

# Better Pseudo-label: Joint Domain-aware Label and Dual-classifier for Semi-supervised Domain Generalization

Ruiqi Wang<sup>a,b,\*\*</sup>, Lei Qi<sup>c,\*\*</sup>, Yinghuan Shi<sup>a,b,\*</sup>, Yang Gao<sup>a,b</sup>

<sup>a</sup>State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China

<sup>b</sup>National Institute of Healthcare Data Science, Nanjing University, Nanjing, China

<sup>c</sup>School of Computer Science and Engineering, Key Lab of Computer Network and Information Integration (Ministry of Education), Southeast University, Nanjing, China

---

## Abstract

With the goal of directly generalizing trained model to unseen target domains, domain generalization (DG), a newly proposed learning paradigm, has attracted considerable attention. Previous DG models usually require a sufficient quantity of annotated samples from observed source domains during training. In this paper, we relax this requirement about full annotation and investigate semi-supervised domain generalization (SSDG) where only one source domain is fully annotated along with the other domains totally unlabeled in the training process. With the challenges of tackling the domain gap between observed source domains and predicting unseen target domains, we propose a novel deep framework via joint domain-aware labels and dual-classifier to produce high-quality pseudo-labels. Concretely, to predict accurate pseudo-labels under domain shift, a domain-aware pseudo-labeling module is developed. Also, considering inconsistent goals between generalization and pseudo-labeling: former prevents overfitting on all source domains while latter might overfit the unlabeled source domains for high accuracy, we employ a dual-classifier to independently perform pseudo-labeling and domain generalization in the training process. When

---

\*Corresponding author: Yinghuan Shi.

\*\*Co-first authors: Ruiqi Wang and Lei Qi.

*Email addresses:* wangrq@mail.nju.edu.cn (Ruiqi Wang), qilei@seu.edu.cn (Lei Qi), syh@nju.edu.cn (Yinghuan Shi), gaoy@nju.edu.cn (Yang Gao)

accurate pseudo-labels are generated for unlabeled source domains, the domain mixup operation is applied to augment new domains between labeled and unlabeled domains, which is beneficial for boosting the generalization capability of the model. Extensive results on publicly available DG benchmark datasets show the efficacy of our proposed SSDG method.

*Keywords:* Semi-supervised learning, Domain generalization, Image recognition, Feature representation

---

## 1. Introduction

Nowadays, with the development of data acquisition, current data are frequently captured from multiple sources (*e.g.*, video, image, text), generated from various contributors (*e.g.*, different artists), or collected from multiple sites (*e.g.*, different data centers), making the distribution shift between different modalities or sites usually occurs [1, 2]. Therefore, due to the distribution shift, the model trained on training data or source domains could perform poorly on test data or target domains. To address this limitation, a new setting namely domain generalization (DG), aiming to train model on observed source domains for directly generalizing to arbitrary unseen target domains, is becoming a hot topic with increasing interests.

According to our investigation, unfortunately, most current DG models belong to supervised setting where multiple fully labeled source domains are the prerequisite before training DG models. As we known, high-quality labels are often expensive and laborious to obtain, which drives us to alleviate the label requirement in the observed source domains.

Formally, we here name our setting—first training the model with both labeled source domains and unlabeled source domains and then performing prediction on unseen target domains—as semi-supervised domain generalization (SSDG in short). This setting owns its practical meaning. For example, in real-world applications, there are a large number of totally unlabeled datasets (*i.e.*, web-crawled datasets, massive data in data center). The advantage of our

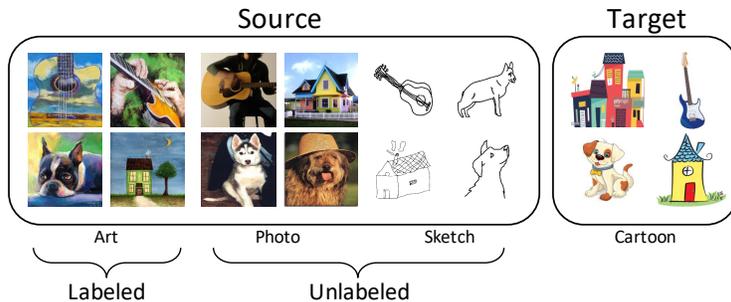


Figure 1: Unlike typical conventional domain generalization, semi-supervised domain generalization takes both the labeled and unlabeled source domains as input, aiming to train an adaptive model for the unseen target domain.

setting is that arbitrary unlabeled domains can be utilized to cooperate with the labeled domains for benefiting domain generalization in a free lunch way. We show this setting in Figure 1. Particularly, in this paper we merely consider the case that only one source domain is fully labeled (along with several unlabeled source domains) in the training stage. The ‘one labeled source domain’ case is more practical because annotating samples is difficult, expensive and time-consuming. And a huge amount of unlabeled data can be easily obtained in real-world applications. However, the ‘one labeled source domain’ case is more challenging because we need to generate pseudo-labels for extensive samples from multiple unlabeled domains with different data distributions.

Since unlabeled samples in source domains are abundant and each unlabeled sample actually belongs to a specific yet unknown class, we consider assigning pseudo-labels to unlabeled samples. Intuitively, the accuracy of pseudo-labels largely affects the final results. The pseudo-labeling technique [3] has shown its effectiveness in conventional semi-supervised learning (SSL) problems by iteratively using higher confident samples to aid subsequent learning on lower confident samples. However, compared with conventional SSL, producing high-quality pseudo-labels in SSDG is much more challenging due to the following two reasons:

1. The domain shift between observed labeled and unlabeled source domains is definitely a negative factor for accurate pseudo-labels, which may lead to a drastic performance degeneration.
2. Since there is the unpredictable domain discrepancy between unlabeled source domains and the unseen target domain, a generalizable model which well fits the target domain may suffer a drop of accuracy on unlabeled source domains.

Considering these two issues, we develop two improvements for accurate pseudo-label prediction.

Firstly, we propose *domain-aware pseudo-labeling method* to improve the quality of pseudo-labels under domain shift. As aforementioned, the domain shift between labeled and unlabeled source domains deteriorates the accuracy of pseudo-labels. In Figure 2, we visualize the feature distribution of the DANN [4] model trained via fully supervised learning on PACS. As observed, the experimental result shows that samples have been well mapped to their categories, whereas inside a typical class, features from different domains intra a class are separated. Therefore, to obtain more accurate pseudo-labels, we iteratively maintain the average feature of the most confident unlabeled samples for each class of each domain in the memory, which is used as *domain-aware class representation*. Afterwards, when assigning pseudo-labels to unlabeled samples, we combine the output probability of the classifier with the similarity to its class representation to decide which class it belongs to.

Secondly, in SSDG, our goal is to improve the generalization ability of the network to adapt to an arbitrary target domain. However, intuitively, a generalizable model could underfit the unlabeled training domains, thus the accuracy of pseudo-labels could decrease. Considering inconsistent goals between generalization and pseudo-labeling—former prevents to overfit source domains while latter might overfit specific domain for high accuracy, we propose to use a dual-classifier to avoid the possible accuracy degradation of pseudo-labels, which leverages the independent classifiers for joint pseudo-label assignment and do-

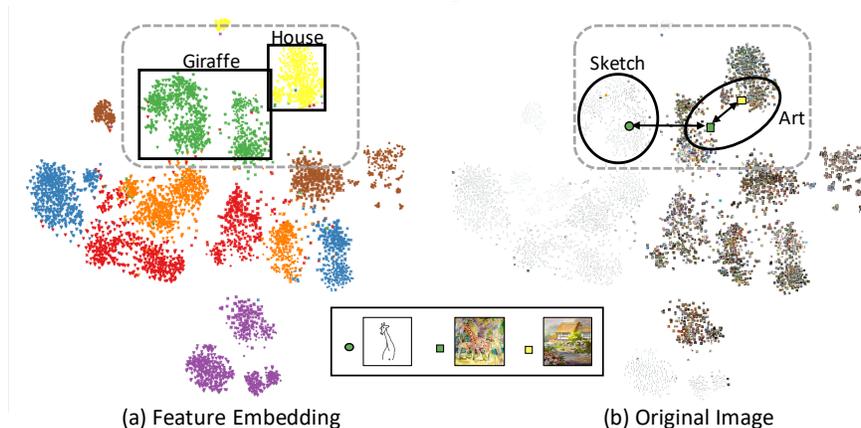


Figure 2: (a): Visualization of feature embeddings from a fully supervised DG model using a domain discriminator to reduce the domain gap on PACS. Note that different colors denote different classes. (b): The original images from PACS, which correspond to features in (a) from the position view. We focus on giraffe and house, Sketch domain, and Art painting domain. As seen, after feature alignment, features are still not strictly domain-invariant, *e.g.*, several giraffe samples from the Art are closer to some house samples from the same domain than a giraffe from the Sketch.

main generalization. In our dual-classifier network, the two branches are trained with different objective functions but a shared feature extractor.

Our contributions can be summarized as follows:

- We propose an effective framework that can be trained in an end-to-end manner to obtain more accurate pseudo-labels of unlabeled data for the semi-supervised domain generalization task.
- We develop the domain-aware pseudo-labeling module to handle the domain shift during generating pseudo-labels. Also, the dual-classifier is proposed to mitigate the conflict between the DG task and the pseudo-label generation.
- Extensive experiments on benchmark datasets, *i.e.*, PACS, OfficeHome, miniDomainNet and VLCS, show the effectiveness of our method compared with several baselines and the state-of-the-art methods.

## 2. Related Works

We review the recent work about unsupervised domain adaptation, domain generalization, semi-supervised learning, and semi-supervised domain generalization.

### 2.1. Unsupervised Domain Adaptation

Unsupervised domain adaptation can effectively transfer knowledge from an annotated source domain to an unlabeled target domain. Mainstream approaches include discrepancy-based [5, 6, 7, 8, 9, 10], adversarial-based [11, 12, 13, 14, 15] and pseudo-labeling-based [16, 17] methods. Lee *et al.* [7] design sliced Wasserstein discrepancy (SWD) to capture the discrepancy between the outputs of task-specific classifiers. Li *et al.* [9] propose maximum density divergence (MDD) to measure the distribution divergence and apply MDD to minimize the inter-domain discrepancy and maximize the intra-class density. Zhang *et al.* [14] introduce a novel Hybrid Adversarial Network (HAN), which achieves a joint adversarial learning with class information and domain alignment. Zhang *et al.* [17] propose to increase the robustness of the model by incorporating high-confidence samples from the target domain. Li *et al.* [15] propose a more practical UDA setting where either the source data or the target data are unknown, and handle the UDA setting by the adversarial attack. In addition, some recent methods aim to address the issue of limited computing power in UDA problems. Li *et al.* [18] propose the Faster Domain Adaptation (FDA) protocol to accelerate unsupervised domain adaptation. Despite UDA being related to DG, UDA has access to the target domain while DG cannot observe the target domain during training.

### 2.2. Domain Generalization

Domain generalization methods can be substantially categorized into data-based methods, feature-based methods, and learning strategy-based methods [19]. The data-based methods aim to generate virtual training data for a more generalizable model, *e.g.*, the methods in [20, 21, 22, 23, 24] enlarge the training

set by image generation and data augmentation techniques which are applied to solve data insufficiency[25]. Feature-based solutions [26, 27, 28, 29] extract domain-agnostic representations on multi-source domains. Another promising technique for domain generalization is meta-learning, such as [30, 31, 32]. Besides, some methods based on other learning strategies (*e.g.*, self-supervision, Ensemble learning) are proposed to obtain the generalizable model, including [33, 34].

### 2.3. Semi-supervised Learning

Current semi-supervised learning methods can be roughly classified into three categories, *i.e.*, entropy regularization based methods, pseudo-label based methods, and consistency regularization based methods. The essence of all these three categories is to force a low-density distribution between different classes [35]. A straightforward way is to add a loss term to directly and explicitly reduce the entropy of the predictions on unlabeled data. Entropy regularization [36] encourages a confident prediction on unlabeled data by minimizing the entropy of the predictions of unlabeled data. Pseudo-label based methods [3, 37, 38] assign approximate classes to unlabeled samples by the inference of the model trained on labeled samples. Consistency regularization shows great success more recently, which includes  $\pi$ -Model [39], Temporal Ensembling [39] and Mean Teacher [40], etc. Besides, a series of holistic approaches to semi-supervised learning have obtained state-of-the-art performance on commonly-studied SSL benchmarks recently. Unsupervised Data Augmentation (UDA) [41] improves consistency loss by substituting simple noising operations with advanced data augmentation, such as RandAugment [42]. MixMatch [43] unifies the existing data augmentation, pseudo-labeling, and mixup to achieve both consistency regularization and entropy regularization. ReMixMatch [44] further improves MixMatch [43]. FeatMatch [45] applies learned feature-based augmentation to consistency loss. FixMatch [46] inherits UDA and ReMixMatch, combines pseudo-labeling and consistency regularization, and finally obtains good performance on SSL benchmarks. Differently, in our SSDG, there is a

data-distribution discrepancy between labeled and unlabeled training data, thus these typical SSL methods could not effectively handle the issue.

#### 2.4. Semi-supervised Domain Generalization

To the best of our knowledge, only very a few works have been proposed for semi-supervised domain generalization problem. DGSML [47] and StyleMatch [48] tackle a new setting in domain generalization problem, where the labeled samples in each domain are not fully labeled. Although we both assign pseudo-labels to unlabeled data, we solve two different scenarios and the challenges we face are totally different. DSDGN [49] solves a semi-supervised domain generalization problem which is similar to us. It applies a Wasserstein generative adversarial network with gradient penalty based adversarial training framework to align feature embedding, and simply adopts the original pseudo-labeling method for unlabeled data. However, this method does not consider the domain shift during pseudo-labeling.

### 3. Method

Unlike the supervised domain generalization (DG) setting, as aforementioned, semi-supervised DG further alleviates the fully-labeled requirement and allows several source domains to be totally unlabeled during training. Formally, we now provide the notations used in our setting. In SSDG, assume we have one labeled source domain in the labeled domain set  $\mathcal{S}_l = \{D^l\}$ , and  $n$  unlabeled source domains in the unlabeled domain set  $\mathcal{S}_u = \{D_1^u, \dots, D_n^u\}$ , and one target domain  $D_t$ . Note that,  $D_t$  is not used in the training process. A training sample in the labeled source domain can be represented as a raw input  $x$ , a semantic label  $y$ , and a domain label  $z$ . Assuming that the number of the training samples in the labeled source domain is  $n_l$ , it can be denoted as  $D^l = \{(x_i^l, y_i^l, z_i^l = 0)\}_{i=1}^{n_l}$ . And an unlabeled domain  $D_j$  can be represented as  $D_j = \{(x_i^u, z_i^u = j)\}_{i=1}^{n_j}$ , when the number of the training samples in  $D_j$  is  $n_j$ .  $C$  stands for the number of categories in the classification dataset. The shared

feature extractor, the predictive classifier, the generalizable classifier, and the domain classifier are denoted by  $F_g$ ,  $F_c$ ,  $F_m$ , and  $F_d$ , respectively. The overview of our method is illustrated in Figure 3.

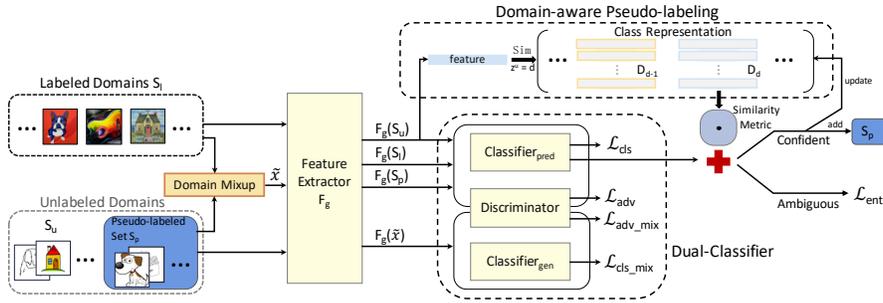


Figure 3: Illustration of our framework, which mainly consists of the feature extractor, domain-aware pseudo-labeling, dual-classifier and domain discriminator. During training, model optimization and pseudo-label prediction are alternate and iterative.

### 3.1. Domain-aware Pseudo-labeling

In the cross-domain scenario, a mixture of samples from all domains are fed into a shared classifier together. However, the data distributions of different domains are significantly different. Accordingly, for different domains, the discriminative characteristics that are critical to classification could be different. For instance, on the PACS dataset, images in *photo* domain are color-specified, while images in *cartoon* domain are not. In *sketch* domain, the color information is totally erased, which makes samples in this domain even harder to distinguish. Thus, even we have applied a discriminator to align features from different domains, the features haven't been perfectly aligned yet. Due to the large and unpredictable domain gap between different domains, the classifier that well fits the labeled domain will generate the poor pseudo-label for unlabeled domains.

In order to alleviate the bias caused by the aforementioned domain gap, we propose domain-aware pseudo-labeling module in our framework. In particular, we first yield the *domain-aware class representation* for each class of each unlabeled domain, which indicates the mean feature of the most highly confident

samples for each class in each domain. Then, when generating the pseudo-label for each unlabeled sample by integrating 1) its predicted probability from the shared predictive classifier and 2) its largest similarity to the class representation from its domain, we can obtain the modified probability for the more reliable pseudo-label.

For an unlabeled sample  $(x^u, z^u)$ , if the domain label  $z^u$  is equal to  $d$ , it means that this sample is from the  $d^{th}$  domain. Using  $M^d \in \mathbb{R}^{D \times C}$  to denote the matrix by gathering the  $D$ -dimensional class representation from total  $C$  classes in the  $d^{th}$  domain.  $\mathbf{sim}(\cdot, \cdot)$  is a similarity measurement function. The similarity is conducted by calculating the cosine similarity between  $F_g(x^u)$  and the class representation of each class, where  $F_g(x^u) \in \mathbb{R}^{1 \times D}$  is the  $D$ -dimensional feature embedding.  $\psi(x^u) \in \mathbb{R}^{1 \times C}$  is the  $C$ -dimension similarity vector obtained after a softmax function. Then we modify the predicted probability  $q(x^u)$  by a correction term to form  $s(x^u)$ . This above process can be formulated as follows:

$$s(x^u) = \gamma q(x^u) + (1 - \gamma)\psi(x^u), \quad (1)$$

where  $q(x^u) = F_c(F_g(x^u))$ ,  $\psi(x^u) = \mathbf{sim}(F_g(x^u), M^d)$ .

Now the pseudo-label  $\widehat{y}^u$  is assigned by  $s(x^u)$  as same as the conventional pseudo-labeling way:

$$\widehat{y}_i^u = \begin{cases} 1 & \text{if } i = \operatorname{argmax}_{i'} s_{i'}(x^u) \text{ and } s_{i'}(x^u) > \delta \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

where  $\delta$  is a threshold. If the final confidence of one sample is larger than the threshold, the pseudo-label is reliable, otherwise it is ambiguous.

The group of representation is updated at each epoch and used at the next epoch. We propose two policies for the domain-aware class representation production. The simplest way is to select the sample with the highest confidence predicted by the classifier as class representation at each epoch. However, the mislabeled sample could arise an accumulation of prediction error at the next epoch [50]. Furthermore, we propose to calculate an ensemble representation with some reliable samples of each class in each unlabeled domain. Specifically,

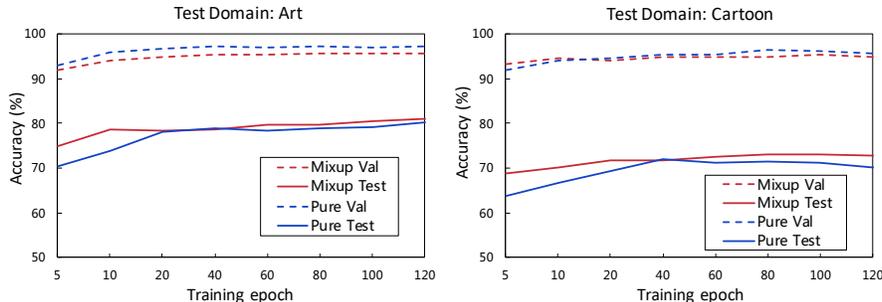


Figure 4: Validation and test accuracy with mixup samples or with pure samples on PACS. Left: Train on Photo, Cartoon and Sketch and test on Art painting. Right: Train on Photo, Art painting and Sketch and test on Cartoon.

we maintain a list of the highest confident samples for each class. We append one sample to the list only if its confidence is higher than the existing highest confidence in the list. And the capacity of each list is  $k$ , we stop appending samples if the size of the list reaches the limit. At the start of each epoch, we calculate the average of the samples (feature embeddings) in the list as class representation. In our experiments  $k$  is set as 100. The comparison of the two schemes will be shown in our experiment.

### 3.2. Dual-classifier

As known, overfitting empirically occurs when a model begins to fit the domain-specific characteristic in the training data rather than learning to generalize from a trend [51]. Consequently, the overfitting issue finally causes performance degradation on unseen test data. This empirical knowledge inspires us that a generalizable model which works well in domain generalization could be non-overfitting for the source training domains.

We train a model with domain mixup [52] and without domain mixup in a supervised way and draw the comparison of both test and validation accuracy. As shown in Figure 4, by domain mixup, the model performs more accurately in the test domain, indicating that the model is more generalizable. However,

the accuracy in the training domains drops, which is displayed by the validation accuracy.

In our approach, we apply cross-domain mixup towards a more generalizable classifier. However, we observe that the generalizable DG classifier is not good at pseudo-labeling. Thus, we apply an auxiliary classifier to achieve the pseudo-labeling, and the classifier for giving the pseudo-labels is trained only by pure samples from source domains. Therefore, we utilize dual-classifier architecture to reduce the conflict between the DG task and the pseudo-label production.

### 3.3. Mining the Knowledge of Unlabeled Domains

When accurate pseudo-labels on unlabeled samples are generated, domain mixup is applied to confident samples with their pseudo-labels and the raw labeled data. For ambiguous samples, since pseudo-labels are not be assigned, the entropy loss is applied to make full use of these samples.

**Confident unlabeled samples.** As inspired by [52], we innovatively interpolate between a labeled domain and an unlabeled domain to achieve inter-domain data augmentation. Our intuition is that the inter-domain samples generated by domain mixup can boost the generalization of networks by introducing additional training domains, which has been verified as an important technique in domain generalization [23, 21, 22]. Assuming that we have already assigned a pseudo-label for an unlabeled sample  $(x^u, z^u)$ , forming  $(x^u, \widehat{y}^u, z^u)$ . And we have a sample from a labeled domain  $(x^l, y^l, z^l)$ .  $\widehat{y}^u$ ,  $z^u$ ,  $y^l$ , and  $z^l$  are all one-hot vectors. Then the operation is formulated as below:

$$\tilde{x} = \lambda x^l + (1 - \lambda)x^u, \quad (3)$$

$$\tilde{y} = \lambda y^l + (1 - \lambda)\widehat{y}^u, \quad (4)$$

where  $\lambda \sim \text{Beta}(\alpha, \alpha)$ , for  $\alpha > 0$ , and  $\lambda \in [0, 1]$ . The hyper-parameter  $\alpha$  controls the strength of interpolation.

Additionally, to further enhance the generalization ability, we also apply mixup to domain labels of labeled and unlabeled samples for training domain

discriminator as follows:

$$\tilde{z} = \lambda z^l + (1 - \lambda)z^u, \quad (5)$$

where  $\lambda$  is same with that used in Eqn. (3) and Eqn. (4). According to these steps, we could finally generate a virtual sample  $(\tilde{x}, \tilde{y}, \tilde{z})$ .

**Ambiguous unlabeled samples.** For the unlabeled samples with low-confidence prediction, since we are unclear about their real labels, it is hard to generate mixed samples by them. In order to further improve the generalization ability of our method by fully leveraging these unlabeled samples, we employ entropy loss to encourage unlabeled samples to be classified into a specific category:

$$L_{ent} = \frac{1}{N_u} \sum_{i=1}^{N_u} H(F_c(F_g(x_i^u))), \quad (6)$$

where  $H(\cdot)$  is the entropy function. By introducing  $L_{ent}$ , the networks are forced to make more confident predictions on ambiguous unlabeled samples.

#### 3.4. Training Procedure

Formally, as aforementioned, we denote  $\mathcal{S}_l$  as the set of all labeled samples,  $\mathcal{S}_u$  as the set of all unlabeled samples. During training, once we select the reliable unlabeled samples with high confident pseudo-labels, forming  $\mathcal{S}_w$ , we move them to the set of samples with pseudo-labels which is denoted as  $\mathcal{S}_p$ . At each training epoch, we assume,  $|\mathcal{S}_l| = N_l$ ,  $|\mathcal{S}_u| = N_u$  and  $|\mathcal{S}_p| = N_p$ . The total number of all training samples is denoted as  $N = N_l + N_u + N_p$  (refer to Algorithm 1).

For training our network, we apply 1) a classification loss  $L_{cls}$  to  $\mathcal{S}_l$  and  $\mathcal{S}_p$  since they are with labels or predicted pseudo-labels, and 2) an adversarial loss  $L_{adv}$  to all sets of samples, *i.e.*,  $\mathcal{S}_l$ ,  $\mathcal{S}_p$ , and  $\mathcal{S}_u$ .

$$L_{cls} = \frac{1}{N_l} \sum_{i=1}^{N_l} \ell(F_c(F_g(x_i^l)), y_i^l) + \frac{1}{N_p} \sum_{i=1}^{N_p} \ell(F_c(F_g(x_i^p)), \hat{y}_i^p), \quad (7)$$

$$L_{adv} = \frac{1}{N} \left( \sum_{i=1}^{N_l} \ell(F_d(F_g(x_i^l)), z_i^l) + \sum_{i=1}^{N_u} \ell(F_d(F_g(x_i^u)), z_i^u) \right)$$

---

**Algorithm 1: Training Process**

---

**Input:** Labeled source  $\mathcal{S}_l$ , unlabeled source  $\mathcal{S}_u$

**Output:** Generalizable model  $F_g$  and  $F_m$

- 1 **Initialize** Networks, Pseudo-label Set  $\mathcal{S}_p = \emptyset$ , Class representation list  $L = [none] \times C$  for each domain;
  - 2 **while** *not end of epoch* **do**
    - 3  $\mathcal{S}_m \leftarrow$  Perform domain mixup on  $\mathcal{S}_l$  and  $\mathcal{S}_p$ ;
    - 4 Training the model using Eqns. (11) and (12);
    - 5 Inference the model on  $\mathcal{S}_u$ , and obtain  $q(x^u)$ ;
    - 6 **if** *no none in L* **then**
      - 7  $s(x^u) \leftarrow$  Calculated by Eqn. (1);
    - 8 **end**
    - 9 Update class representation by  $q(x^u)$ ;
    - 10 Assign pseudo-labels by  $s(x^u)$ ;
    - 11 Recognize confident and ambiguous samples by Eqn. (2), and the confident set is denoted as  $\mathcal{S}_{u'}$ ;
    - 12 Update  $\mathcal{S}_u \leftarrow \mathcal{S}_u - \mathcal{S}_{u'}$ ,  $\mathcal{S}_p \leftarrow \mathcal{S}_p \cup \mathcal{S}_{u'}$ ;
  - 13 **end**
- 

$$+ \sum_{i=1}^{N_p} \ell(F_d(F_g(x_i^p)), z_i^p), \quad (8)$$

where  $\ell(\cdot, \cdot)$  is the cross-entropy loss. The loss on samples mixed up by  $\mathcal{S}_p$  and  $\mathcal{S}_l$  is defined as:

$$L_{cls-mix} = \frac{1}{N_l} \sum_{i=1}^{N_l} \sum_{j=1}^{N_p} \ell(F_m(F_g(\tilde{x})), \tilde{y}), \quad (9)$$

$$L_{adv-mix} = \frac{1}{N_l} \sum_{i=1}^{N_l} \sum_{j=1}^{N_p} \ell(F_d(F_g(\tilde{x})), \tilde{z}). \quad (10)$$

The training objective of our semi-supervised DG model can be described as

follows:

$$\min_{F_g, F_c, F_m} L_{cls} + L_{cls\_mix} + w^t(-L_{adv} - L_{adv\_mix} + L_{ent}), \quad (11)$$

$$\min_{F_d} L_{adv} + L_{adv\_mix}, \quad (12)$$

where the weight function  $w^t$  ramps up from zero to one during the training procedure. The whole training procedure is described in Algorithm 1.

## 4. Experiments

### 4.1. Experimental Setting

**Datasets.** To evaluate the effectiveness of our method for the semi-supervised domain generalization, we conduct several experiments on four benchmark datasets. **PACS** [53] contains 7 categories of images from 4 domains (Photo, Art painting, Cartoon and Sketch). **OfficeHome** [54] consists of images from 4 different domains (Artistic, Clip art, Product and Real-World). For each domain, this dataset involves 65 object categories found typically in office and home. To verify that our method also has promotion when data is abundant, we validate our method on **miniDomainNet**. It is a subset of DomainNet aggregated by [55], which contains 140,006 images from 4 domains (Clipart, Painting, Real, and Sketch), covering 126 classes in the raw dataset. **VLCS** [56] includes five categories of images from four domains (*i.e.*, Caltech 101, PASCAL VOC, LabelMe and SUN09). For each dataset, we select two domains as the unlabeled source domains, one domain as the labeled source domain, and leave the remaining one for test. We test all 12 combinations of domains on each dataset and report the average accuracy. To ensure a fair comparison with other methods, all experiments utilize the same scheme of data division.

**Implementation Details.** In all experiments, we use ResNet [57] as the backbone, and we start with a pre-trained model and fine-tune on source domains with the batch size of 128. Since the test domain is unavailable during training. For PACS and OfficeHome, we train 120 epochs and select the model

obtained in the final epoch for test. For miniDomainNet, the network converges at an earlier epoch due to the huge amount of data. So we just train 60 epochs and save the final model. We apply a SGD optimizer with momentum 0.9, and we set the initial learning rate to 1e-3 in the case of PACS, OfficeHome and miniDomainNet, and 1e-4 in the case of VLCS. The learning rate is divided by 10 at the 30-th and the 50-th epochs. All experiments are implemented in Pytorch with  $4 \times 11$  GB RTX 2080Ti GPUs.

#### 4.2. Comparison with Other Methods

On PACS, we compare our method with the following baselines and the state-of-the-art semi-supervised learning approaches (*i.e.*, FixMatch [46], FeatMatch [45] and AdaMatch [58]), domain adaptation (DA) approaches (*i.e.*, SymNets [59], SRDC [60], CGDM [61], ATDOC [62] and FixBi [63]), domain generalization (DG) methods (*i.e.*, L2D [64], JiGen [65] and RSC [66]) and ‘DA+DG’ methods on the classification accuracy of the target domain using ResNet-18 and ResNet-50. And we also compare ours with the baselines and the state-of-the-art methods on OfficeHome and miniDomainNet using ResNet-18. Finally, we report our performance, the baseline results on VLCS using ResNet-18.

- **Baseline**

- **SupOne:** Train a plain model on one labeled source domain in a supervised way and the other unlabeled source domains are not used.
- **DSDGN:** Implement the semi-supervised DG method proposed in [49] using the same network structure as our framework.

- **DA:** The labeled source domain is used as the source domain in the DA methods, and the mixture of the two unlabeled source domains is used as the target domain. Finally, the unseen target domain is used for testing.
- **DG:** For single-DG methods, *i.e.* L2D, the unlabeled domains are not used. For DG methods, *i.e.* JiGen and RSC, unlabeled samples are used during training.

Table 1: Experimental results of accuracy (%) on PACS based on ResNet-18 and ResNet-50. The title in the first row indicates the name of the target domain and the title in the second row is the name of the labeled domain among training source domains. **P, A, C, S** are Photo, Art painting, Cartoon, Sketch, respectively. Note that the best performance is in **bold**.

Method		Photo			Art Painting			Cartoon			Sketch			Avg.
		A	C	S	P	C	S	P	A	S	P	A	C	
ResNet-18														
Baseline	SupOne	95.69	86.17	40.66	64.60	68.75	25.93	25.90	58.53	38.10	32.55	49.05	61.29	53.94
	DSDGN [67]	<b>96.29</b>	88.38	36.94	<b>66.46</b>	70.36	26.12	47.27	67.96	50.64	46.35	59.35	61.77	59.82
SSL	FixMatch [46]	95.39	85.63	59.40	65.38	68.41	51.37	45.22	62.20	55.55	53.32	61.52	<b>76.94</b>	65.03
	FeatMatch [45]	95.33	82.57	56.89	66.06	72.02	58.69	47.10	65.57	57.30	64.39	72.59	74.40	67.74
	AdaMatch [58]	94.61	49.64	41.20	69.43	<b>81.01</b>	42.48	64.12	65.53	60.28	54.38	68.30	76.27	63.94
DG	L2D [64]	95.51	86.65	47.25	64.75	73.78	49.95	38.35	68.77	63.31	41.23	67.24	70.20	63.92
	JiGen [65]	95.39	83.53	47.43	66.26	69.53	33.54	34.68	66.98	54.18	40.83	57.78	62.61	59.40
	RSC [66]	86.23	84.79	46.05	63.38	68.75	33.15	<b>66.42</b>	47.06	57.34	<b>66.15</b>	70.83	60.09	62.52
DA	SymNets [59]	88.42	76.44	28.36	61.80	48.27	31.86	30.21	<b>74.24</b>	45.69	16.36	53.41	58.70	51.15
	SRDC [60]	91.50	80.36	38.32	55.37	71.68	31.40	50.43	69.28	54.73	21.41	54.75	60.15	56.62
	CGDM [61]	95.15	75.69	47.01	62.26	64.86	31.69	39.38	59.56	50.73	12.57	41.84	59.86	53.38
	ATDOC [62]	90.06	78.86	<b>79.52</b>	56.20	53.52	46.78	57.68	60.54	55.16	46.63	51.01	42.30	59.86
	FixBi [63]	87.90	85.75	48.02	45.41	68.51	47.95	39.38	62.07	50.51	27.46	40.85	52.97	54.73
DA+DG	SymNets+RSC	92.28	<b>92.46</b>	40.90	69.46	72.87	46.70	60.96	30.59	51.28	63.78	<b>73.05</b>	18.91	59.44
	SRDC+RSC	92.46	91.74	50.36	66.36	70.61	<b>65.87</b>	58.23	57.21	<b>65.49</b>	64.98	66.38	62.26	67.66
	CGDM+JiGen	93.87	91.66	58.54	65.75	75.17	58.86	46.01	72.54	54.54	61.40	68.57	61.98	67.41
	ATDOC+RSC	92.81	73.41	77.01	<b>70.46</b>	61.28	59.81	60.20	65.61	64.29	49.55	65.64	72.18	67.69
	FixBi+RSC	92.99	87.31	50.48	66.46	73.24	57.67	40.19	62.63	53.67	26.22	60.96	70.25	61.84
Ours	94.37	91.02	66.53	69.92	75.68	55.37	54.22	71.46	57.94	65.69	71.27	71.06	<b>70.38</b>	
ResNet-50														
Baseline	SupOne	<b>98.14</b>	86.71	36.89	73.39	72.46	30.37	34.47	65.87	45.86	34.03	56.76	68.39	58.61
	DSDGN [67]	97.90	93.65	38.62	71.88	79.15	39.60	45.95	70.78	53.11	41.46	68.36	70.12	64.22
SSL	FixMatch [46]	95.59	94.07	58.04	68.16	85.28	55.90	54.47	76.58	66.55	46.09	73.40	77.70	70.99
	FeatMatch [45]	97.84	93.71	52.22	74.37	86.87	58.74	66.13	71.08	65.78	<b>71.54</b>	73.84	<b>81.04</b>	74.43
	AdaMatch [58]	96.95	93.29	50.84	71.29	<b>88.13</b>	48.39	67.02	66.17	59.98	62.65	67.54	77.32	70.80
DG	L2D [64]	97.13	92.28	52.04	73.63	78.08	52.10	43.30	74.40	69.75	44.08	66.48	76.18	68.29
	JiGen [65]	97.31	89.76	53.23	72.27	76.27	42.04	42.28	70.01	64.59	45.71	59.05	69.00	65.13
	RSC [66]	<b>91.56</b>	88.80	59.10	64.31	68.70	34.42	<b>68.56</b>	47.78	65.49	65.46	74.12	62.13	65.87
DA	SymNets [59]	91.52	75.13	52.99	64.77	60.84	38.97	38.67	67.41	51.49	29.55	56.71	62.41	57.54
	SRDC [60]	92.52	93.89	31.74	65.14	77.98	42.63	44.45	72.40	57.85	24.33	53.09	60.22	59.69
	CGDM [61]	97.37	76.35	69.46	72.27	64.01	54.10	51.02	62.03	57.51	28.94	43.15	63.60	61.65
	ATDOC [62]	92.87	82.81	88.38	39.70	76.07	55.86	55.50	68.00	57.68	38.69	52.00	53.53	63.42
	FixBi [63]	90.24	92.10	55.27	64.21	74.51	52.05	37.59	70.99	65.10	45.48	54.57	60.88	63.58
DA+DG	SymNets+RSC	94.79	94.56	40.27	77.87	74.52	50.18	61.12	44.06	54.93	62.24	<b>76.99</b>	30.25	63.48
	SRDC+RSC	95.93	<b>95.03</b>	64.25	75.39	78.93	60.40	59.60	74.27	67.41	58.18	72.10	74.40	72.99
	CGDM+JiGen	97.60	94.91	73.03	78.02	83.89	73.44	49.87	64.97	63.35	57.04	64.72	73.63	72.87
	ATDOC+RSC	95.57	80.30	<b>90.06</b>	62.30	81.15	<b>76.81</b>	58.32	72.99	<b>74.87</b>	51.06	71.06	78.37	74.41
	FixBi+RSC	95.99	93.35	54.79	69.97	77.88	55.57	35.07	70.99	67.41	67.29	71.37	75.06	69.56
Ours	97.61	93.53	66.35	<b>78.03</b>	86.98	62.45	59.17	<b>76.88</b>	69.37	65.61	74.09	78.52	<b>75.72</b>	

- **DA+DG:** We consider the labeled domain as the source domain and other unlabeled domains as the unlabeled target domain. Then we use a UDA method to generate pseudo-labels for unlabeled domains. Finally, both the fully-labeled domain and pseudo-labeled unlabeled domains are fed into a supervised DG method to train the final model.

Table 1 displays the results on PACS. Compared with the baselines, our method achieves outstanding performance by significant margins with both smaller and larger network architectures. Specifically, our method improves the performance of SupOne by +16.44% and +17.11% with ResNet-18 and ResNet-50, respectively. This shows that ours significantly improves the performance

by leveraging unlabeled data. Compared with DSDGN, ours gains +10.56% with ResNet-18 and gains +11.50% with ResNet-50. The results indicate that our method makes more efficient use of unlabeled samples, thus our proposed domain-aware pseudo-labeling and other modules outperform the naive pseudo-labeling method in DSDGN.

Compared with the recent SOTA semi-supervised methods, ours also shows priority in the SSDG task. Specifically, ours outperforms the best SSL methods by +2.64% using ResNet-18 and +1.29% using ResNet-50. This shows that conventional semi-supervised methods are not superior in solving SSDG problems. By considering the domain shift between the labeled source domain and unlabeled source domains, our proposed method can achieve better performance in a clear margin compared with semi-supervised approaches.

From Table 1 we also observe that using ResNet-18 as a backbone architecture, ours outperforms all DA methods (*i.e.*, SymNets, SRDC, CGDM, AT-DOC and FixBi) by large margins: +19.23%, +13.76%, +17.00%, +10.52% and +15.65%. Ours is also clearly better than DA methods when ResNet-50 is employed as a backbone architecture. This main reason is that these DA methods do not address the domain gap between source domains and the unseen target domain because the unseen target domain cannot be employed during training in the SSDG task.

Additionally, as seen in Table 1, our method achieves better performance than both single-DG and DG methods. For example, our method improves L2D by +6.46% using ResNet-18 and +7.43% using ResNet-50. These results show that our method improves the performance of single-DG with free unlabeled data. Ours also outperforms the best DG methods by +7.86% using ResNet-18 and +9.85% using ResNet-50. These results imply that our method makes more efficient use of unlabeled samples, *i.e.*, the proposed domain-aware pseudo-labeling module and the proposed dual-classifier surpass “self-challenging” mechanism (RSC) and “solving a jigsaw puzzle” task (JiGen) in utilizing unlabeled samples.

Moreover, we can observe in Table 1 that compared with the SOTA ‘DA+DG’

Table 2: Experimental results of accuracy (%) on OfficeHome. The title in the first row indicates the name of the target domain and the title in the second row is the name of the labeled domain in training source. **A**, **C**, **P**, **R** stand for Art, Clipart, Product, Real, respectively. The best performance is **bold**.

Method		Art			Clipart			Product			Real			Avg.
		C	P	R	A	P	R	A	C	R	A	C	P	
Baseline	SupOne	43.63	38.20	54.84	38.95	37.34	44.10	54.88	54.02	70.94	64.08	57.54	62.98	51.79
	DSDGN [67]	45.49	41.37	56.49	38.40	37.46	43.98	55.37	54.99	72.11	64.56	55.93	63.44	52.47
SSL	FixMatch [46]	44.53	42.14	58.11	43.42	42.09	45.67	56.24	56.39	70.80	65.38	56.43	<b>65.71</b>	53.91
	FeatMatch [45]	49.98	47.34	55.67	41.95	42.96	47.33	54.11	53.71	68.19	61.12	54.99	64.33	53.47
DG	L2D [64]	42.14	38.25	56.72	40.75	41.42	<b>51.33</b>	47.04	51.78	69.90	56.90	54.57	63.58	51.20
	JiGen [65]	32.06	30.49	42.85	35.51	34.71	41.70	42.60	47.38	62.78	52.51	49.28	54.26	43.84
	RSC [66]	40.08	37.19	53.99	39.56	38.37	45.43	48.13	57.12	69.45	63.46	52.05	60.80	50.47
DA	SymNets [59]	39.64	45.20	53.65	33.29	32.30	37.96	50.96	48.52	68.17	58.73	51.11	60.12	48.30
	SRDC [60]	<b>50.31</b>	<b>52.37</b>	<b>58.96</b>	37.18	37.27	38.24	58.68	57.87	68.35	61.26	59.08	64.61	53.68
	CGDM [61]	46.14	44.91	51.92	38.01	37.69	42.97	57.22	58.43	69.43	64.86	59.63	63.69	52.91
DA+DG	SymNets+RSC	30.74	37.54	51.76	33.24	39.25	43.76	42.69	46.18	65.71	49.02	41.27	60.16	45.11
	SRDC+RSC	43.51	45.74	49.40	39.43	<b>43.30</b>	39.40	54.04	55.64	69.09	56.92	55.73	63.35	51.30
	CGDM+JiGen	49.36	48.66	53.59	41.28	39.52	43.43	<b>60.35</b>	<b>61.63</b>	71.09	63.21	<b>61.35</b>	64.03	54.79
	Ours	47.55	46.07	58.01	<b>44.33</b>	42.34	47.90	57.56	57.83	<b>72.43</b>	<b>65.48</b>	59.74	65.09	<b>55.36</b>

methods using ResNet-18 as backbone, ours shows advantages. Specifically, compared with ‘SymNets+RSC’, our method improves the average accuracy by +10.94%. Compared with ‘SRDC+RSC’, ours gains +2.72%, and compared with ‘CGDM+JiGen’, there is an improvement of +2.97%. Compared with ‘ATDOC+RSC’, ours improves +2.69%, and compared with ‘FixBi+RSC’, ours gains +8.54%. Moreover, using ResNet-50 as backbone, compared with ‘SymNets+RSC’, ‘SRDC+RSC’, ‘CGDM+JiGen’, ‘ATDOC+RSC’ and ‘FixBi+RSC’, ours improves the accuracy by +12.24%, +2.73%, +2.85%, +1.31% and +6.16%, respectively. The reason why ‘DA+DG’ methods are not superior to ours is that ‘DA+DG’ methods are not an end-to-end deep framework, thus the labeling process of UDA methods cannot be improved by the process of DG, which means that the performance of DG is largely influenced by that of UDA. Particularly, when the UDA model is unreliable, the performance of the whole model would become inferior.

Experimental results on OfficeHome are shown in Table 2. It is worth noting that OfficeHome has a relatively smaller domain shift than PACS. As seen in this table, our method gains +3.57% compared with SupOne and improves the average accuracy by +2.89% compared with DSDGN, which thanks to the efficacy of the proposed domain-aware pseudo-labeling scheme and dual-classifier. Furthermore, compared with the SOTA SSL, DG, DA and ‘DA+DG’ methods,

Table 3: Experimental results of accuracy (%) on miniDomainNet. The title in the first row indicates the name of the target domain, and the title in the second row is the name of the labeled domain in source domains. **C**, **P**, **R** and **S** stand for Clipart, Painting, Real and Sketch, respectively. The best performance is **bold**.

Method		Clipart			Painting			Real			Sketch			Avg.
		P	R	S	C	R	S	C	P	S	C	P	R	
Baseline	SupOne	47.45	47.83	53.54	37.91	48.43	42.69	48.53	64.90	47.35	43.21	42.12	42.33	47.19
	DSDGN [67]	47.29	48.16	53.25	39.27	50.88	46.34	44.78	64.46	50.21	42.18	44.50	42.61	47.83
SSL	FixMatch [46]	50.79	56.89	57.46	48.75	60.38	57.18	52.47	63.08	56.54	38.09	42.22	40.20	52.00
	FeatMatch [45]	52.70	52.17	52.06	<b>48.91</b>	53.38	37.05	37.40	<b>66.06</b>	41.14	45.57	<b>49.43</b>	47.17	48.59
DG	L2D [64]	50.23	52.75	53.90	37.47	<b>56.49</b>	42.40	45.69	59.16	45.81	46.31	45.59	45.37	48.43
	JiGen [65]	42.38	46.87	44.27	27.81	43.22	26.23	37.39	54.01	37.44	33.45	30.18	33.06	38.03
	RSC [66]	42.41	48.20	47.84	28.00	52.65	33.07	37.63	56.48	35.49	41.92	42.68	40.59	42.25
DA	SymNets [59]	39.51	50.76	47.73	35.11	52.35	46.10	43.42	58.95	48.48	37.14	33.44	32.21	43.77
	SRDC [60]	46.40	46.64	52.00	42.18	52.11	50.53	55.71	65.49	56.06	41.75	39.91	38.30	48.92
	CGDM [61]	50.31	49.71	56.03	42.64	53.41	48.57	52.27	64.52	54.76	43.39	40.70	42.20	49.88
DA+DG	SymNets+RSC	47.66	48.09	47.24	39.21	51.01	52.28	40.14	52.57	49.68	38.80	40.17	39.04	45.49
	SRDC+RSC	48.74	45.86	55.25	41.12	55.37	52.23	48.23	59.33	52.94	39.73	43.34	43.04	48.77
	CGDM+JiGen	49.79	<b>58.34</b>	56.75	45.20	54.73	51.91	<b>55.76</b>	62.19	<b>58.77</b>	46.96	45.83	46.27	52.71
	Ours	<b>54.29</b>	55.98	<b>58.16</b>	45.36	56.22	<b>53.08</b>	51.51	64.55	55.05	<b>48.70</b>	49.38	<b>49.91</b>	<b>53.52</b>

Table 4: Experimental results of accuracy (%) on VLCS. The title in the first row indicates the name of the target domain and the title in the second row is the name of the labeled domain among training source domains. C, L, V, S stand for CALTECH, LABELME, VOC, SUN, respectively. The best performance is in **bold**.

Method	CALTECH			LABELME			VOC			SUN			Avg.
	L	V	S	C	V	S	C	L	S	C	L	V	
SupOne	68.83	98.02	72.23	53.09	57.27	59.64	44.52	59.87	61.37	35.44	<b>50.24</b>	73.43	61.16
DSDGN[67]	70.18	<b>98.30</b>	<b>74.35</b>	48.01	65.85	59.22	49.01	56.40	61.08	38.73	<b>50.24</b>	70.35	61.81
Ours	<b>79.01</b>	97.39	72.93	<b>65.85</b>	<b>66.08</b>	<b>61.11</b>	<b>51.08</b>	<b>60.38</b>	<b>63.12</b>	<b>39.64</b>	49.30	<b>74.40</b>	<b>65.02</b>

our method increases +1.45%, +4.16%, +1.68% and +0.57%, respectively.

We compare our method with baselines and the SOTA methods on miniDomainNet in Table 3. As seen, our method achieves +6.33% and +5.69% gain on miniDomainNet compared with SupOne and DSDGN, respectively. Besides, our method outperforms all the SOTA methods (*i.e.*, SSL, DG, DA and ‘DA+DG’) by large margins: +1.52%, +5.09%, +3.64% and +0.81%. The results indicate that our method is also effective on the large-scale dataset.

In Table 4 we report the performance of ours and baselines on VLCS. We notice that SupOne achieves an average accuracy of 61.16%, which is only trained on the labeled source domain in a supervised way. DSDGN achieves slight improvement compared with SupOne. Our method increases the average accuracy by +3.86% compared with SupOne on VLCS, which could be attributed

to that we assign accurate pseudo-labels to unlabeled source domains and the inter-domain mixup improves the generalization ability of the model.

### 4.3. Ablation Study

Table 5: Accuracy (%) of ablation study on PACS. **P, A, C, S** stand for Photo, Art painting, Cartoon, Sketch, respectively. Baseline represents our method without both DAPL and DC. The best performance is **bold**.

Method	Photo			Art Painting			Cartoon			Sketch			Avg.
	A	C	S	P	C	S	P	A	S	P	A	C	
Baseline	92.52	87.49	48.80	66.80	72.51	45.51	51.07	67.53	51.83	52.07	53.47	59.36	62.41
Ours w/o DAPL	92.64	88.38	49.40	69.24	69.14	45.02	54.10	68.94	53.41	58.77	64.01	65.20	64.85
Ours w/o DC	93.29	85.69	60.96	66.31	73.68	46.83	<b>56.83</b>	68.13	52.30	54.26	56.05	63.32	64.80
DBSCAN[68]	93.47	85.57	57.78	66.21	70.07	43.07	49.62	67.19	48.81	44.31	54.11	54.42	61.22
Agglomerative[69]	85.81	89.28	62.28	65.82	69.63	52.59	51.07	67.66	50.77	46.02	56.15	56.07	62.76
Ours	<b>94.37</b>	<b>91.02</b>	<b>66.53</b>	<b>69.92</b>	<b>75.68</b>	<b>55.37</b>	54.22	<b>71.46</b>	<b>57.94</b>	<b>65.69</b>	<b>71.27</b>	<b>71.06</b>	<b>70.38</b>

In order to verify the contributions of domain-aware pseudo-labeling (DAPL) and dual-classifier (DC), we conduct an ablation study on PACS, as reported in Table 5. “Ours w/o DAPL” replaces DAPL by naive pseudo-labeling, and “Ours w/o DC” conducts prediction and generalization by the same classifier in the training stage. Without DAPL or DC, the generalization ability of our method dramatically drops. The results show that our modules are crucial to improving the accuracy of the classification on unseen target domains. Besides, we show the pseudo-label accuracy in Figure 5. As seen, our method can obtain more accurate pseudo-labels for unlabeled data compared with the baseline (*i.e.*, our method removes both DAPL and DC). In order to further evaluate the effectiveness of DAPL, we conduct an ablation study to replace DAPL with other state-of-the-art cluster methods ( *i.e.*, DBSCAN [68] and Agglomerative Clustering [69]) for comparison. As seen, ours outperforms the best clustering-based methods by +7.62%. The results indicate that domain-aware pseudo-labeling is more effective compared with clustering methods.

In addition, as mentioned above, we propose to utilize a dual-classifier (a predictive classifier for producing pseudo-labels and a generalizable classifier) to avoid the possible accuracy degradation of pseudo-labels, which leverages the

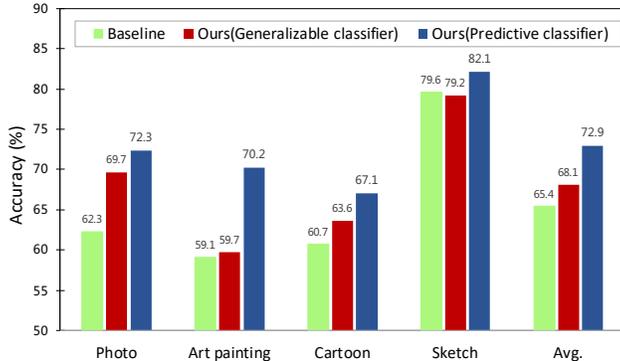


Figure 5: Pseudo-label accuracy (%) on PACS.

independent classifier for joint pseudo-label assignment and domain generalization. The predictive classifier is trained by original samples and the generalizable classifier is trained by mixed samples. We also compare the pseudo-label accuracy of the two classifiers on PACS in Figure 5. As seen, the predictive classifier gains 4.8% in average accuracy compared with the generalizable classifier. This shows that the generalizable model causes a drop in the accuracy of pseudo-labels. And the dual-classifier module is effective for mitigating this problem. It is worth mentioning that our method picks up pseudo labels with an accuracy of over 70% in a 7-class classification task using ResNet-18.

As described in Sec. 3.3, domain mixup is applied to promote the generalization ability of our model on the unseen target domain. And we only choose confident unlabeled samples to mix up with the labeled ones. We study the influence of different designs and show the results in Table 6. As seen, our method shows better performance compared with “w/o mixup” and the model that mixes up all unlabeled samples. We also perform ablation study on the losses in the training objective, *i.e.*,  $L_{cls\_mix}$ ,  $L_{adv}$ ,  $L_{adv\_mix}$  and  $L_{ent}$ . The results show that ours outperforms the scheme without  $L_{cls\_mix}$ ,  $L_{adv}$ ,  $L_{adv\_mix}$  and  $L_{ent}$  by 7.24%, 5.29% , 1.36% and 2.45%, respectively, demonstrating the effectiveness of each loss function. In summary, in these experiments, we confirm

Table 6: Experimental results of accuracy(%) with different designs of Domain Mixup and different loss functions on PACS. The titles indicate the target domains, and each column is the average accuracy of three combinations of source domains (*e.g.*, A, C, S under Photo).

Target	Photo	Art	Cartoon	Sketch	Avg.
w/o Mixup	80.48	63.48	62.88	60.03	66.72
MixupAll	79.50	63.51	59.24	63.16	66.35
w/o $L_{cls.mix}$	77.61	61.63	57.95	55.38	63.14
w/o $L_{adv}$	79.44	65.90	59.88	55.15	65.09
w/o $L_{adv.mix}$	81.66	67.58	61.96	64.87	69.02
w/o $L_{ent}$	82.69	62.44	59.43	67.15	67.93
Ours	83.97	66.99	61.21	69.34	70.38

that the main modules in our framework are useful.

#### 4.4. Further Analysis

**Sensitivity of Hyper-parameters.** Here we discuss the sensitivity to hyper-parameters of our method on PACS, including  $\alpha$  in mixup and  $(\gamma, \delta)$  in domain-aware pseudo-labeling. To simplify the analysis, We select 4 combinations. To be specific, “Photo (Sketch)”, “Art (Photo)”, “Cartoon (Art)” and “Sketch (Cartoon)” are chosen, where “Photo(Sketch)” represents the case that Sketch is labeled for training and Photo is the target. We train our model with  $\alpha$  in [0.1, 0.2, 0.4, 0.8, 1.0]. With the increase of  $\alpha$ , mixup is more likely to generate more confusing samples. The results are reported in Figure 6. The effect of hyperparameter  $\alpha$  on testing accuracy does not show similar trends in each experiment. Photo (Sketch) and Sketch (Cartoon) show large variance. And an optimal value for all experiments is between 0.2 and 0.8. As for  $\gamma$ , we set four appropriate value pairs for the weight  $\gamma$  and the threshold  $\delta$  in domain-aware pseudo-labeling, specifically, (0.05, 0.2), (0.1, 0.24), (0.2, 0.3) and (0.3, 0.36). It can be observed that the performance is not sensitive to  $\gamma$  and smaller value is slightly better. We set  $\alpha$  and  $\gamma$  as 0.2 and 0.1 for all combinations in our experiments.

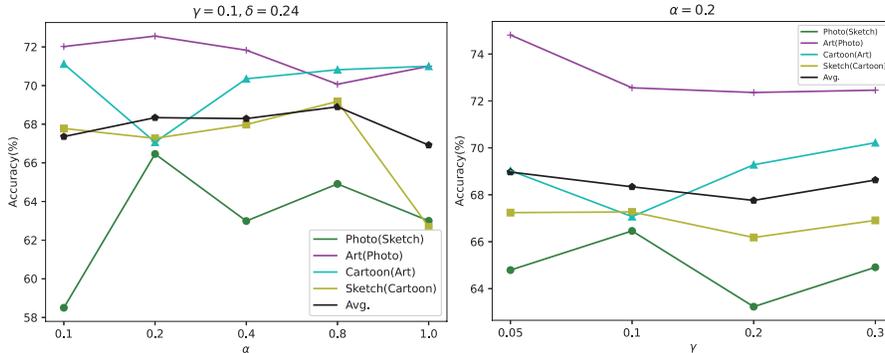


Figure 6: The sensitivity analysis of  $\alpha$  in mixup (left) and  $\gamma$  in domain-aware pseudo-labeling (right) on PACS.

Table 7: Compare different policies of selecting class representation in domain-aware pseudo-labeling.

Target	Photo	Art	Cartoon	Sketch	Avg.
One	83.06	66.57	61.11	67.95	69.67
Ensemble	83.97	66.99	61.21	69.34	70.38

**Different Schemes for Class Representation.** We propose two policies of selecting class representation for domain-aware pseudo-labeling in our framework. The results is shown in Table 7, “One” stands for picking up the most confident one unlabeled sample from each class as class representation at every epoch, and “Ensemble” means calculating the average of several samples as domain-aware class representation. When applying “Ensemble” policy, we choose one sample only if its confidence is higher than the existing highest confidence in the same class. We hold all chosen samples and calculate the average at every epoch. This experiment shows that “Ensemble” policy is more fault-tolerant and its performance is better.

**Visualization of Feature Distributions.** Figure 7 visualizes the feature distributions of three source domains on PACS, including one labeled domain and two unlabeled domains. As seen, the selected domain-aware class representation samples are accurate according to the corresponding original images.

Meanwhile, it can be observed that there are still some misclassified unlabeled samples, especially in Sketch, due to the large domain gap.

## 5. Conclusion

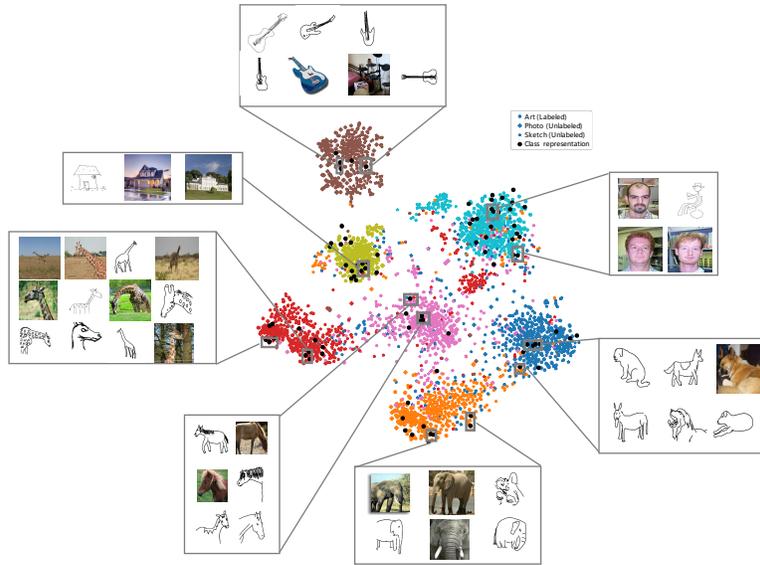
In this paper, we address the problem of semi-supervised domain generalization via producing better pseudo-labels for unlabeled data. Firstly, we propose domain-aware pseudo-labeling for picking up more accurate pseudo labels for unlabeled data by domain-based modification. Then, a dual-classifier network structure is employed to promote the generalization of the model and the accuracy of pseudo-labels. Finally, utilizing the accurate pseudo-labels in unlabeled domains, we apply domain mixup to them and enforce entropy regularization on ambiguous samples. Extensive experiments on benchmark datasets validate the efficacy of our framework.

## Acknowledgement

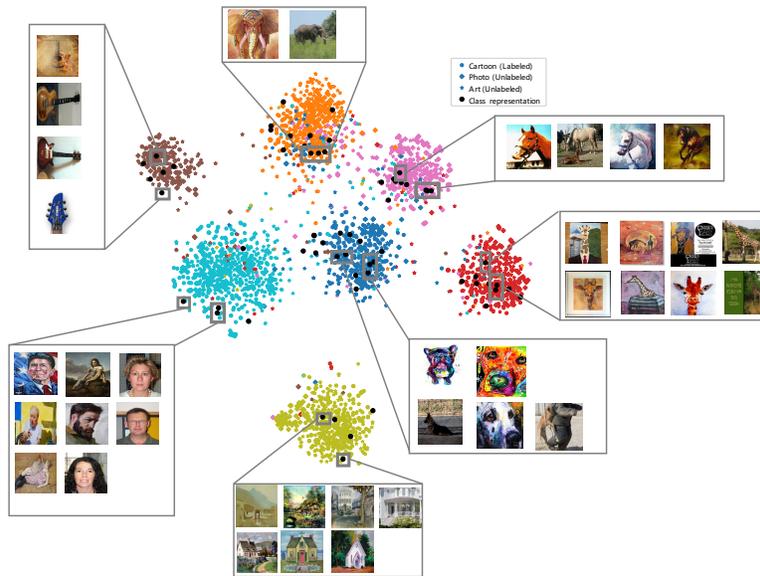
This work was supported by NSFC (62222604, 62192783), CAAI-Huawei MindSpore Project (CAAIXSJLJJ-2021-042A), China Postdoctoral Science Foundation Project (2021M690609), Jiangsu Natural Science Foundation Project (BK20210224), and CCF-Lenovo Bule Ocean Research Fund.

## References

- [1] X. Han, L. Qi, Q. Yu, Z. Zhou, Y. Zheng, Y. Shi, Y. Gao, Deep symmetric adaptation network for cross-modality medical image segmentation, *IEEE Transactions on Medical Imaging (TMI)* (2021) 1–1.
- [2] C.-X. Ren, X.-L. Xu, Z. Lei, A deep and structured metric learning method for robust person re-identification, *Pattern Recognition (PR)* (2019) 106995.



(a) Labeled:Art Unlabeled:Photo,Sketch



(b) Labeled:Cartoon Unlabeled:Photo,Art

Figure 7: The t-SNE [70] visualization of feature distribution in source domains. The selected class representation of all unlabeled domains are marked and the original images are shown. Two combinations of PACS are involved.

- [3] L. Dong-Hyun, Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks, in: International Conference on Machine Learning Workshop (ICMLW), 2013.
- [4] Y. Ganin, V. Lempitsky, Unsupervised domain adaptation by backpropagation, in: International Conference on Machine Learning (ICML), 2015.
- [5] M. Long, Y. Cao, J. Wang, M. I. Jordan, Learning transferable features with deep adaptation networks, in: International Conference on Machine Learning (ICML), 2015.
- [6] J. Jiang, X. Wang, M. Long, J. Wang, Resource efficient domain adaptation, in: ACM International Conference on Multimedia (ACMMM), 2020.
- [7] C.-Y. Lee, T. Batra, M. H. Baig, D. Ulbricht, Sliced wasserstein discrepancy for unsupervised domain adaptation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- [8] B. Gholami, P. Sahu, O. Rudovic, K. Bousmalis, V. Pavlovic, Unsupervised multi-target domain adaptation: An information theoretic approach, IEEE Transactions on Image Processing (TIP) (2020) 3993–4002.
- [9] J. Li, E. Chen, Z. Ding, L. Zhu, K. Lu, H. T. Shen, Maximum density divergence for domain adaptation, IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) (2021) 3918–3930.
- [10] Y.-H. Liu, C.-X. Ren, X.-L. Xu, K.-K. Huang, Bures joint distribution alignment with dynamic margin for unsupervised domain adaptation (2022). [arXiv:arXiv:2203.06836](https://arxiv.org/abs/2203.06836).
- [11] Y. Ganin, V. Lempitsky, Unsupervised domain adaptation by back propagation, in: International Conference on Machine Learning (ICML), 2015.
- [12] E. Tzeng, J. Hoffman, K. Saenko, T. Darrell, Adversarial discriminative domain adaptation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.

- [13] W. Zhang, W. Ouyang, W. Li, D. Xu, Collaborative and adversarial network for unsupervised domain adaptation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.
- [14] C. Zhang, Q. Zhao, Y. Wang, Hybrid adversarial network for unsupervised domain adaptation, Information Sciences (2020) 44–55.
- [15] J. Li, Z. Du, L. Zhu, Z. Ding, K. Lu, H. T. Shen, Divergence-agnostic unsupervised domain adaptation by adversarial attacks, IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) (2021) 1–1.
- [16] C. Chen, W. Xie, T. Xu, W. Huang, Y. Rong, X. Ding, Y. Huang, J. Huang, Progressive feature alignment for unsupervised domain adaptation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- [17] Y. Zhang, B. D. Davison, Adversarial continuous learning in unsupervised domain adaptation, in: International Conference on Pattern Recognition (ICPR), 2021.
- [18] J. Li, M. Jing, H. Su, K. Lu, L. Zhu, H. T. Shen, Faster domain adaptation networks, IEEE Transactions on Knowledge and Data Engineering (TKDE) (2021) 1–1.
- [19] J. Wang, C. Lan, C. Liu, Y. Ouyang, W. Zeng, T. Qin, Generalizing to unseen domains: A survey on domain generalization (2021). [arXiv:2103.03097](https://arxiv.org/abs/2103.03097).
- [20] M. M. Rahman, C. Fookes, M. Baktashmotlagh, S. Sridharan, Multi-component image translation for deep domain generalization, in: IEEE Winter Conference on Applications of Computer Vision(WACV), 2019.
- [21] K. Zhou, Y. Yang, T. Hospedales, T. Xiang, Learning to generate novel domains for domain generalization, in: European Conference on Computer Vision (ECCV), 2020.

- [22] K. Zhou, Y. Yang, T. M. Hospedales, T. Xiang, Deep domain-adversarial image generation for domain generalisation., in: AAAI Conference on Artificial Intelligence (AAAI), 2020.
- [23] R. Volpi, H. Namkoong, O. Sener, J. C. Duchi, V. Murino, S. Savarese, Generalizing to unseen domains via adversarial data augmentation, in: Advances in Neural Information Processing Systems (NeurIPS), 2018.
- [24] Y. Han, P. Zhang, W. Huang, Y. Zha, G. D. Cooper, Y. Zhang, Robust visual tracking based on adversarial unlabeled instance generation with label smoothing loss regularization, Pattern Recognition (PR) 97 (2020) 107027.
- [25] Z. Li, W. Zhao, F. Shi, L. Qi, X. Xie, Y. Wei, Z. Ding, Y. Gao, S. Wu, J. Liu, Y. Shi, D. Shen, A novel multiple instance learning framework for covid-19 severity assessment via data augmentation and self-supervised learning, Medical Image Analysis (MedIA) (2021) 101978.
- [26] Y. Li, M. Gong, X. Tian, T. Liu, D. Tao, Deep domain generalization via conditional invariant adversarial networks, in: European Conference on Computer Vision (ECCV), 2018.
- [27] T. Matsuura, T. Harada, Domain generalization using a mixture of multiple latent domains., in: AAAI Conference on Artificial Intelligence (AAAI), 2020.
- [28] H. Li, S. Jialin Pan, S. Wang, A. C. Kot, Domain generalization with adversarial feature learning, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.
- [29] S. Seo, Y. Suh, D. Kim, G. Kim, J. Han, B. Han, Learning to optimize domain specific normalization for domain generalization (2020). [arXiv: 1907.04275](https://arxiv.org/abs/1907.04275).

- [30] Y. Balaji, S. Sankaranarayanan, R. Chellappa, Metareg: Towards domain generalization using meta-regularization, in: *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [31] Q. Dou, D. C. de Castro, K. Kamnitsas, B. Glocker, Domain generalization via model-agnostic learning of semantic features, in: *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [32] Generalizable model-agnostic semantic segmentation via target-specific normalization, *Pattern Recognition (PR)* 122 (2022) 108292.
- [33] M. Mancini, S. R. Bulò, B. Caputo, E. Ricci, Best sources forward: domain generalization through source-specific nets, in: *International Conference on Image Processing (ICIP)*, 2018.
- [34] F. M. Carlucci, A. D’Innocente, S. Bucci, B. Caputo, T. Tommasi, Domain generalization by solving jigsaw puzzles, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [35] X. Zhai, A. Oliver, A. Kolesnikov, L. Beyer, S4l: Self-supervised semi-supervised learning, in: *IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [36] Y. Grandvalet, Y. Bengio, Semi-supervised learning by entropy minimization, in: *Advances in Neural Information Processing Systems (NeurIPS)*, 2005.
- [37] Z. Donyavi, S. Asadi, Diverse training dataset generation based on a multi-objective optimization for semi-supervised classification, *Pattern Recognition (PR)* (2020) 107543.
- [38] Y. Shi, J. Zhang, T. Ling, J. Lu, Y. Zheng, Q. Yu, L. Qi, Y. Gao, Inconsistency-aware uncertainty estimation for semi-supervised medical image segmentation, *IEEE Transactions on Medical Imaging (TMI)* (2021) 1–1.

- [39] S. Laine, T. Aila, Temporal ensembling for semi-supervised learning, in: International Conference on Learning Representations (ICLR), 2017.
- [40] A. Tarvainen, H. Valpola, Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results, in: Advances in Neural Information Processing Systems (NeurIPS), 2017.
- [41] Q. Xie, Z. Dai, E. Hovy, M.-T. Luong, Q. V. Le, Unsupervised data augmentation for consistency training (2020). [arXiv:1904.12848](#).
- [42] E. D. Cubuk, B. Zoph, J. Shlens, Q. V. Le, Randaugment: Practical automated data augmentation with a reduced search space (2019). [arXiv:1909.13719](#).
- [43] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, C. Raffel, Mixmatch: A holistic approach to semi-supervised learning (2019). [arXiv:1905.02249](#).
- [44] D. Berthelot, N. Carlini, E. D. Cubuk, A. Kurakin, K. Sohn, H. Zhang, C. Raffel, Remixmatch: Semi-supervised learning with distribution alignment and augmentation anchoring (2020). [arXiv:1911.09785](#).
- [45] C.-W. Kuo, C.-Y. Ma, J.-B. Huang, Z. Kira, Featmatch: Feature-based augmentation for semi-supervised learning, in: European Conference on Computer Vision (ECCV), 2020.
- [46] K. Sohn, D. Berthelot, C.-L. Li, Z. Zhang, N. Carlini, E. D. Cubuk, A. Kurakin, H. Zhang, C. Raffel, Fixmatch: Simplifying semi-supervised learning with consistency and confidence, in: Advances in Neural Information Processing Systems (NeurIPS), 2020.
- [47] H. Sharifi-Noghabi, H. Asghari, N. Mehrasa, M. Ester, Domain generalization via semi-supervised meta learning (2020). [arXiv:2009.12658](#).
- [48] K. Zhou, C. C. Loy, Z. Liu, Semi-supervised domain generalization with stochastic stylematch (2021). [arXiv:2106.00592](#).

- [49] Y. Liao, R. Huang, J. Li, Z. Chen, W. Li, Deep semisupervised domain generalization network for rotary machinery fault diagnosis under variable speed, *IEEE Transactions on Instrumentation and Measurement (TIM)* (2020) 8064–8075.
- [50] N. Anh, Y. Jason, C. Jeff, Deep neural networks are easily fooled: High confidence predictions for unrecognizable images, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [51] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, Dropout: A simple way to prevent neural networks from overfitting, *Journal of Machine Learning Research (JMLR)* 15 (1) (2014) 1929–1958.
- [52] H. Zhang, M. Cisse, Y. Dauphin, D. Lopez-Paz, mixup: Beyond empirical risk minimization, in: *International Conference on Learning Representations (ICLR)*, 2018.
- [53] D. Li, Y. Yang, Y.-Z. Song, T. M. Hospedales, Deeper, broader and artier domain generalization, in: *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [54] H. Venkateswara, J. Eusebio, S. Chakraborty, S. Panchanathan, Deep hashing network for unsupervised domain adaptation, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [55] K. Zhou, Y. Yang, Y. Qiao, T. Xiang, Domain adaptive ensemble learning, *IEEE Transactions on Image Processing (TIP)* 30 (2021) 8008–8018.
- [56] C. Fang, Y. Xu, D. N. Rockmore, Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias, in: *IEEE International Conference on Computer Vision (ICCV)*, 2013.
- [57] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

- [58] D. Berthelot, R. Roelofs, K. Sohn, N. Carlini, A. Kurakin, Adamatch: A unified approach to semi-supervised learning and domain adaptation, in: International Conference on Learning Representations (ICLR), 2022.
- [59] Y. Zhang, H. Tang, K. Jia, M. Tan, Domain-symmetric networks for adversarial domain adaptation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- [60] H. Tang, K. Chen, K. Jia, Unsupervised domain adaptation via structurally regularized deep clustering, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020.
- [61] Z. Du, J. Li, H. Su, L. Zhu, K. Lu, Cross-domain gradient discrepancy minimization for unsupervised domain adaptation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2021.
- [62] J. Liang, D. Hu, J. Feng, Domain adaptation with auxiliary target domain-oriented classifier, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2021.
- [63] J. Na, H. Jung, H. J. Chang, W. Hwang, Fixbi: Bridging domain spaces for unsupervised domain adaptation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2021.
- [64] Z. Wang, Y. Luo, R. Qiu, Z. Huang, M. Baktashmotlagh, Learning to diversify for single domain generalization, in: IEEE International Conference on Computer Vision (ICCV), 2021.
- [65] F. M. Carlucci, A. D’Innocente, S. Bucci, B. Caputo, T. Tommasi, Domain generalization by solving jigsaw puzzles, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- [66] Z. Huang, H. Wang, E. P. Xing, D. Huang, Self-challenging improves cross-domain generalization, in: European Conference on Computer Vision (ECCV), 2020.

- [67] Y. Liao, R. Huang, J. Li, Z. Chen, W. Li, Deep semisupervised domain generalization network for rotary machinery fault diagnosis under variable speed, *IEEE Transactions on Instrumentation and Measurement (TIM)* (2020) 8064–8075.
- [68] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, A density-based algorithm for discovering clusters in large spatial databases with noise, in: *International Conference on Knowledge Discovery and Data Mining (KDD)*, 1996.
- [69] J. H., W. Jr., Hierarchical grouping to optimize an objective function, *Journal of the American Statistical Association (JASA)* (1963) 236–244.
- [70] L. van der Maaten, G. Hinton, Viualizing data using t-sne, *Journal of Machine Learning Research (JMLR)* (2008) 2579–2605.