracy gains of BoDA *vs.* ERM over all domain-class pairs on TerraInc-MLT.

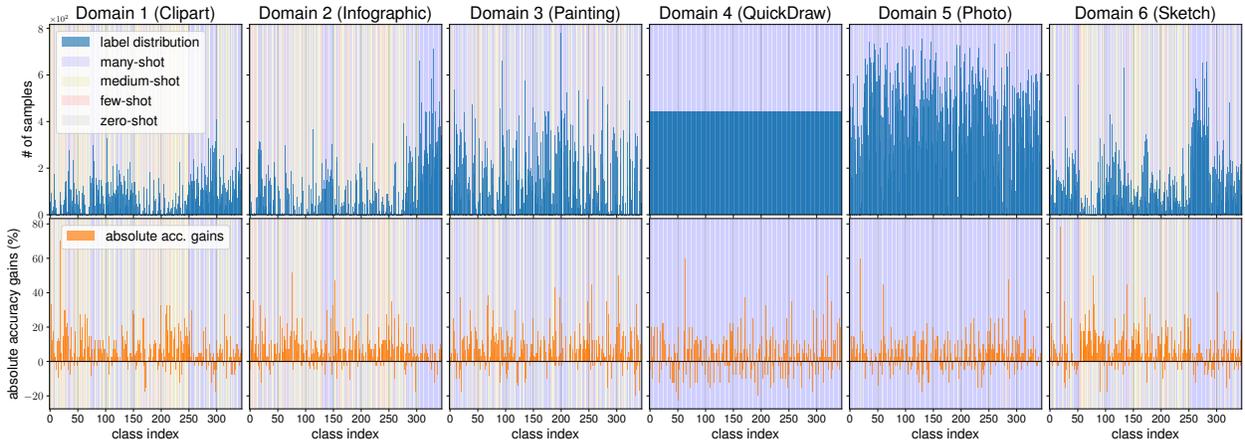


Figure 13: The absolute accuracy gains of BoDA *vs.* ERM over all domain-class pairs on DomainNet-MLT.

### H.3 Robustness to Diverse Skewed Label Distributions

We investigate how BoDA performs under arbitrary label imbalance across domains, especially when the cross-domain label distributions are both *imbalance* and *divergent*. We again employ the *Digits-MLT* dataset, and manually vary the label proportions for each domain.

As Fig. 14 demonstrates, when the label distributions for two domains are balanced and identical, both ERM and BoDA maintains discriminative representations. If the label distributions become imbalanced but still identical across domains, ERM is still able to align similar classes in the two domains, but with majority classes being closer in terms of transferability than minority classes. In contrast, BoDA maintains consistent transferability regardless of number of samples within each class. Finally, as the label distributions become further mismatched across domains, ERM is not able to align the domains and produces a clear gap; by contrast, BoDA maintains consistent and transferable representations even under severe data imbalance. As a result, BoDA substantially boosts the performance upon ERM, with an average gains of 6.4% across all label configurations.

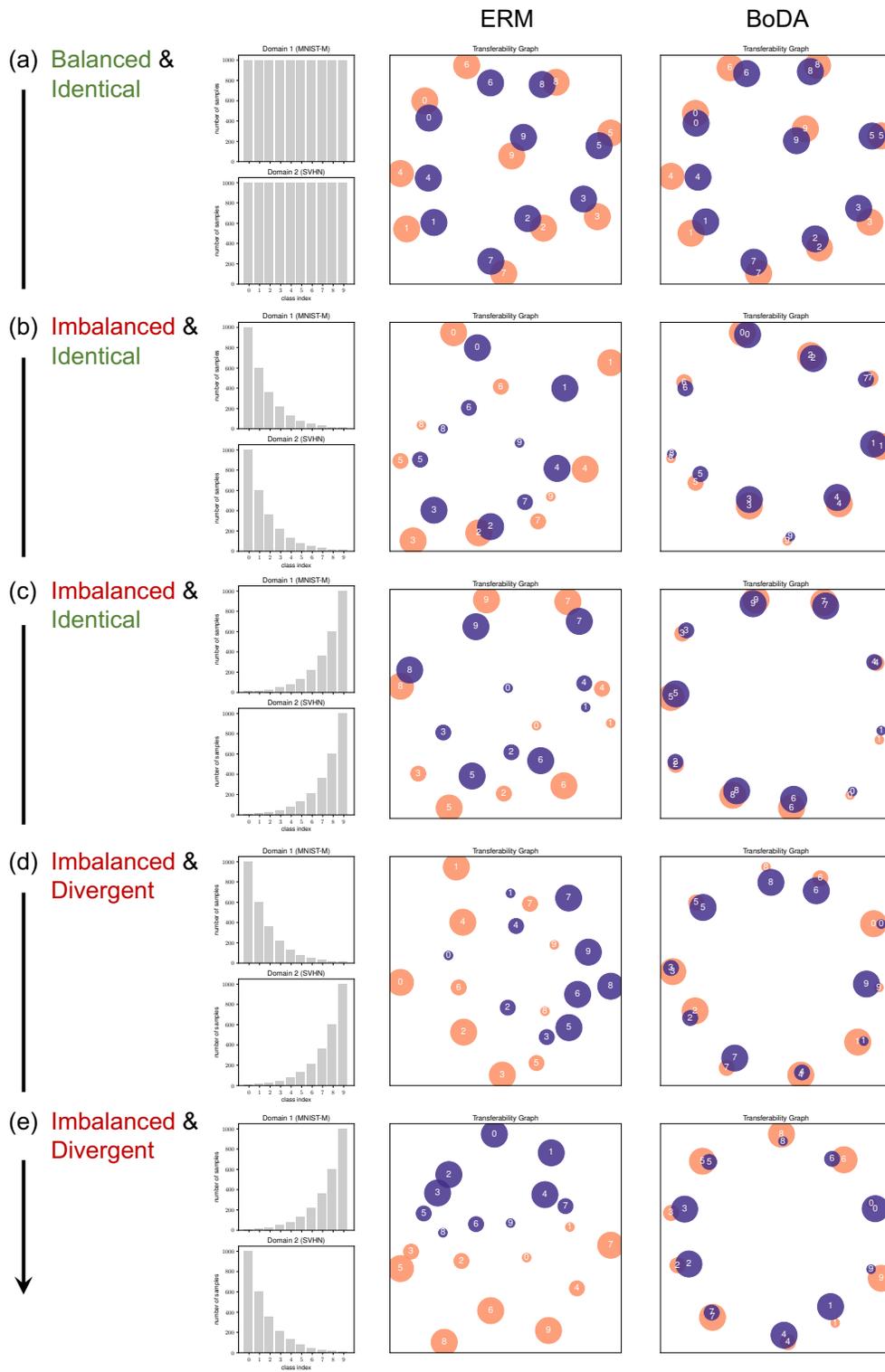


Figure 14: The evolving patterns of the transferability graph of BoDA *vs.* ERM across different label configurations on Digits-MLT. Label distributions for two domains are (a) balanced and identical; (b)(c) imbalanced and identical; (d)(e) imbalanced and divergent. BoDA maintains consistent and transferable representations across all label configurations, and leads to much better test accuracy.

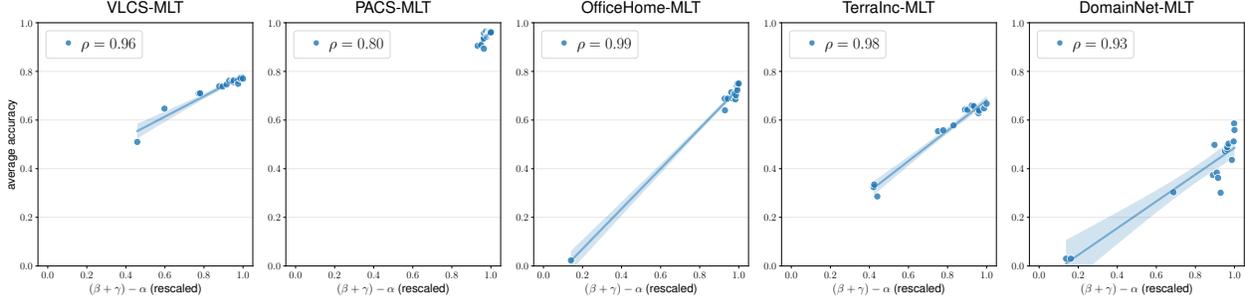


Figure 15: Correspondence between  $(\beta + \gamma) - \alpha$  quantity and test accuracy across different MDLT datasets. Each point within each plot corresponds to a model trained with ERM using different hyperparameters.

#### H.4 Transferability *vs.* Generalization on More Datasets

We provide further results on transferability statistics *vs.* generalization on real MDLT datasets, in addition to results on `Digits-MLT` as we showed in the main text.

Specifically, on all five MDLT datasets, we train 20 ERM models with varying hyperparameters, calculate the  $(\alpha, \beta, \gamma)$  statistics for each model, and plot its classification accuracy against  $(\beta + \gamma) - \alpha$ . Fig. 15 reveals similar and consistent findings, that the  $(\alpha, \beta, \gamma)$  statistics characterize model performance in MDLT. Across all datasets, the  $(\beta + \gamma) - \alpha$  quantity displays a very strong correlation with test performance across the entire range, suggesting that the  $(\alpha, \beta, \gamma)$  statistics govern the success of learning in MDLT.

#### H.5 Additional Visualization of Feature Discrepancy

We provide additional results for understanding BoDA, i.e., how BoDA calibrates the feature statistics. Fig. 16 shows the feature discrepancy of BoDA *vs.* ERM across different label configurations on `Digits-MLT`. In addition to the mean distance we showed in the main text, we show also the feature covariance distance between training and test data, and plot them for both domains. Similarly, solid lines plot the distance between training and test data from the same domain-class pairs. Dashed lines plot the distance between test data from a particular domain-class pair and the training data with which it shares the same class but differs in the domain. The figure also shows regions with different data densities using colors blue, yellow, red.

As the figure confirms, across different label distributions, BoDA consistently learns better representations especially for the tail data (i.e., the red regions), where the feature mean/covariance distance between training and test data becomes smaller and more aligned across domains. Comparing BoDA with ERM further demonstrates that BoDA maintains consistent and transferable representations with smaller feature discrepancy.

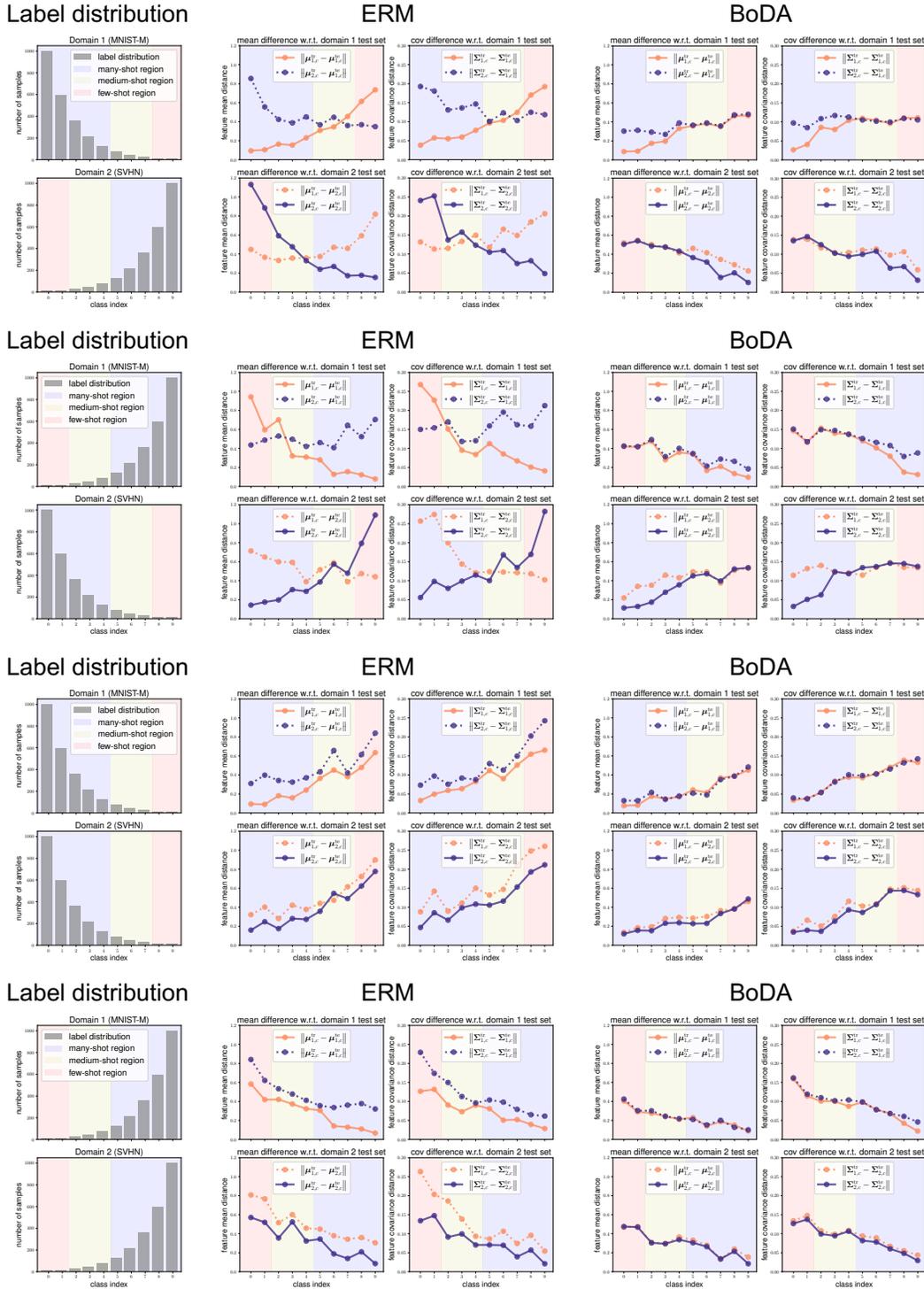


Figure 16: Feature discrepancy of BoDA *vs.* ERM across different label configurations on Digits-MLT. Each row plots a per-domain label distribution, and the feature mean / covariance distance between training and test data on each domain for both ERM and BoDA. BoDA enables better learned tail ( $d, c$ ) with smaller feature discrepancy.