

Invertible Projection and Conditional Alignment for Multi-Source Blended-Target Domain Adaptation

Yuwu Lu^{1,2}, Haoyu Huang¹, Waikung Wong^{2*}, Xue Hu¹

¹South China Normal University, Guangzhou, China

²Hong Kong Polytechnic University, Hong Kong, China

{luyuwu2008, hyhuang99, hx1430940232}@163.com, calvin.wong@polyu.edu.hk

Abstract

Multi-source domain adaptation (MSDA), which utilizes multiple source domains to align the distribution of a single target domain, is a popular and challenging setting in domain adaptation (DA). However, existing MSDA approaches are difficult to obtain sufficient target domain knowledge, which serves as the transfer object. Furthermore, the target distributions are confused in the real world, i.e., the model cannot obtain the domain labels of target domains. To tackle these problems, we consider a more realistic DA setting Multi-Source Blended-Target Domain Adaptation (MBDA) and propose an *Invertible Projection and Conditional Alignment* (IPCA) method. Specifically, to reduce the impact of the distribution discrepancy, we construct an invertible projection for the source and blended-target domains. Then, we adopt a projection consistency regularization to our model, which makes the model more robust on the domain-specific parts. In addition, because the labels of the blended-target domain are unseen, we introduce conditional discrepancy to obtain the domain-level discriminative information and guide the classifier to serve as the discriminator, which is suitable for MBDA setting. Extensive experiment results on the ImageCLEF-DA, Office-Home, and DomainNet datasets validate the effectiveness of our method.

Code — <https://github.com/hyhuang99/IPCA>

Introduction

Recently, deep neural networks (DNNs) have attracted extensive attention due to their promising performance, especially in computer vision tasks (He et al. 2023; Huang et al. 2023; Ni et al. 2023). However, the performance of DNNs is based on the setting that training data and test data are drawn from the same distribution. In a real-world environment, the training data and test data are usually drawn from different distributions, which may lead to model degradation due to domain shift (Pan and Yang 2010) between the training and test data. To address this issue, unsupervised domain adaptation (UDA) (Ganin and Lempitsky 2015; Long et al. 2018) has been proposed.

The standard UDA setting is to learn domain-invariant representations from a single source domain (training data)

*Corresponding author.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

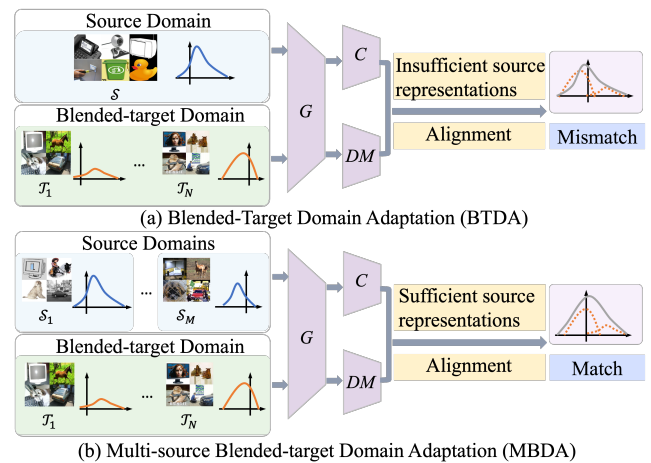


Figure 1: The comparison of (a) BTDA and (b) MBDA. F denotes the feature generator, C denotes the classifier, and DM denotes the discrepancy matching technique. Due to the lack of source representations in BTDA, the model is difficult to align distribution.

and a single target domain (test data). Due to insufficient feature information from a single source domain and a single target domain, the standard UDA methods (Wu et al. 2018; Ganin et al. 2016) are difficult to learn domain-invariant representations. Therefore, to obtain more feature information, multi-source domain adaptation (MSDA) attracts more attention (Peng et al. 2019; Zhao et al. 2018). Existing MSDA methods can be roughly divided into two categories. The first category is metric-based methods (Peng et al. 2019; Zhu, Zhuang, and Wang 2019), which measure and minimize the distribution discrepancy between multiple source and target domains. The second category is adversarial learning-based methods (Zhao et al. 2018; Li et al. 2021b), which group multiple source and target domains into multiple adversarial learning groups and perform the min-max game between the feature generator and the relevant discriminator to align feature implicitly.

Although existing MSDA methods have made significant progress, they are still facing some challenges. Current MSDA works are mainly based on the assumption that the

transfer tasks are from multiple source domains to a single target domain. However, in many real-world applications, there is a large amount of unlabeled target data drawn from different distributions. Due to the distribution discrepancy between different target domains, MSDA may not be the best solution to adapt the model to multiple target domains. Based on the analysis from Yang et al. (2022), the underlying data distributions of these unlabeled domains are not totally different; instead, they exhibit a certain degree of similarity. Thus, these domains can be used to construct a more effective and robust model for the task of cross-multiple domains, which is the objective of multi-target domain adaptation (MTDA). Combining the advantages of MSDA, multi-source and multi-target (MMDA) has been explored (Wang et al. 2021; Lu et al. 2024). Besides, as shown in Figure 1a, multiple target domains are mixed into a large domain in some cases, which is a scenario called blended-target domain adaptation (BTDA) (Chen et al. 2019; Xu, Wang, and Ling 2023) that transfers knowledge from a single source domain to a blended-target domain. But the sufficient feature representations from multiple source domains are largely overlooked in BTDA.

In this paper, as shown in Figure 1, we focus on a more realistic DA scenario that is multi-source blended-target domain adaptation (MBDA) and propose a novel method called Invertible Projection and Conditional Alignment (IPCA). Specifically, we construct an invertible projection for the source and blended-target domains, which aims at projecting the relevant domain features to more correlative latent feature spaces when facing the distribution discrepancy from different domains. Then, these feature spaces are utilized to optimize their original feature spaces. In addition, to further explore the domain-specific attributes and reduce the impact of domain-irrelevant information between multi-source and blended-target domains, we utilize a projection consistency regularization for the original and projected feature spaces. Finally, we utilize the conditional kernel bures (CKB) metric to measure the discrepancy between source and target domains and construct an adversarial learning strategy with the guidance of CKB. The adversarial learning strategy of our method is without the requirement of domain labels for the target domain, which is suitable for the MBDA scenario.

The contributions of this paper are threefold:

- A more realistic and challenge scenario, MBDA, is considered in the paper, which aims to transfer knowledge from multiple source domains to the blended-target domain without the domain labels of the sub-target domains.
- We propose a novel IPCA approach to address the MBDA problem, which utilizes the invertible projection consistency regularization to reduce the impact of domain discrepancy from multiple source and blended-target domains.
- We construct a CKB-guided adversarial learning strategy that matches the requirement of MBDA scenario. Extensive experiments demonstrate that our approach effectively tackles the new MBDA issue and achieves the su-

perior performance comparing with the state-of-the-art methods.

Related Works

Single-source and Single-target DA (SSDA)

SSDA is the standard UDA scenario, where a model is adapted from a single source domain to a single target domain. SSDA methods can be roughly categorized into metric learning-based methods (Long et al. 2019; Xie et al. 2023) and adversarial learning-based methods (Ganin et al. 2016; Chen et al. 2022). For metric learning-based methods, maximum mean discrepancy (MMD) is a common metric technique in SSDA, such as DAN (Long et al. 2019), which utilized multi-kernel MMD to reduce distribution differences between domains. CAF (Xie et al. 2023) utilized Wasserstein distance and sliced Wasserstein distance (SWD) (Lee et al. 2019) to minimize domain gaps. Adversarial learning-based methods optimize a feature generator and domain discriminator to align domains in a unified feature space. DANN (Ganin et al. 2016) introduced the gradient reverse layer for adversarial training, while DALN (Chen et al. 2022) replaced the domain discriminator with nuclear Wasserstein distance (NWD) to guide domain discrimination. Since the limited feature representations from SSDA, we consider the more challenging and realistic multi-source and multi-target domain setting in UDA, where data come from different distributions.

Multiple Domains DA

The existing DA methods in multiple domains setting can be roughly divided into four parts: multi-source domain adaptation (MSDA) (Zhao et al. 2018; Li et al. 2021b), multi-target domain adaptation (MTDA) (Yang et al. 2022), blended-target domain adaptation (BTDA) (Xu, Wang, and Ling 2023; Chen et al. 2019), and multi-source multi-target domain adaptation (MMDA) (Wang et al. 2021). MDAN (Zhao et al. 2018) constructed multiple domain discriminators to learn domain-invariant features across all domains and analyzed the generalization bounds of the MSDA classification problem. DRT (Li et al. 2021b) constructed a dynamic model to reduce the negative impact of the domain gap in MSDA. In MTDA, HGAN (Yang et al. 2022) utilized heterogeneous graph attention network to transfer multiple domains semantic information and align single source and multiple target domain distributions. BTDA is a more realistic scenario than MTDA. The BTDA models cannot know which distribution the target domains are drawn from, i.e., the domain labels of target domains are unseen. For example, MCDA (Xu, Wang, and Ling 2023) constructed a categorical domain discriminator and utilized low-level features to address the BTDA issue. In MMDA, the model aims to obtain more feature information from multiple source domains and more transfer objects from multiple target domains. AMDA (Wang et al. 2021) constructed an intra-domain and inter-domain attention module to explore the transferable knowledge between multiple domains and utilized multiple discriminators to align multi-source and multi-target domains. Unlike the above methods, our method

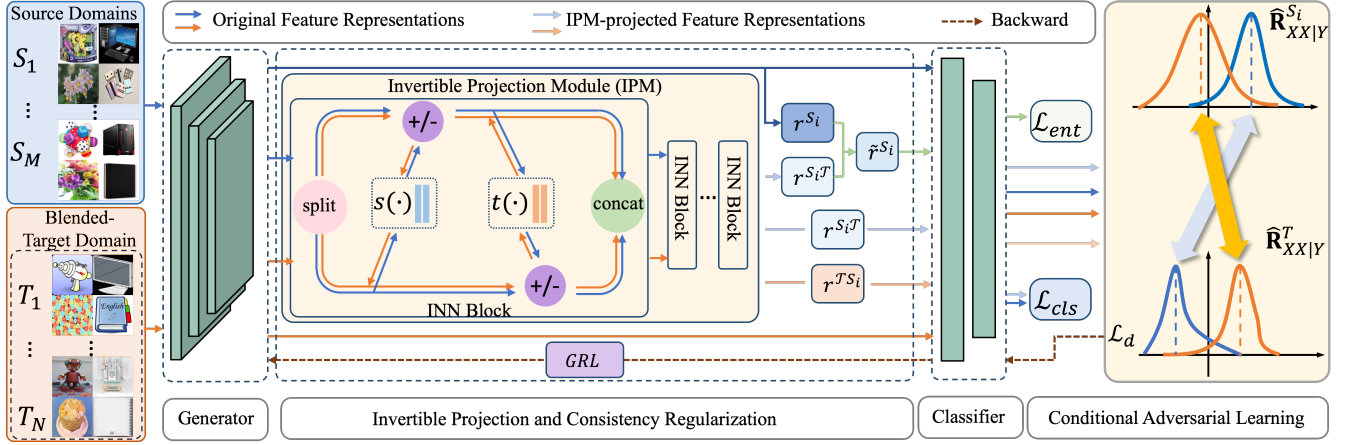


Figure 2: Overview of our proposed method. *GRL* is the gradient reverse layer. The invertible projection module is implemented by K INN blocks, $s(\cdot)$ and $t(\cdot)$ are constructed with two layers fully connected neural network. \mathcal{L}_{cls} is the cross-entropy loss on the source domains. \mathcal{L}_{ent} is the conditional entropy maximization loss on \tilde{r}^{S_i} to regularize the variation parts between r^{S_i} and $r^{S_i^T}$. \mathcal{L}_d is the CKB-guided adversarial learning loss that can be coupled with classifier to implicitly serve as a discriminator without domain labels.

does not require the domain labels of the sub-target domains and focuses on a more realistic MBDA scenario to adapt the model to mixture distribution.

Method

Preliminaries

In the MBDA setting, we have M labeled source domains $\{\mathcal{S}_i\}_{i=1}^M$ and a blended-target domain \mathcal{T} combined with N unlabeled sub-target domains $\{\mathcal{T}_i\}_{i=1}^N$. The samples of i -th source domain selected from distribution $P_{\mathcal{S}_i}(\mathcal{X}, \mathcal{Y})$ can be defined as $\mathcal{S}_i = \{(x_j^{S_i}, y_j^{S_i})\}_{j=1}^{n_{S_i}}$, which contains n_{S_i} samples. Besides, the samples of the blended-target domain selected from distribution $P_{\mathcal{T}}(\mathcal{X})$ can be defined as $\mathcal{T} = \{(x_j^T)\}_{j=1}^{n_{\mathcal{T}}}$ which contains $n_{\mathcal{T}}$ samples. Here, the distribution of blended-target domain is the combination of the sub-target domain distributions $\{P_{\mathcal{T}_j}(\mathcal{X})\}_{j=1}^N$, i.e., $P_{\mathcal{T}}(\mathcal{X}) = \sum_{j=1}^N \pi_j P_{\mathcal{T}_j}(\mathcal{X})$, where $\pi \in [0, 1]$ and $\sum_{j=1}^N \pi_j = 1$. The objective of MBDA is to learn an adaptive model on $\{\mathcal{S}_i\}_{i=1}^M$ and \mathcal{T} , that can generalize well on the unlabeled blended-target domain. Considering the confused target distribution space in MBDA, we construct an invertible projection for the multiple sources and blended-target domains and adopt a projection consistency regularization to optimize the classification model. The conditional distribution-guided discriminator-free adversarial learning strategy effectively guides the classifier to serve as the discriminator, which is suitable for MBDA since the adversarial learning strategy is without the requirement of domain labels. The overview of IPCA is illustrated in Fig. 2.

Invertible Projection and Consistency

Invertible Projection Module (IPM). In our training steps, the data obtained from multiple source and blended-target domains are first handled by a shared feature generator G

and mapped to the corresponding feature space. The analysis in Lugmayr *et al.* (2020) demonstrates that the Invertible Neural Network (INN) can connect two distributions due to its information preservation property. Considering the confused distribution of blended-target domain and the distribution discrepancy between multiple sources and blended-target domains, we utilize K INN blocks to construct the IPM so as to obtain better feature representations. Specifically, the source feature space is transformed to the blended-target feature space in the forward projection step of IPM, i.e., from $f_{\mathcal{S}}$ to $f_{\mathcal{T}}$. Thus, the input feature of k -th block of IPM can be defined as $z_{1:d}^k$, d denotes the dimension of the feature. The input feature $z_{1:d}^k$ is split into two parts $[z_{1:d/2}^k, z_{1+d/2:d}^k]$ and the linear neural networks $s(\cdot)$ and $t(\cdot)$ are utilized to respectively handle $z_{1:d/2}^k$ and $z_{1+d/2:d}^k$. The output of k -th block of IPM is given as follows:

$$\begin{aligned} z_{1:d/2}^{k+1} &= z_{1:d/2}^k + s(z_{1+d/2:d}^k), \\ z_{1+d/2:d}^{k+1} &= z_{1+d/2:d}^k + t(z_{1:d/2}^{k+1}). \end{aligned} \quad (1)$$

In the inverse projection step of IPM, the blended-target domain feature space is transformed to the source feature space, i.e., from $f_{\mathcal{T}}$ to $f_{\mathcal{S}}$. According to Eq. (1), the output of k -th block of IPM is defined as follows:

$$\begin{aligned} z_{1+d/2:d}^k &= z_{1+d/2:d}^{k+1} - t(z_{1:d/2}^{k+1}), \\ z_{1:d/2}^k &= z_{1:d/2}^{k+1} - s(z_{1+d/2:d}^k). \end{aligned} \quad (2)$$

After being projected by the IPM, the feature spaces of multiple sources and blended-target domains can be well connected. Hence, we obtain two better feature spaces $f_{\mathcal{S}\mathcal{T}} = G_{IPM}^{forward}(f_{\mathcal{S}})$ and $f_{\mathcal{T}\mathcal{S}} = G_{IPM}^{inverse}(f_{\mathcal{T}})$, where $G_{IPM}^{forward}$ and $G_{IPM}^{inverse}$ respectively denote the forward and inverse projections of Eqs. (1) and (2).

Projection Consistency. Although we can successfully

project the feature spaces of multiple source and blended-target domains to obtain better feature representations, it still does not guarantee that the model will maintain robustness in the domain-specific distribution. Inspired by self-supervised learning that makes use of consistency regularization for training, we also adopt the technique of conditional entropy maximization on the original and IPM-projected feature spaces. Specifically, given a sample x^{S_i} from i -th source domain, we can obtain two different representations of the same sample as $r^{S_i} = G(x^{S_i})$ and $r^{S_i\mathcal{T}} = G_{IPM}^{forward}(G(x^{S_i}))$.

During the training steps, we define \tilde{r}^{S_i} as the variation parts in the IPM-projected representation relative to the original representation. Thus, the relation between the IPM-projected representation $r^{S_i\mathcal{T}}$ and original representation r^{S_i} is:

$$r^{S_i\mathcal{T}} = \lambda r^{S_i} + (1 - \lambda)\tilde{r}^{S_i}, \quad (3)$$

$$\tilde{r}^{S_i} = \frac{r^{S_i\mathcal{T}} - \lambda r^{S_i}}{1 - \lambda}, \quad (4)$$

where λ denotes the proportion of r^{S_i} in $r^{S_i\mathcal{T}}$ and $\lambda \in (0, 1)$. To make sure that \tilde{r}^{S_i} does not contain too much category discriminative information, we maximize the conditional entropy of \tilde{r}^{S_i} , and the objective function is formulated as:

$$\mathcal{L}_{ent}^{S_i}(\tilde{r}^{S_i}) = -\frac{1}{n_{S_i}} \sum_{j=1}^{n_{S_i}} \sum_{c=1}^C H^{(c)}(\tilde{r}^{S_i}) \log H^{(c)}(\tilde{r}^{S_i}), \quad (5)$$

where H denotes the category classifier and $H^{(c)}(\tilde{r}^{S_i})$ indicates the probability of predicting variation part \tilde{r}^{S_i} to class c . By maximizing Eq. (5), the predictions of \tilde{r}^{S_i} will be regularized to have equal probability of being classified into each category.

In addition, considering the proportion of \tilde{r}^{S_i} in the IPM-projected representation, $r^{S_i\mathcal{T}}$ is not fixed in the training process. In the early stages of model training, the model is sensitive to domain-specific parts. Thus, the hyperparameter λ that controls the proportion of \tilde{r}^{S_i} should be a tiny value due to the discrepancy between r^{S_i} and $r^{S_i\mathcal{T}}$ is large. As the model training goes on, the model becomes less sensitive to domain-specific parts because the model lays more emphasis on domain-invariant parts. Hence, the value of λ should be increased. Based on the above-mentioned analysis, we utilize an annealing strategy proposed in (Ganin et al. 2016) to control the value of λ as:

$$\lambda = \lambda' [1 - (1 + \tau \frac{t}{T})^{-\omega}], \quad (6)$$

where $\tau = 10$, $\omega = 0.75$, and λ' is the initial value of λ . T is the total iterations of training and t is the current value of iteration.

Different from the current feature space transform methods (Zhou et al. 2023; He et al. 2023), our method further considers the influence of domain-specific parts between original and projected feature spaces in the model training. The augmented feature spaces generated by IPM can further decrease the negative impact of domain-irrelevant information.

CKB-guided Conditional Adversarial Learning

It is widespread to utilize adversarial learning technique to align distribution in DA (Ganin et al. 2016; Li et al. 2021a; Long et al. 2018). A straightforward way to implement adversarial learning in DA is to introduce an additional domain discriminator to perform the min-max game with feature generator. However, the current adversarial learning methods (Long et al. 2018) may not be suitable in the MBDA scenario due to the hybrid target distribution of the blended-target domain. Inspired by the DALN (Chen et al. 2022), the category classifier H can be reused as the domain discriminator when a function can utilize the outputs of the classifier to represent the cross-domain discrepancy. Thus, we utilize CKB to measure the conditional distribution discrepancy between multiple source and blended-target domains and guide adversarial learning in MBDA. Following (Luo and Ren 2021), the CKB loss can be defined as:

$$\mathcal{L}_{CKB}(x_{S_i}, x_{\mathcal{T}}) = \hat{d}_{CKB}^2(\hat{\mathbf{R}}_{XX|Y}^{S_i}, \hat{\mathbf{R}}_{XX|Y}^{\mathcal{T}}), \quad (7)$$

where \hat{d}_{CKB}^2 denotes the empirical estimation of the CKB metric and $\hat{\mathbf{R}}_{XX|Y}^{S_i/\mathcal{T}}$ denotes the conditional covariance operator of i -th source domain and the blended-target domain, respectively. Specifically, the conditional covariance operator $\hat{\mathbf{R}}_{XX|Y}^{S_i/\mathcal{T}}$ can be defined as follows:

$$\hat{\mathbf{R}}_{XX|Y}^{S_i/\mathcal{T}} = \frac{1}{n_{S_i/\mathcal{T}}} \Phi_{S_i/\mathcal{T}} \mathbf{H}_{n_{S_i/\mathcal{T}}} \mathbf{C}_{S_i/\mathcal{T}} (\Phi_{S_i/\mathcal{T}} \mathbf{H}_{n_{S_i/\mathcal{T}}} \mathbf{C}_{S_i/\mathcal{T}})^{\top}, \quad (8)$$

where $\Phi_{S_i/\mathcal{T}} = [G(x_1^{S_i/\mathcal{T}}), \dots, G(x_{n_{S_i/\mathcal{T}}}^{S_i/\mathcal{T}})]$ denotes the feature map matrix of i -th source domain and blended-target domain, respectively. $\mathbf{H}_{n_{S_i/\mathcal{T}}} = \mathbf{I}_{n_{S_i/\mathcal{T}}} - \frac{1}{n_{S_i/\mathcal{T}}} (1, \dots, 1)_{n_{S_i/\mathcal{T}}} (1, \dots, 1)_{n_{S_i/\mathcal{T}}}^{\top}$ is the $n \times n$ centering matrix. $\mathbf{C}_{S_i/\mathcal{T}}$ denotes the conditional information contained in i -th source or blended-target domain. Thus, the empirical estimation of the CKB metric can be reformulated as:

$$\begin{aligned} & \hat{d}_{CKB}^2(\hat{\mathbf{R}}_{XX|Y}^{S_i}, \hat{\mathbf{R}}_{XX|Y}^{\mathcal{T}}) \\ &= \text{etr}[\mathbf{G}_X^{S_i}(\epsilon n_{S_i} \mathbf{I}_{n_{S_i}} + \mathbf{G}_Y^{S_i})^{-1}] + \text{etr}[\mathbf{G}_X^{\mathcal{T}}(\epsilon n_{\mathcal{T}} \mathbf{I}_{n_{\mathcal{T}}} + \mathbf{G}_Y^{\mathcal{T}})^{-1}] \\ & \quad - \frac{2}{\sqrt{n_{S_i} n_{\mathcal{T}}}} \|(\mathbf{H}_{n_{\mathcal{T}}} \mathbf{C}_{\mathcal{T}})^{\top} \mathbf{K}_{XX}^{\mathcal{T}S_i} (\mathbf{H}_{n_{\mathcal{T}}} \mathbf{C}_{\mathcal{T}})\|_* \end{aligned} \quad (9)$$

where $\epsilon > 0$ is the regularization parameter and $\|\cdot\|_*$ denotes the Nuclear norm. $\mathbf{K}_{XX}^{\mathcal{T}S_i}$ denotes the explicit kernel matrix and $(\mathbf{K}_{XX}^{\mathcal{T}S_i})_{jk} = k_{\mathcal{X}}(x_j^{\mathcal{T}}, x_k^{S_i})$. $\mathbf{G}_{X/Y}^{S_i} = \mathbf{H}_{n_{S_i}} \mathbf{K}_{XX/Y}^{S_i} \mathbf{H}_{n_{S_i}}$ and $\mathbf{G}_{X/Y}^{\mathcal{T}} = \mathbf{H}_{n_{\mathcal{T}}} \mathbf{K}_{XX/Y}^{\mathcal{T}} \mathbf{H}_{n_{\mathcal{T}}}$ are the centralized kernels of i -th source or blended-target domain samples, where $(\mathbf{K}_{XX}^{S_i/\mathcal{T}})_{jk} = k_{\mathcal{X}}(x_j^{S_i/\mathcal{T}}, x_k^{S_i/\mathcal{T}})$ and $(\mathbf{K}_{YY}^{S_i/\mathcal{T}})_{jk} = k_{\mathcal{Y}}(x_j^{S_i/\mathcal{T}}, x_k^{S_i/\mathcal{T}})$. Considering that the ground-truth target domain labels are unseen in MBDA, we use the pseudo-labels $\tilde{y}^{\mathcal{T}}$ obtained by classifier with softmax layer $\phi(\cdot)$ and $\tilde{y}^{\mathcal{T}} = \phi(H(G(x^{\mathcal{T}})))$.

In the adversarial learning strategy of our method, the category classifier not only serves category-oriented alignment

but also domain-oriented alignment. Note that when classifier serves the domain-oriented alignment between multiple sources and blended-target domains, we construct a min-max game between feature generator G and classifier H using a gradient reverse layer (GRL) (Ganin and Lempitsky 2015). Then, the CKB-guided adversarial learning strategy is defined as follows:

$$\min_G \max_H \mathcal{L}_{CKB}(r^{S_i}, r^{\mathcal{T}}). \quad (10)$$

Considering that we have the extra feature spaces $f_{S\mathcal{T}}$ and $f_{\mathcal{T}S}$ that are generated by $G_{IPM}^{forward}$ and $G_{IPM}^{inverse}$, we combine the relevant feature space when aligning the domain discrepancy. The adversarial loss in Eq. (10) can be further reformulated as:

$$\min_{G, G_{IPM}} \max_H \mathcal{L}_{com}^{S_i} = \mathcal{L}_{CKB}(r^{S_i}, r^{\mathcal{T}}) + \mathcal{L}_{CKB}(r^{S_i\mathcal{T}}, r^{\mathcal{T}}) + \mathcal{L}_{CKB}(r^{\mathcal{T}S_i}, r^{\mathcal{T}}). \quad (11)$$

Since the IPM can maintain the topological structure of features, based on the feature space combination by Eq. (11), the model is effective to achieve the cross-domain lossless feature transformation.

Overall Loss Function

The overall loss function which optimizes the model to transfer knowledge from multi-source to blended-target domains is defined as:

$$\mathcal{L}_{IPCA} = \mathcal{L}_{cls} - \alpha \mathcal{L}_d - \beta \mathcal{L}_{ent}, \quad (12)$$

where α and β are hyper-parameters.

The first term of Eq. (12) is the supervised learning on the labeled multi-source domains to make sure that the classification model can classify samples to the correct category. Thus, we utilize the cross-entropy loss function \mathcal{L}_{ce} to guide the supervised training, and define it as follows:

$$\min_G \mathcal{L}_{cls} = \sum_{i=1}^M \frac{1}{n_{S_i}} \sum_{j=1}^{n_{S_i}} (\mathcal{L}_{ce}(H(r_j^{S_i}), y_j^{S_i}) + \mathcal{L}_{ce}(H(r_j^{S_i\mathcal{T}}), y_j^{S_i})), \quad (13)$$

where $r_j^{S_i} = G(x_j^{S_i})$ and $r_j^{S_i\mathcal{T}} = G_{IPM}^{forward}(G(x_j^{S_i}))$.

The second term \mathcal{L}_d in Eq. (12) is the adversarial loss to perform min-max game between feature generator and classifier using the original and IPM-projected feature representations, the expression of which is denoted as:

$$\min_{G, G_{IPM}} \max_H \mathcal{L}_d = \sum_{i=1}^M \mathcal{L}_{com}^{S_i}. \quad (14)$$

The third term \mathcal{L}_{ent} in Eq. (12) is the loss function that enhances the model robustness of domain-specific parts \tilde{r}^{S_i} , we maximize the conditional entropy of the variation parts between original r^{S_i} and IPM-projected representations $r^{S_i\mathcal{T}}$ based on Eq. (5):

$$\max_H \mathcal{L}_{ent} = \sum_{i=1}^M \frac{1}{n_{S_i}} \sum_{j=1}^{n_{S_i}} \mathcal{L}_{ent}^{S_i}(\tilde{r}_j^{S_i}). \quad (15)$$

After training by the above-mentioned loss functions, our model can effectively transfer knowledge from multi-source domains to the blended-target domain.

Experiments

Experimental Setup

Datasets. We evaluated and compared the state-of-the-art (SOTA) methods with our method on three popular DA datasets (i.e., **DomainNet**, **Office-Home**, and **ImageCLEF-DA**). The **DomainNet** (Peng et al. 2019) is the largest dataset in DA that contains 0.6 million images from 345 categories in 6 domains: Clipart (c), Infograph (i), Painting (p), Quickdraw (q), Real-world (r), and Sketch (s). In our experiments, we sampled 126 categories and 4 domains (c, p, r, and s) to evaluate our model, following the protocol in (Zhou et al. 2021a). The **Office-Home** (Venkateswara et al. 2017) is a challenge dataset with label imbalance, which contains 15,500 images in total. It consists of 65 common categories in 4 domains: Art (Ar), Clipart (CI), Products (Pr), and Real-world (Rw). The **ImageCLEF-DA** (Caputo et al. 2014) contains a total of 2,400 images, including 12 common categories in 4 domains: Bing (B), Caltech (C), ImageNet (I), and Pascal (P).

Protocols. To highlight the challenge in MBDA setting, we cannot anymore use the standard protocols from the above three datasets. Thus, for SSDA setting, one column denotes one SSDA task, e.g., $r \rightarrow c$ in Table 1a. For MSDA setting, two domains are selected as sources and one domain is selected as target, e.g., $r+s \rightarrow c$ and $r+s \rightarrow p$ in Table 1a. For MTDA/BTDA setting, one domain is selected as source and two domains are selected as targets, e.g., $r \rightarrow c+p$ and $s \rightarrow c+p$ in Table 1a. For MMDA/MBDA setting, two domains are sources, and the other domains are targets, e.g., $r+s \rightarrow c+p$ in Table 1a. For a clearer comparison, all classification results in Tables 1 and 2 are the average of accuracies of two target domains. Results obtained from the released code of the corresponding method are marked with “*”, while the remaining results are obtained from their original papers. The best results are bolded.

Competitors. We compared our method in four different settings of DA methods. The first setting is SSDA methods, i.e., DANN (Ganin et al. 2016), DALN (Chen et al. 2022), and SCDA (Li et al. 2021a). The second setting of methods is MSDA, including MDAN (Zhao et al. 2018) and DIDA (Deng et al. 2022). The third setting of methods is MTDA/BTDA, including MTDA (Gholami et al. 2020) and MCDA (Xu, Wang, and Ling 2023). The last setting of methods is MMDA, including AMDA (Wang et al. 2021), DGWA (Lu et al. 2024), and HTA (Wu et al. 2023).

Implementation Details. We implemented and evaluated our method on the PyTorch (Paszke et al. 2019) platform; the number of PyTorch is 1.13.1. The number of INN blocks which contains in the IPM is $K = 5$. For fair comparison, all experiments on three datasets utilize the same backbone network, ResNet-50 (He et al. 2016), and run on a Nvidia GeForce RTX-4090 GPU. The version of CUDA is 11.7. The batch size of all experiments in the training step is set to 32. The optimizer is Stochastic Gradient Descent (SGD) with a momentum parameter of 0.9 and a weight decay of $1e-3$. In addition, the learning rate is set to $1e-3$ and updated by the LambdaLR (Paszke et al. 2019) during the training process.

Source Target	r+s	s+p	p+r	c+s	r+c	c+p	Avg.	Source Target	Rw+Pr	Cl+Rw	Pr+Cl	Rw+Ar	Ar+Pr	Cl+Ar	Avg.
	c+p	c+r	c+s	p+r	p+s	r+s			Ar+Cl	Ar+Pr	Ar+Rw	Cl+Pr	Cl+Rw	Pr+Rw	
DANN	31.4	39.7	26.8	29.3	31.3	31.2	31.6	DANN	53.5	61.9	53.5	55.6	57.1	60.1	57.6
MDAN	54.5	59.0	45.0	58.8	51.7	61.0	54.5	MDAN	51.9	64.9	60.3	59.4	58.2	62.4	59.5
MTDA	52.4	48.7	45.5	53.3	51.5	52.0	50.5	MTDA	55.4	69.1	61.2	61.5	55.9	70.4	62.2
AMDA	65.8	67.8	56.7	65.1	58.9	66.4	63.4	SCDA	64.1	74.7	70.0	68.3	68.7	77.6	70.1
DALN*	61.2	69.2	64.1	63.5	59.3	64.8	63.7	AMDA	61.4	77.0	72.3	67.4	64.9	77.4	70.0
MCDA*	62.2	68.7	61.7	63.4	61.2	65.4	63.8	MCDA*	63.6	74.9	70.0	68.7	68.1	78.1	70.6
DGWA*	66.4	71.3	63.4	67.5	64.6	70.2	65.0	DGWA	63.6	78.2	73.7	70.3	66.7	78.5	71.3
IPCA	70.0	75.7	66.3	70.8	65.0	72.7	70.1	IPCA	65.6	79.9	75.2	70.1	71.8	79.3	73.7

(a) DomainNet

(b) Office-Home

Table 1: Accuracy (%) on the (a) DomainNet and the (b) Office-Home datasets for MBDA (with ResNet-50 as a backbone).

Source Target	I+P	P+C	C+I	B+P	I+B	B+C	Avg.
	B+C	B+I	B+P	C+I	C+P	I+P	
DANN	76.4	72.4	69.1	87.9	82.9	79.3	77.9
AMDA	78.8	77.3	71.7	92.3	85.2	83.8	81.5
SCDA*	78.9	77.2	71.8	93.9	85.5	85.0	82.1
DIDA*	78.9	77.9	72.0	92.2	86.8	85.1	82.2
HTA	79.3	78.2	72.3	92.8	85.6	84.9	82.2
MCDA*	77.4	79.1	70.3	91.8	86.2	83.8	81.4
DGWA	79.7	79.1	72.7	93.8	86.0	84.5	81.7
IPCA	79.9	80.2	73.1	94.8	87.6	85.8	83.6

Table 2: Accuracy (%) on the ImageCLEF-DA (with ResNet-50 as a backbone).

Comparisons to the State-of-the-Art

DomainNet: Table 1a presents the results on the DomainNet dataset. It can be observed that IPCA can achieve significant progress compared with the state-of-the-art (SOTA) methods in most of the experimental tasks and achieve the highest average classification accuracy (**70.1%**). Note that as a BTDA method, although MCDA can obtain sufficient target information to support the adaptation, its performance is still significantly lower than that of our method due to the limited feature representations in the single source domain. Compared to the MMDA method, AMDA, although it can obtain sufficient representations from multiple source domains and access the domain labels of target domains, is still a suboptimal solution because it tries to align all feature spaces to a common space. Unlike MCDA and AMDA, IPCA considers the better source and targets feature space to utilize invertible projection and align the conditional distribution by the CKB metric without the requirement of domain labels of the target domain, which is suitable for MBDA.

Office-Home: The results on the Office-Home dataset are reported in Table 1b, where our proposed method generally outperforms the current SOTA approaches and achieves the best classification accuracy (**73.7%**). Note that there is a huge class imbalance on the Office-Home dataset; e.g., the sample number of domain Rw (34,856) is far more numerous than other domains. Our method can optimize the feature representations by IPM, which can transfer more useful

Source Target	r+s	s+p	p+r	c+s	r+c	c+p	Avg.
	c+p	c+r	c+s	p+r	p+s	r+s	
w/o IPM	68.3	74.2	64.8	68.3	62.5	70.4	68.1
w/o \mathcal{L}_{ent}	69.9	75.1	65.6	69.5	63.8	71.4	69.2
w/o \mathcal{L}_d	65.8	72.2	61.7	66.8	60.2	68.3	65.8
IPCA	70.0	75.7	66.3	70.8	65.0	72.7	70.1

(a) Effectiveness of each components

Source Target	I+P	P+C	C+I	B+P	I+B	B+C	Avg.
	B+C	B+I	B+P	C+I	C+P	I+P	
RA	78.9	78.5	71.9	93.3	86.7	85.5	82.5
pAdaIN	79.6	79.1	72.5	93.6	87.0	85.8	82.9
MixStyle	79.2	78.8	72.7	93.0	87.2	85.5	82.7
IPCA	79.9	80.2	73.1	94.8	87.6	85.8	83.6

(b) Comparisons of different reconstruction methods

Table 3: Ablation results produced by IPCA (with ResNet-50 as a backbone).

knowledge. These obtained improvements are mainly produced by the invertible projection learning and conditional adversarial alignment.

ImageCLEF-DA: Table 2 illustrates the experimental results of our method and other SOTA methods on the ImageCLEF-DA dataset. We can observe that our proposed method exceeds the SOTA approaches in most of the classification tasks and achieves the highest average classification accuracy (**83.6%**). These results show that our method is beneficial to the MBDA setting and effectively adapts the model to the confused target distribution.

Further Analysis

Ablation Study. Table 3a presents the ablation study on the DomainNet dataset, evaluating the core components of our method: 1) without IPM projection (w/o IPM), 2) without conditional entropy maximization loss \mathcal{L}_{ent} (w/o \mathcal{L}_{ent}), and 3) without conditional adversarial alignment loss \mathcal{L}_d (w/o \mathcal{L}_d). The full method achieves the best performance. Removing \mathcal{L}_d causes a significant 4.3% performance drop,

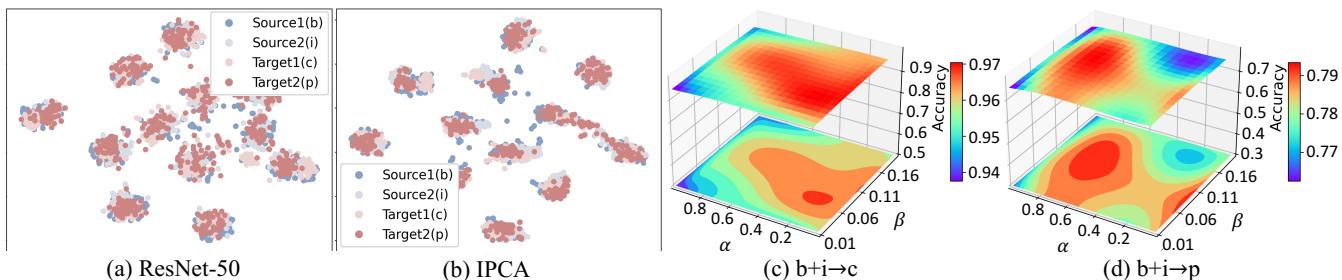


Figure 3: (a, b) t-SNE visualizations of the b+i→c+p task, (c, d) confusion matrices, and (e, f) analysis of the IPCA parameters (zoom in for better visualizations).

Targets	r+s	q+s	q+r	p+s	p+r	p+q	i+s	i+r	i+q	i+p	c+s	c+r	c+q	c+p	c+i	Avg.
DGWA*	58.3	30.5	38.5	50.9	68.6	33.0	36.5	44.5	28.3	38.8	53.6	62.3	36.7	59.0	42.3	44.1
IPCA(4S2T)	59.5	31.7	39.7	52.1	69.8	34.2	37.7	45.7	29.5	40.0	54.8	63.5	37.9	60.2	43.5	46.7
Sources	r+s	q+s	q+r	p+s	p+r	p+q	i+s	i+r	i+q	i+p	c+s	c+r	c+q	c+p	c+i	Avg.
DGWA*	39.9	50.4	48.3	35.2	39.0	50.2	38.8	42.4	48.5	39.3	42.3	37.4	46.0	37.5	41.6	40.5
IPCA(2S4T)	41.2	52.1	49.7	36.5	40.0	51.7	39.8	44.5	50.7	40.2	44.1	38.6	47.3	39.7	41.8	43.7
3S3T	c+i+p	c+i+q	c+i+r	c+p+q	c+p+r	c+p+s	c+p+s	c+r+s	i+p+q	i+p+r	i+p+s	i+q+r	i+r+s	p+q+r	p+q+s	Avg.
	q+r+s	p+r+s	p+q+s	i+r+s	i+q+s	i+q+r	i+p+r	i+p+q	c+r+s	c+q+s	c+q+r	c+p+s	c+p+q	c+i+s	c+i+r	
DGWA*	38.9	49.3	45.9	52.2	29.7	33.5	47.7	31.8	56.0	39.1	38.1	48.8	42.0	43.3	51.2	41.7
IPCA	40.1	50.8	47.3	53.7	30.6	34.5	49.1	32.7	57.5	40.3	39.2	50.3	43.2	44.7	52.7	44.5

Table 4: Accuracy (%) on the default version of DomainNet dataset (with ResNet-101 as a backbone).

highlighting its importance in aligning conditional distributions across sources and the blended target domains. Results for “w/o IPM” confirm its effectiveness in MBDA, while \mathcal{L}_{ent} enhances robustness on domain-specific features.

In Table 3b, we compare the performance of IPM with other feature reconstruction methods, including Random Augmentation (RA), pAdaIN (Nurriel, Benaim, and Wolf 2021), and MixStyle (Zhou et al. 2021b). The results show that IPM is better suited for the MBDA scenario than other reconstruction methods.

Feature Visualization. The t-SNE (Maaten and Hinton 2008) feature visualization of our method is shown in Figure 3b, alongside the ResNet-50 visualization in Figure 3a. Dots of different colors represent features from different domains. The source-only model (ResNet-50) clusters well on source features but struggles with target features. In contrast, our method generates highly discriminative features, clustering same-class features better and separating different-class features more effectively. These results further validate the effectiveness of the IPM mapping and conditional adversarial alignment in MBDA.

Hyper-parameter Sensitivity. Figures 3e and 3f report the sensitivity of our method to two loss function hyper-parameters α and β . We performed the experiments on the ImageCLEF-DA dataset with classification task b+i→c+p. The hyper-parameter choices are $\alpha = \{0.01, 0.05, 0.1, 0.5, 1.0\}$ and $\beta = \{0.01, 0.05, 0.1, 0.15, 0.2\}$. According to the results, our method is not sensitive to β , but is a little bit sensitive

to α . In general, the best choice of hyper-parameters is $\alpha = 1.0$ and $\beta = 0.1$.

Scalability. To evaluate IPCA’s scalability, we conduct experiments on the default version of DomainNet dataset, which includes 6 domains with 345 categories. As shown in Table 4, the experimental tasks are categorized into four-sources-two-targets (4S2T), two-sources-four-targets (2S4T), and three-sources-three-targets (3S3T). Compared to DGWA, IPCA achieves significant performance gains, demonstrating that IPCA is more suitable to the multiple domains scenario.

Conclusion

In this paper, we focus on the MBDA issue and propose an IPCA approach. Considering the confused distribution in the blended-target domain, we construct the invertible projection for multi-source and blended-target domains to obtain the better feature space and adopt a conditional entropy maximization strategy to enhance model robustness in domain-specific parts. A CKB-guided conditional adversarial learning strategy is applied to learn domain-invariant representations without the requirements of domain labels. Extensive experiments show that our method achieves the superior performance on three domain adaptation benchmarks.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (Grant No. 62176162) and the

References

- Caputo, B.; Müller, H.; Martinez-Gomez, J.; Villegas, M.; Acar, B.; Patricia, N.; Marvasti, N.; Üsküdarlı, S.; Paredes, R.; and Cazorla, M. 2014. Imageclef 2014: Overview and Analysis of the Results. In *ICCLEF*, 192–211.
- Chen, L.; Chen, H.; Wei, Z.; Jin, X.; Tan, X.; Jin, Y.; and Chen, E. 2022. Reusing the Task-specific Classifier as a Discriminator: Discriminator-free Adversarial Domain Adaptation. In *CVPR*, 7171–7181.
- Chen, Z.; Zhuang, J.; Liang, X.; and Lin, L. 2019. Blending-target Domain Adaptation by Adversarial Meta Adaptation Networks. In *CVPR*, 2248–2258.
- Deng, Z.; Zhou, K.; Li, D.; He, J.; Song, Y.-Z.; and Xiang, T. 2022. Dynamic Instance Domain Adaptation. *IEEE TIP*, 31: 4585–4597.
- Ganin, Y.; and Lempitsky, V. 2015. Unsupervised Domain Adaptation by Backpropagation. In *ICML*, 1180–1189.
- Ganin, Y.; Ustinova, E.; Ajakan, H.; Germain, P.; Larochelle, H.; Laviolette, F.; Marchand, M.; and Lempitsky, V. 2016. Domain-adversarial Training of Neural Networks. *JMLR*, 17(1): 2096–2130.
- Gholami, B.; Sahu, P.; Rudovic, O.; Bousmalis, K.; and Pavlovic, V. 2020. Unsupervised Multi-Target Domain Adaptation: An Information Theoretic Approach. *IEEE TIP*, 29: 3993–4002.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *CVPR*, 770–778.
- He, Q.; Xiao, S.; Ye, M.; Zhu, X.; Neri, F.; and Hou, D. 2023. Independent Feature Decomposition and Instance Alignment for Unsupervised Domain Adaptation. In *IJCAI*, 819–827.
- Huang, S.; Li, H.; Wang, Y.; Zhu, H.; Dai, J.; Han, J.; Rong, W.; and Liu, S. 2023. Discovering Sounding Objects by Audio Queries for Audio Visual Segmentation. In *IJCAI*, 875–883.
- Lee, C.-Y.; Batra, T.; Baig, M. H.; and Ulbricht, D. 2019. Sliced Wasserstein Discrepancy for Unsupervised Domain Adaptation. In *CVPR*, 10285–10295.
- Li, S.; Xie, M.; Lv, F.; Liu, C. H.; Liang, J.; Qin, C.; and Li, W. 2021a. Semantic Concentration for Domain Adaptation. In *ICCV*, 9102–9112.
- Li, Y.; Yuan, L.; Chen, Y.; Wang, P.; and Vasconcelos, N. 2021b. Dynamic Transfer for Multi-Source Domain Adaptation. In *CVPR*, 10998–11007.
- Long, M.; Cao, Y.; Cao, Z.; Wang, J.; and Jordan, M. I. 2019. Transferable Representation Learning with Deep Adaptation Networks. *IEEE TPAMI*, 41(12): 3071–3085.
- Long, M.; Cao, Z.; Wang, J.; and Jordan, M. I. 2018. Conditional Adversarial Domain Adaptation. In *NeurIPS*, 1640–1650.
- Lu, Y.; Huang, H.; Zeng, B.; Lai, Z.; and Li, X. 2024. Multi-Source and Multi-Target Domain Adaptation Based on Dynamic Generator with Attention. *IEEE TMM*, 26: 6891–6905.
- Lugmayr, A.; Danelljan, M.; Gool, L. V.; and Timofte, R. 2020. Srflo: Learning the Super-resolution Space with Normalizing Flow. In *ECCV*, 715–732.
- Luo, Y.; and Ren, C. 2021. Conditional Bures Metric for Domain Adaptation. In *CVPR*, 13984–13993.
- Maaten, L. V. D.; and Hinton, G. 2008. Visualizing Data Using T-SNE. *JMLR*, 9: 2579–2605.
- Ni, Z.-L.; Yang, F.; Wen, S.; and Zhang, G. 2023. Dual Relation Knowledge Distillation for Object Detection. In *IJCAI*, 1276–1284.
- Nuriel, O.; Benaim, S.; and Wolf, L. 2021. Permuted Adain: Reducing the Bias Towards Global Statistics in Image Classification. In *CVPR*, 9482–9491.
- Pan, S. J.; and Yang, Q. 2010. A Survey on Transfer Learning. *IEEE TKDE*, 22(10): 1345–1359.
- Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; and Antiga, L. 2019. Pytorch: An Imperative Style, High-performance Deep Learning Library. In *NeurIPS*, 8026–8037.
- Peng, X.; Bai, Q.; Xia, X.; Huang, Z.; Saenko, K.; and Wang, B. 2019. Moment Matching for Multi-Source Domain Adaptation. In *CVPR*, 1406–1415.
- Venkateswara, H.; Eusebio, J.; Chakraborty, S.; and Panchanathan, S. 2017. Deep Hashing Network for Unsupervised Domain Adaptation. In *CVPR*, 5018–5027.
- Wang, Y.; Zhang, Z.; Hao, W.; and Song, C. 2021. Attention Guided Multiple Source and Target Domain Adaptation. *IEEE TIP*, 30: 892–906.
- Wu, Z.; Han, X.; Lin, Y.-L.; Uzunbas, M. G.; Goldstein, T.; Lim, S.; and Davis, L. S. 2018. Dcan: Dual Channel-wise Alignment Networks for Unsupervised Scene Adaptation. In *ECCV*, 518–534.
- Wu, Z.; Meng, M.; Liang, T.; and Wu, J. 2023. Hierarchical Triple-Level Alignment for Multiple Source and Target Domain Adaptation. *Appl. Intell.*, 53: 3766–3782.
- Xie, B.; Li, S.; Lv, F.; Liu, C. H.; Wang, G.; and Wu, D. 2023. A Collaborative Alignment Framework of Transferable Knowledge Extraction for Unsupervised Domain Adaptation. *IEEE TKDE*, 35(7): 6518–6533.
- Xu, P.; Wang, B.; and Ling, C. 2023. Class Overwhelms: Mutual Conditional Blended-Target Domain Adaptation. In *AAAI*, 3036–3044.
- Yang, X.; Deng, C.; Liu, T.; and Tao, D. 2022. Heterogeneous Graph Attention Network for Unsupervised Multiple-Target Domain Adaptation. *IEEE TPAMI*, 44(4): 1992–2003.
- Zhao, H.; Zhang, S.; Wu, G.; Moura, J. M. F.; Costeira, J. P.; and Gordon, G. J. 2018. Adversarial Multiple Source Domain Adaptation. In *NeurIPS*, 8559–8570.
- Zhou, K.; Yang, Y.; Qiao, Y.; and Xiang, T. 2021a. Domain Adaptive Ensemble Learning. *IEEE TIP*, 30: 8008–8018.

Zhou, K.; Yang, Y.; Qiao, Y.; and Xiang, T. 2021b. Domain Generalization with Mixstyle. In *ICLR*.

Zhou, L.; Ye, M.; Zhu, X.; Xiao, S.; Fan, X.-Q.; and Neri, F. 2023. Homeomorphism Alignment for Unsupervised Domain Adaptation. In *ICCV*, 18669–18710.

Zhu, Y.; Zhuang, F.; and Wang, D. 2019. Aligning Domain-Specific Distribution and Classifier for Cross-Domain Classification from Multiple Sources. In *AAAI*, 5989–5996.