# Heterogeneous Domain Adaptation via Soft Transfer Network

Yuan Yao
School of Computer
Science and Technology,
Harbin Institute of
Technology, Shenzhen,
China
yaoyuan@stu.hit.edu.cn

Yu Zhang
Department of Computer
Science and Engineering,
Southern University of
Science and Technology,
Shenzhen, China
yu.zhang.ust@gmail.com

Xutao Li
School of Computer
Science and Technology,
Harbin Institute of
Technology, Shenzhen,
China
lixutao@hit.edu.cn

Yunming Ye*
School of Computer
Science and Technology,
Harbin Institute of
Technology, Shenzhen,
China
yeyunming@hit.edu.cn

## ABSTRACT

Heterogeneous domain adaptation (HDA) aims to facilitate the learning task in a target domain by borrowing knowledge from a heterogeneous source domain. In this paper, we propose a Soft Transfer Network (STN), which jointly learns a domain-shared classifier and a domain-invariant subspace in an end-to-end manner, for addressing the HDA problem. The proposed STN not only aligns the discriminative directions of domains but also matches both the marginal and conditional distributions across domains. To circumvent negative transfer, STN aligns the conditional distributions by using the soft-label strategy of unlabeled target data, which prevents the hard assignment of each unlabeled target data to only one category that may be incorrect. Further, STN introduces an adaptive coefficient to gradually increase the importance of the soft-labels since they will become more and more accurate as the number of iterations increases. We perform experiments on the transfer tasks of image-to-image, text-to-image, and text-to-text. Experimental results testify that the STN significantly outperforms several state-of-the-art approaches.

## CCS CONCEPTS

• **Computing methodologies → Transfer learning**; *Object recognition*; Semi-supervised learning settings.

## KEYWORDS

Heterogeneous domain adaptation, soft-label, adaptive coefficient, subspace learning

---
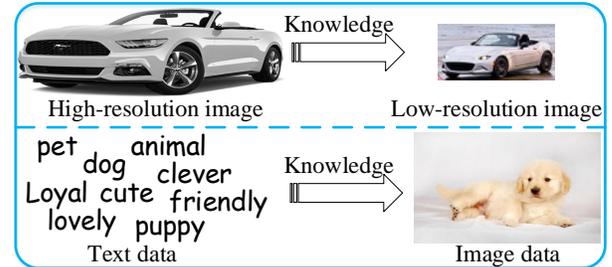
*Corresponding Author.

---

**Figure 1: We aim to transfer knowledge between different feature spaces. Top row: transferring the knowledge from high-resolution images to low-resolution ones. Bottom row: transferring the knowledge from text data to image data.**

## 1 INTRODUCTION

The success of supervised learning mostly depends on the availability of sufficient labeled data. However, in many real-world applications, it is often prohibitive and time-consuming to acquire enough labeled data. To alleviate this concern, *domain adaptation* (DA) models [8, 9, 30, 43] are proposed to facilitate the learning task in a label-scarce domain, *i.e.*, the target domain, by borrowing knowledge from a label-rich and related domain, *i.e.*, the source domain. They have achieved a series of successes in fields such as visual recognition [12, 25, 27, 28, 34, 35, 42, 53] and text categorization [44, 52]. As a rule, most existing studies [4, 5, 12, 24–28, 42] assume that the source and target domains share the same feature space. In reality, however, it is often not easy to seek for a source domain with the same feature space as the target domain of interest. Hence, we focus in this paper on a more general and challenging scenario where the source and target domains are drawn from different feature spaces, which is known as *heterogeneous domain adaptation* (HDA). For example, the source and target images are characterized by diverse resolutions (*e.g.*, the top row in Figure 1); similarly, the source domain is textual whereas the target one is visual (*e.g.*, the bottom row in Figure 1).

To bridge between two heterogeneous domains, existing HDA works typically choose to project data from one domain to another [17, 18, 21, 36, 37, 44, 52], or find a domain-invariant subspace [14, 19, 23, 32, 41, 45, 49]. On one hand, projections are learned in [14, 17, 18, 23] by performing the classifier adaptation strategy (*i.e.*, training a domain-shared classifier with labeled cross-domain data). However, the type of methods does not explicitly minimize the distributional discrepancy across domains. On the other hand, another line of studies [36, 41] only considers the distribution matching strategy (*i.e.*, reducing the distributional divergency between domains), but ignores the classifier adaptation one. Thus, this kind
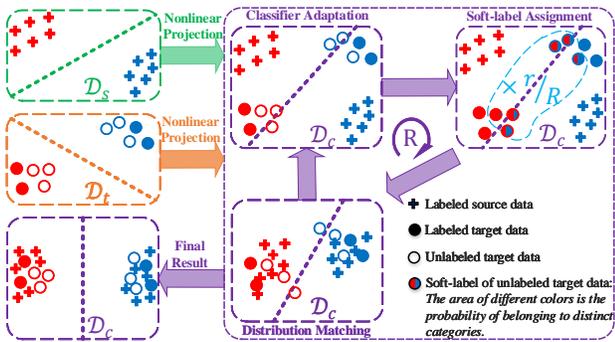
**Figure 2: An illustrative example of STN. Here, $\mathcal{D}_s$, $\mathcal{D}_t$, and $\mathcal{D}_c$ denote the source domain (in green), the target domain (in orange), and the common subspace (in purple), respectively. $r/R$ is the coefficient of the soft-label strategy to adaptively increase its importance, where $R$ is the total number of iterations and $r$ is the index of current iteration. Three key strategies are designed in STN for learning a domain-invariant $\mathcal{D}_c$. The classifier adaptation strategy aligns the discriminative directions of domains. The soft-label strategy avoids the hard assignment of each unlabeled target data to only one class that may be incorrect. The distribution matching strategy reduces the distributional divergence.**

of methods does not directly align the discriminative directions of domains. Moreover, several works [19, 37] derive the projections by iteratively performing the classifier adaptation and distribution matching. However, their performance is not stable since the iterative combination is a bit heuristic. In addition, they use trained classifier to assign each unlabeled target data to only one class that may be incorrect during the distribution matching, which is risky to negative transfer. Although Xiao and Guo [45] jointly consider the classifier adaptation and distribution matching in a unified framework, it does not further reduce the divergence between conditional distributions, which may be more important to the classification performance than aligning the marginal distributions. Note that all the above methods are shallow learning models that cannot learn more powerful feature representations. Although many deep learning models [4, 12, 24, 25, 28] are developed for domain adaptation, they cannot be directly applied to solve the HDA problem. Only few works [6, 33] solve the HDA problem by utilizing deep learning models. However, [33] requires the co-occurrence cross-domain data and [6] does not consider minimizing the distributional divergency across domains.

To overcome the limitations of the above methods, we propose a *Soft Transfer Network* (STN) for solving the HDA problem. The core rationale is depicted in Figure 2. In STN, we first project both the source and target data into a common subspace by using multiple nonlinear transformations. Then, a domain-invariant subspace is learned by jointly performing three key strategies: classifier adaptation, soft-label assignment, and distribution matching. The first stategy trains a domain-shared classifier with projected labeled source and target data by aligning the discriminative directions of domains. The second one assigns the *soft-label* (*i.e.*, a probability distribution over all the categories) to each unlabeled target data, which prevents the hard assignment of each unlabeled target data to only one class that may be incorrect. In addition, we observe that

the soft-labels will be more and more accurate as more iterations are executed. Thus, an adaptive coefficient is used to gradually increase the importance of the soft-labels. The last one matches both the marginal and conditional distributions between domains. Finally, we obtain a domain-shared classifier and a domain-invariant subspace for categorizing unlabeled target data.

The major contributions of this paper are three-fold. **(1)** To the best of our knowledge, there is no existing work that utilizes the *soft-label* strategy of unlabeled target data to address the HDA problem. **(2)** We propose a STN to learn a domain-shared classifier and a domain-invariant subspace in an end-to-end manner. **(3)** Extensive experimental results are presented on the transfer tasks of image-to-image, text-to-image, and text-to-text to verify the effectiveness of the proposed STN.

## 2 RELATED WORK

As aforementioned, existing HDA models either project data from one domain to the other or derive a domain-invariant subspace. Accordingly, they can be grouped into two categories, namely asymmetric transformation method and symmetric transformation method.

**Asymmetric transformation method**. Kulis *et al.* [21] present an Asymmetric Regularized Cross-domain Transformation (ARC-t) to project data from one domain to the other by learning an asymmetric nonlinear transformation. Zhou *et al.* [52] introduce a Sparse Heterogeneous Feature Representation (SHFR) to project the weights of classifiers in the source domain into the target one via a sparse transformation. Hoffman *et al.* [17, 18] propose a Max-Margin Domain Transform (MMDT), which maps the target data into the source domain by training a domain-shared support vector machine. Similarly, Xiao and Guo [44] develop a Semi-Supervised Kernel Matching for Domain Adaptation (SSKMDA) to transform the target data to similar source data, where a domain-shared classifier is simultaneously trained. Tsai *et al.* [36] propose a Label and Structure-consistent Unilateral Projection (LS-UP) model. This method maps the source data into the target domain with the aim of reducing the distributional difference and preserving the structure of projected data. They also propose a Cross-Domain Landmark Selection (CDLS) [37], where the representative source and target data are identified while reducing the distributional divergence between domains. Very recently, a few HDA approaches are presented based on the theory of Optimal Transport (OT) [40]. For example, Ye *et al.* [50] develop a Metric Transporation on HEterogeneous REpresentations (MAPHERE) approach, which jointly learns an asymmetric transformation and an optimal transportation. Yan *et al.* [48] propose a Semi-supervised entropic Gromov-Wasserstein (SGW) discrepancy to transport source data into the target domain. Although this method utilizes the supervision information to guide the learning of optimal transport, it does not consider unlabeled target data when matching the conditional distributions between domains.

**Symmetric transformation method**. Shi *et al.* [32] develop a Heterogeneous spectral MAPping (HeMAP) to learn a pair of projections based on spectral embeddings. Wang and Mahadevan [41] propose a Domain Adaptation using Manifold Alignment (DAMA) to find projections by preserving both the topology of each domain

and the discriminative structure. Duan *et al.* [14] put forward a Heterogeneous Feature Augmentation (HFA) approach. This method works by first augmenting the projected data with the original data features and then training a domain-shared support vector machine with the augmented features. Subsequently, HFA is generalized to a semi-supervised extension named SHFA [23] for effectively leveraging unlabeled target data. Xiao and Guo [45] present a Subspace Co-Projection with ECOC (SCP-ECOC) for HDA to learn projections by both training a domain-shared classifier and reducing the discrepancy on marginal distributions. Yao *et al.* [49] develop a Semi-supervised Domain Adaptation with Subspace Learning (SDASL) to learn projections by simultaneously minimizing the classification error, preserving the locality information of the original data, and considering the manifold structure of the target domain. Hsieh *et al.* [19] introduce a Generalized Joint Distribution Adaptation (G-JDA) model. This method reduces the difference in both the marginal and conditional distributions between domains while learning projections. Yan *et al.* [47] present a Discriminative Correlation Analysis (DCA) to jointly learn a discriminative correlation subspace and a target domain classifier. Recently, Li *et al.* [22] put forward a Progressive Alignment (PA) approach. This method first learns a common subspace by dictionary-shared sparse coding, and then mitigates the distributional divergence between domains. Besides, some HDA methods are developed based on deep learning models. For instance, Shu *et al.* [33] propose a weakly-shared Deep Transfer Network (DTN) to address the learning of heterogeneous cross-domain data, but it requires the co-occurrence cross-domain data. Recently, Chen *et al.* [6] put forward a Transfer Neural Trees (TNT) for HDA, which jointly addresses cross-domain feature projection, adaptation, and recognition in an end-to-end network. However, TNT does not explicitly minimize the distributional divergence across domains.

Moreover, our work is related to a few homogeneous DA studies that utilize the pseudo-label strategy of unlabeled target data. Specifically, Long *et al.* [27] propose a Joint Distribution Adaptation (JDA) to align both the marginal and conditional distributions between domains. As a follow-up, in [26] they further develop a Adaptation Regularization based Transfer Learning (ARTL) framework. This method works by simultaneously minimizing the structural risk, reducing the distributional divergence, and acquiring the manifold information. Yan *et al.* [46] introduce a weighted MMD to mitigate the effect of class weight bias, and further propose a Weighted Domain Adaptation Network (WDAN) to address the learning of homogeneous cross-domain data. Zhang *et al.* [51] present a Joint Geometrical and Statistical Alignment (JGSA) to reduce both the distributional and geometrical divergence across domains. Wang *et al.* [42] propose a Manifold Embedded Distribution Alignment (MEDA) approach, which learns a domain-invariant classifier by both minimizing the structural risk and performing dynamic distribution alignment. Ding *et al.* [11] put forward a Graph Adaptive Knowledge Transfer (GAKT) to jointly learn target labels and domain-invariant features. However, these methods either learn a feature transformation shared by cross-domain data [26, 27, 42, 46], or learn two feature transformations with the same size to constrain the divergence between transformations [11, 51], which thus cannot directly deal with the HDA problem.
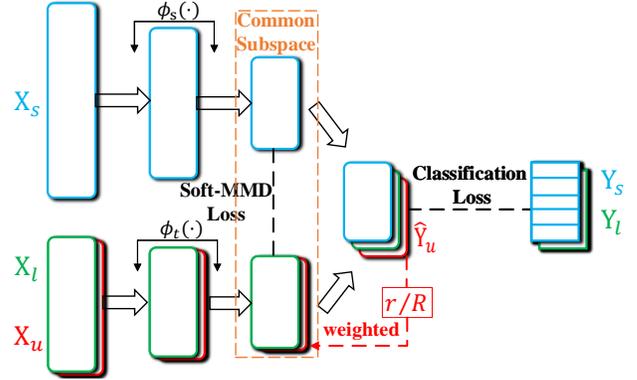


**Figure 3: The architecture of the proposed STN, where $\{X_s, Y_s\}$ (in blue) and $\{X_l, Y_l\}$ (in green) are labeled source and target data, respectively, $X_u$ (in red) is unlabeled target data, $\widehat{Y}_u = f(\phi_t(X_u))$ (in red) is the soft-label of $X_u$, $\phi_s(\cdot)$ and $\phi_t(\cdot)$ are feature projection networks of the source and target domains, respectively. Our model consists of two losses: a *classification loss* for aligning the discriminative directions of domains, and a *soft-MMD loss* to reduce the distributional divergence across domains.**

Overall, our work belongs to the symmetric transformation method. Inspired by the above-mentioned studies, we propose to jointly learn a domain-shared classifier and a domain-invariant subspace in an end-to-end manner, which not only aligns the discriminative directions of domains but also reduces the distributional divergence between domains.

## 3 SOFT TRANSFER NETWORK

In this section, we present the proposed STN. We begin with the definitions and terminologies. The source domain is denoted by $\mathcal{D}_s = \{(x_i^s, y_i^s)\}_{i=1}^{n_s}$, where $x_i^s \in \mathbb{R}^{d_s}$ is the $i$-th source domain data with $d_s$-dimensional features, and $y_i^s \in \mathcal{Y} = \{1, 2, \cdots, C\}$ is its associated class label with $C$ as the number of classes. Similarly, let $\mathcal{D}_t = \mathcal{D}_l \cup \mathcal{D}_u = \{(x_i^l, y_i^l)\}_{i=1}^{n_l} \cup \{x_i^u\}_{i=1}^{n_u}$ be the target domain, where $x_i^l (x_i^u) \in \mathbb{R}^{d_t}$ is the $i$-th labeled (unlabeled) target domain data with $d_t$-dimensional features, and $y_i^l \in \mathcal{Y}$ is the corresponding class label. Note that in the HDA problem, we have $d_s \neq d_t$, $n_s \gg n_l$, and $n_u \gg n_l$. The goal is to design a heterogeneous transfer network for predicting the labels of unlabeled target domain data in $\mathcal{D}_u$.

For simplicity of presentation, we denote by $X_s = [x_1^s, \cdots, x_{n_s}^s]^T \in \mathbb{R}^{n_s \times d_s}$ the data matrix in $\mathcal{D}_s$ and by $Y_s = [g_{y_1^s}, \cdots, g_{y_{n_s}^s}]^T \in \mathbb{R}^{n_s \times C}$ the corresponding label matrix, where $g_y \in \mathbb{R}^C$ denotes a one-hot vector with the $y$-th element being one. Analogously, let $X_l = [x_1^l, \cdots, x_{n_l}^l]^T \in \mathbb{R}^{n_l \times d_t}$ be the data matrix in $\mathcal{D}_l$ with an associated label matrix $Y_l = [g_{y_1^l}, \cdots, g_{y_{n_l}^l}]^T \in \mathbb{R}^{n_l \times C}$ and $X_u = [x_1^u, \cdots, x_{n_u}^u]^T \in \mathbb{R}^{n_u \times d_t}$ be the data matrix in $\mathcal{D}_u$. We also denote by $X_t = [X_l; X_u] \in \mathbb{R}^{n_t \times d_t}$ the data matrix in $\mathcal{D}_t$, where $n_t = n_l + n_u$ is the total number of the target domain data.

The key challenge of HDA arises in that the feature spaces of the source and target domains are different. Thus, we cannot directly transfer knowledge from the source domain to the target one. To

approach this problem, we first assume there is a common subspace with the dimensionality $d$, and then build two projection networks $\phi_s(\cdot) : \mathbb{R}^{d_s} \to \mathbb{R}^d$ and $\phi_t(\cdot) : \mathbb{R}^{d_t} \to \mathbb{R}^d$, which are used to transform source and target data into the subspace, respectively. We aim to find the optimal projection networks $\phi_s(\cdot)$ and $\phi_t(\cdot)$, which can form a domain-invariant subspace and lead to an effective transfer strategy. As depicted in Figure 3, we achieve this by optimizing an objective function with two components: a *classification loss* for aligning the discriminative directions of domains, and a *Soft-MMD loss* to reduce the distributional divergence across domains. The details of each component will be described in the following sections.

## 3.1 Classification Loss

We apply the classification loss to train a domain-shared classifier $f(\cdot)$ with projected labeled source and target data. It can be used to align the discriminative directions of domains. Further, we propose to jointly learn $f(\cdot)$, $\phi_s(\cdot)$, and $\phi_t(\cdot)$ under the *Structural Risk Minimization* (SRM) framework [39], and formulate the classification loss as

$$\mathcal{L}_c\left[Y_a, f(X_a)\right] = \frac{1}{n_s + n_l}\mathcal{L}\left[Y_a, f(X_a)\right] + \tau\left(\|\phi_s\|^2 + \|\phi_t\|^2 + \|f\|^2\right),$$
(1)

where $X_a = \left[\phi_s(X_s); \phi_t(X_l)\right] \in \mathbb{R}^{(n_s+n_l)\times d}$ denotes projected labeled data matrix in two domains, $Y_a = \left[Y_s; Y_l\right] \in \mathbb{R}^{(n_s+n_l)\times C}$ is the associated label matrix, $\mathcal{L}[\cdot, \cdot]$ is the cross-entropy function, $f(\cdot)$ is the softmax function, and $\tau$ is a positive regularization parameter. According to the SRM principle, the domain-shared classifier $f(\cdot)$ can accurately predict unknown target data by minimizing Eq. (1), given the domain-invariant subspace.

## 3.2 Soft-MMD Loss

As mentioned above, we expect to discover a domain-invariant subspace by enforcing the discriminative structure of the source domain to be consistent with that of the target domain. If there are sufficient labeled data in both domains, we can make it by only minimizing Eq. (1). However, the labeled target data is very scarce, which is not enough to identify a domain-invariant subspace. To eliminate this problem, we design a *Soft Maximum Mean Discrepancy* (Soft-MMD) loss, which can be applied to match both the marginal and conditional distributions between domains. As a result, the Soft-MMD loss is formulated as

$$\mathcal{L}_s\left[\phi_s(X_s), \phi_t(X_t)\right] = Q_m\left[\phi_s(X_s), \phi_t(X_t)\right] + Q_c\left[\phi_s(X_s), \phi_t(X_t)\right],$$
(2)

where $Q_m[\cdot, \cdot]$ acts as the divergence between marginal distributions of projected cross-domain data, and $Q_c[\cdot, \cdot]$ stands for that between conditional distributions.

We first detail how to formulate $Q_m[\cdot, \cdot]$. The empirical *Maximum Mean Discrepancy* (MMD) [15] has been proven to be a powerful tool for measuring the divergency between marginal distributions. Its exact idea is to calculate the distance between the centroids of both the source and target data, and then uses the distance to model the discrepance on marginal distributions. Thus, we define $Q_m[\cdot, \cdot]$ as

$$Q_m\left[\phi_s(X_s), \phi_t(X_t)\right] = \left\|\frac{1}{n_s}\sum_{i=1}^{n_s}\widetilde{x}_i^s - \frac{1}{n_t}\sum_{i=1}^{n_t}\widetilde{x}_i^t\right\|^2,$$
(3)

where $\widetilde{x}_i^s$ and $\widetilde{x}_i^t$ are the $i$-th projected source and target data, respectively.

Next, we describe how to design $Q_c[\cdot, \cdot]$. Since directly modelling the conditional distribution (*i.e.*, $\mathcal{P}(Y|X)$) is difficult, we turn to explore the class-conditional distribution (*i.e.*, $\mathcal{P}(X|Y)$) instead [19, 27, 37], and apply the divergency between class-conditional distributions to approximate that between conditional distributions. Hence, the divergence on conditional distributions can be approximated by calculating the sum of the distance between centroids of both the source and target data in each category. However, a large amount of target data has no label information. To make full use of unlabeled target data, [27] first proposes to use the *hard-label* strategy of unlabeled target data, which can be performed by applying a base classifier trained on labeled cross-domain data to unlabeled target data. Although this srategy may boost the adaptation performance, we note that it enforces each unlabeled target data to be assigned the class with the highest predicted probability, which may have two potential risks: (i) some unlabeled target data are aligned to incorrect class centroids since their hard-labels may be incorrect; and (ii) because of issue (i), negative transfer may occur. To alleviate these risks, we propose to adopt the *soft-label* strategy of unlabeled target data instead, which avoids the hard assignment of each unlabeled target data to only one class that may be incorrect. Concretely, the soft-label of the $i$-th projected unlabeled target data $\widetilde{x}_i^u$, *i.e.*, $\widehat{y}_i^u = f(\phi_t(x_i^u))$, is a $C$-dimensional vector, where the $k$-th element $\widehat{y}_{k,i}^u$ indicates the probability of $\widetilde{x}_i^u$ belonging to class $k$. We use $\widehat{y}_{k,i}^u$ as the weight of $\widetilde{x}_i^u$ to calculate the centroid of the $k$-th class in the target domain. Accordingly, a preliminary divergence between class-conditional distributions in the two domains is designed as

$$Q_c'\left[\phi_s(X_s), \phi_t(X_t)\right] = \sum_{k=1}^{C}\left\|\frac{1}{n_s^k}\sum_{i=1}^{n_s^k}\widetilde{x}_{k,i}^s - \frac{\sum_{i=1}^{n_l^k}\widetilde{x}_{k,i}^l + \sum_{i=1}^{n_u}\widehat{y}_{k,i}^u\widetilde{x}_i^u}{n_l^k + \sum_{i=1}^{n_u}\widehat{y}_{k,i}^u}\right\|^2,$$
(4)

where $\widetilde{x}_{k,i}^s$ and $\widetilde{x}_{k,i}^l$ are the $i$-th projected labeled source and target data of class $k$, respectively, and $n_s^k$, $n_l^k$ are the number of labeled source and target data of class $k$, respectively. If $\widehat{y}_{k,i}^u$ adopts the hard-label strategy as in [27], then Eq. (4) reduces to the hard-label approach [27]. Therefore, the proposed soft-label approach is a generalization of the hard-label approach [27].

Besides the *soft-label* strategy introduced in Eq. (4), we present another *iterative* weighting mechanism to further circumvent the risks mentioned above. We observe that the performance of $f(\cdot)$ gradually improves as the number of iterations increases. Thus, the value of $\widehat{y}_{k,i}^u$ will become more and more precise and reliable as more iterations are performed. As a result, we itroduce an adaptive coefficient to gradually increase the importance of $\widehat{y}_{k,i}^u$ during adaptation, and based on $Q_c'[\cdot, \cdot]$ formulate $Q_c[\cdot, \cdot]$ as

$$Q_c\left[\phi_s(X_s), \phi_t(X_t)\right] = \sum_{k=1}^{C}\left\|\frac{1}{n_s^k}\sum_{i=1}^{n_s^k}\widetilde{x}_{k,i}^s - \frac{\sum_{i=1}^{n_l^k}\widetilde{x}_{k,i}^l + \sum_{i=1}^{n_u}\alpha_i^{(r)}\widetilde{x}_i^u}{n_l^k + \sum_{i=1}^{n_u}\alpha_i^{(r)}}\right\|^2,$$
(5)

where $\alpha_i^{(r)}$ is defined by

$$\alpha_i^{(r)} = \frac{r * \widehat{y}_{k,i}^u}{R}, \tag{6}$$

which is used to adaptively increase the importance of $\widehat{y}_{k,i}^u$. Here, $R$ is the total number of iterations, and $r$ is the index of current iteration.

## 3.3 The Overall Objective of STN

In summary, to safely transfer knowledge across heterogeneous domains, the ideal model should be able to find the domain-shared classifier, *i.e.*, optimal $f(\cdot)$, and the domain-invariant subspace, *i.e.*, optimal $\phi_s(\cdot)$ and $\phi_t(\cdot)$. To this end, we integrate all of the components aforementioned into an end-to-end network, and obtain the overall objective function of STN:

$$\min_{\phi_s, \phi_t, f} \mathcal{L}_c \left[ \mathbf{Y}_a, f(\mathbf{X}_a) \right] + \beta \mathcal{L}_s \left[ \phi_s(\mathbf{X}_s), \phi_t(\mathbf{X}_t) \right], \tag{7}$$

where $\beta$ is a tradeoff parameter to balance the importance between $\mathcal{L}_c[\cdot, \cdot]$ and $\mathcal{L}_s[\cdot, \cdot]$. The proposed STN can jointly learn $f(\cdot)$, $\phi_s(\cdot)$, and $\phi_t(\cdot)$ in an end-to-end network by minimizing Eq. (7). Furthermore, STN adopts both the soft-label and the iterative weighting strategies to match the distributions. These are why STN can be expected to perform quite well as reported in the next section.

## 3.4 Comparison with Existing Studies

We now compare the proposed STN with some existing studies. To our knowledge, the methods related to STN include MMDT [17, 18], HFA [14], SHFA [23], G-JDA [19], CDLS [37], SGW [48], SCP-ECOC [45], JDA [27], JGSA [51], ARTL [26], MEDA [42], WDAN [46], and GAKT [11]. Among these methods, the first seven ones belong to heterogeneous DA techniques, while the last six methods are homogeneous DA ones. However, our work substantially distinguishes from them in the following aspects:

- **Comparison with heterogeneous DA studies.** **(1)** MMDT, HFA, and SHFA only consider the classifier adaptation. **(2)** G-JDA, CDLS, and SGW iteratively perform the classifier adaptation and distribution matching. Moreover, G-JDA and CDLS adopt the hard-label strategy to align the conditional distributions, and SGW does not take it into account. **(3)** SCP-ECOC only matches the marginal distributions.

- **Comparison with homogeneous DA studies.** **(1)** JDA and JGSA take the classifier adaptation and distribution matching as two independent tasks. **(2)** JDA, JGSA, ARTL, and MEDA utilize the hard-label strategy to reduce the conditional distribution divergence. **(3)** WDAN estimates the weights of the source data based on the hard-label strategy. **(4)** GAKT adopts the graph Laplacian regularization to perform the soft-label strategy and neglects the iterative weighting mechanism.

## 4 EXPERIMENTS

In this section, we empirically evaluate the proposed STN on the transfer tasks of image-to-image, text-to-image, and text-to-text.



**Figure 4: Sample images of the *bike* category from the Office+Caltech-256 dataset.**

## 4.1 Datasets and Settings

The **image-to-image** transfer task is performed on the **Office+Caltech-256** dataset [16, 31]. The former comprises 4,652 images with 31 categories collected from three distinct domains: Amazon (A), Webcam (W), and DSLR (D), and the latter includes 30,607 images of 256 objects from Caltech-256 (C). Following the settings in [6, 19, 37], we choose 10 overlapping categories of these two datasets to construct the Office+Caltech-256 dataset. Furthermore, we adopt two different feature representations for this type of task: 800-dimension *SURF* [3] and 4096-dimension *DeCAF*$_6$ [13]. Some sample images of the category of *bike* are shown in Figure 4.

The **text-to-image transfer** task is conducted on the **NUS-WIDE+ImageNet** dataset [7, 10]. The former consists of 269,648 images and the associated tags from Flickr, while the latter includes 5,247 synsets and 3.2 million images. We use the tag information of NUS-WIDE (N) and the image data of ImageNet (I) as the domains for text and image, respectively. Following [6], we select 8 shared categories of the two datasets to generate the NUS-WIDE+ImageNet dataset. Moreover, we follow (Chen et al. 2016) to extract the 4-th hidden layer from a 5-layer neural network as the 64-dimensional feature representation for text data, and to extract *DeCAF*$_6$ feature as the 4096-dimensional feature representation for image data.

The **text-to-text transfer** task is executed on the **Multilingual Reuters Collection** dataset [2]. This dataset contains about 11,000 articles from 6 categories written in five different languages: English (E), French (F), German (G), Italian (I), and Spanish (S). We follow [14, 19, 23, 37] to represent all the articles by BOW with TF-IDF features, and then to reduce the dimensions of features by performing PCA with 60% energy preserved. The final dimensions *w.r.t.* E, F, G, I, and S are 1,131, 1,230, 1,417, 1,041, and 807, respectively.

We implement the proposed STN based on the TensorFlow framework [1]. For the sake of fair comparison, we fix the parameter settings of STN for all the tasks. Both $\phi_s(\cdot)$ and $\phi_t(\cdot)$ are two-layer neural networks, which adopt Leaky ReLU [29] as the activation function. We use the Adam optimizer [20] with a learning rate of 0.001, and empirically set hyper-parameters $\beta = 0.001, \tau = 0.001$. In addition, the dimension of the common subspace $d$ and the number of iterations $R$ are set to 256 and 300, respectively.

## 4.2 Evaluations

We evaluate the proposed STN against eight state-of-the-art supervised learning and HDA methods, *i.e.*, SVMt, NNt, MMDT [17, 18], SHFA [23], G-JDA [19], CDLS [37], SGW [48], and TNT [6]. Among these methods, SVMt and NNt train a support vector machine and a neural network with only the labeled target data, respectively. MMDT, SHFA, G-JDA, CDLS, and SGW are the shallow HDA methods, while TNT is the deep HDA one. Furthermore, we note that SHFA, G-JDA, CDLS, SGW, and TNT are semi-supervised

Table 1: Classification accuracies (%) of all the methods on all the image-to-image transfer tasks.

| $\mathcal{D}_s \to \mathcal{D}_t$ | SVMt | NNt | MMDT | SHFA | G-JDA | CDLS | SGW | TNT | STN |
|---|---|---|---|---|---|---|---|---|---|
| C→A | 89.13±0.39 | 89.6±0.33 | 87.06±0.47 | 85.49±0.51 | 92.49±0.12 | 86.34±0.74 | 89.03±0.37 | 92.35±0.17 | **93.03±0.16** |
| W→A | | | 87±0.47 | 88.83±0.45 | 92.28±0.15 | 87.51±0.44 | 89.02±0.37 | 92.99±0.14 | **93.11±0.16** |
| A→C | 79.64±0.46 | 81.03±0.5 | 75.62±0.57 | 71.16±0.73 | 86.6±0.17 | 78.73±0.49 | 79.88±0.53 | 85.79±0.42 | **88.21±0.16** |
| W→C | | | 75.44±0.59 | 79.66±0.52 | 84.82±0.38 | 77.3±0.71 | 79.85±0.53 | 86.28±0.51 | **87.22±0.45** |
| A→W | 89.34±0.94 | 91.13±0.73 | 89.28±0.77 | 88.11±1.01 | 94.09±0.67 | 91.57±0.81 | 90.26±0.84 | 91.26±0.72 | **96.68±0.43** |
| C→W | | | 89.11±0.76 | 89.47±0.9 | 92.64±0.54 | 88.6±0.8 | 90.26±0.84 | 92.98±0.75 | **96.38±0.38** |
| A→D | 92.6±0.71 | 92.99±0.63 | 91.65±0.83 | 95.16±0.36 | 90.67±0.65 | 94.45±0.59 | 93.43±0.67 | 92.04±0.76 | **96.42±0.43** |
| C→D | | | 91.46±0.85 | 94.25±0.5 | 88.62±0.76 | 90.43±0.79 | 93.43±0.67 | 92.67±0.8 | **96.06±0.5** |
| W→D | | | 91.77±0.83 | 95.31±0.63 | 95.87±0.41 | 92.72±0.75 | 93.43±0.67 | 94.09±0.88 | **96.38±0.57** |
| Avg. | 87.68±0.63 | 88.69±0.55 | 86.49±0.68 | 87.49±0.62 | 90.9±0.43 | 87.52±0.68 | 88.73±0.61 | 91.17±0.44 | **93.72±0.36** |

Table 2: Classification accuracies (%) of all the methods on the text-to-image transfer task.

| $\mathcal{D}_s \to \mathcal{D}_t$ | SVMt | NNt | MMDT | SHFA | G-JDA | CDLS | SGW | TNT | STN |
|---|---|---|---|---|---|---|---|---|---|
| text→image | 66.85±0.96 | 67.68±0.8 | 53.21±0.69 | 64.06±0.61 | 75.76±0.65 | 70.96±0.83 | 68.01±0.8 | 77.71±0.57 | **78.46±0.58** |

HDA methods, which learn from both labeled and unlabeled cross-domain data. As stated in Section 2, homogeneous DA methods cannot directly apply to the HDA problem, thus we do not include them in the comparison. Following [6, 19, 23, 37, 48], we use the classification accuracy as the evaluation metric.

**Image-to-image transfer:** We first perform the task of image-to-image transfer on the Office+Caltech-256 dataset. In this kind of task, the source and target data are drawn from not only different feature representations but also distinct domains. As for the source domain, we take images in *SURF* features and utilize all images as the labeled data. For the target domain, we represent images with *DeCAF$_6$* features and randomly sample 3 images per class as the labeled data. The remaining images in the target domain are used as the testbed. Moreover, D is only viewed as the target domain due to the limited amount of images. Table 1 presents the average classification accuracies of 20 random trials.

From the results, we can make several insightful observations. **(1)** The proposed STN consistently achieves the highest accuracies on all the tasks. The average classification accuracy of STN is **93.72%**, which makes the improvement over the best supervised learning method, *i.e.*, NNt, and the best HDA method, *i.e.*, TNT, by **5.03%** and **2.55%**, respectively. These results clearly demonstrate the superiority of STN. **(2)** STN performs significantly better than MMDT and SHFA. The reason is that they only consider the classifier adaptation strategy but neglect the distribution matching strategy which can reduce the distributional divergence across domains. Although G-JDA and CDLS combine both strategies for HDA, their performance is still worse than STN. One reason is because they iteratively perform the strategies of classifier adaptation and distribution matching. The iterative combination is a bit heuristic and may lead to unstable performance. Another important reason is that they use the hard-label strategy of unlabeled target data during the distribution matching. The hard-labels may be not correct, which may result in limited performance improvement and negative transfer. The performance of SGW is worse than STN. An important reason is that SGW does not utilize unlabeled target data to align the conditional distributions between domains. STN yields better performance than TNT with one reason that TNT does not explicitly minimize the distributional discrepancy between domains. **(3)** MMDT and SHFA do not always outperform

the supervised learning method, *i.e.*, SVMt. One possible explanation is that the distributional divergence across domains may be large. Although they both train a domain-shared classifier to align the discriminative directions of domains, they do not minimize the distributional difference very well due to the limited amount of the labeled target data. Thus, it is risky to overfit the target data that result in negative transfer. In addition, SHFA exceeds MMDT with one major reason that the former utilizes the unlabeled target data while the latter does not. **(4)** G-JDA, CDLS, and SGW are better or comparable than SVMt. They iteratively perform the strategies of classifier adaptation and distribution matching, which confirms that these strategies are both meaningful and useful. In addition, we note that CDLS achieves worse performance than G-JDA. A possible explanation is that the instance weighting scheme used in CDLS is not always beneficial and may hurt the performance. **(5)** STN and TNT perform better than all the other methods, which verifies deep models are effective for addressing the HDA problem in the image domain.

**Text-to-image transfer:** We then conduct the text-to-image transfer task on the NUS-WIDE+ImageNet dataset. It is very hard because there is no co-occurrence text and image data for learning. According to [6], we treat the text (*i.e.*, N) and image (*i.e.*, I) datasets as the source and target domains, respectively. For the source domain, we choose 100 texts per category as the labeled data. As for the target domain, we randomly sample 3 images as the labeld data from each class, and the rest images are considered as the testbed. Table 2 shows the average classification accuracies in 20 random trials.

From the results, we can make the following important observations. **(1)** The proposed STN achieves the best performance on this task. The classification accuracy of STN is **78.46%**, which outperforms the best supervised learning method, *i.e.*, NNt, and the best HDA method, *i.e.*, TNT, by **10.78%** and **0.75%**, respectively. These results further corroborate the effectiveness of STN. **(2)** Both STN and TNT significantly exceed the other methods, which again verifies deep networks is helpful for learning heterogeneous cross-domain data. **(3)** The performance of MMDT is quite poor. This observation is consistent with [6] with one possible reason that the distributional divergence between domains may be very large since the source domain is textual while the target one is visual. **(4)** We

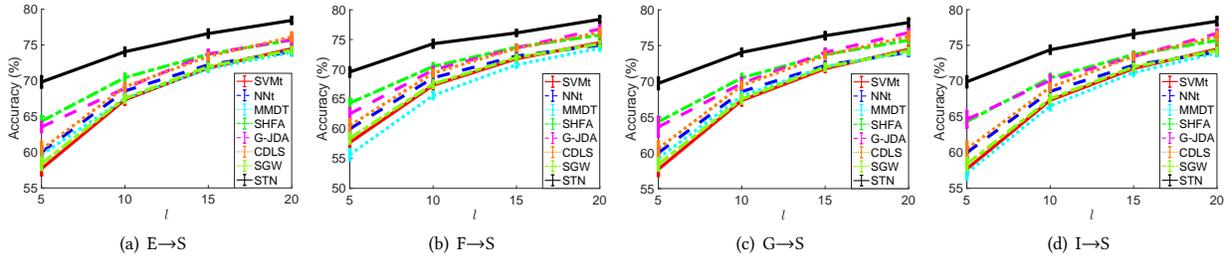**(a) E→S**  **(b) F→S**  **(c) G→S**  **(d) I→S**

**Figure 5: Classification accuracies (%) of all the methods with different numbers of labeled target data per category (*i.e.*, $l = 5, 10, 15, 20$) on all the text-to-text transfer tasks.**
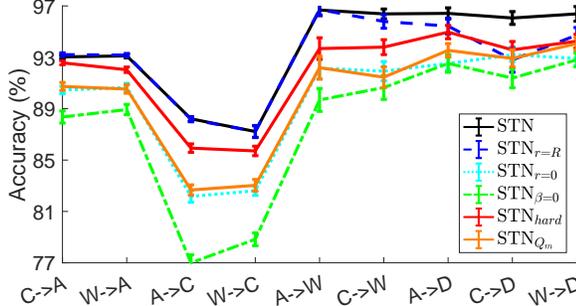


**Figure 6: Performance comparison of STN and its variants on all the image-to-image transfer tasks. Here, $\text{STN}_{r=R}$ removes the iterative weighting mechanism, $\text{STN}_{r=0}$ ignores the unlabeled target data, $\text{STN}_{\beta=0}$ ablates the soft-MMD loss, $\text{STN}_{hard}$ adopts the hard-label strategy of unlabeled target data, and $\text{STN}_{Q_m}$ neglects the divergence between conditional distributions.**

have the same observation as the image-to-image transfer task that SHFA and G-JDA are better than MMDT and CDLS, respectively.

**Text-to-text transfer:** We now execute the text-to-text transfer task on the Multilingual Reuters Collection dataset. In this type of task, we regard articles in distinct languages as data in different domains. According to the settings in [14, 19, 23, 37], we consider E, F, G and I as the source domains, and S as the target one. For the source domain, we randomly pick up 100 articles per category as the labeled data. As for the target domain, we randomly select $l$ (*i.e.*, $l = 5, 10, 15, 20$) and 500 articles per class as the labeled and unlabeled data, respectively. We explore how the number of labeled target data per category $l$ affects the performance and plot the average classification accuracies of 20 random trials in Figure 5. We do not report the results of TNT as they are much worse than the other methods (we note that the original paper of TNT also does not report the results on this type of task, please refer to [6] for details).

From the results, we have several interesting observations. **(1)** The performance of all the methods are boosted with the increase of the number of labeled target data. This observation is intuitive and reasonable. **(2)** The proposed STN substantially outperforms all the baseline methods on all the tasks. In particular, the average classification accuracy of STN with 5 labeled target data on all the tasks is **69.75%**, which outperforms the best supervised learning method, *i.e.*, NNt, and the best HDA method, *i.e.*, SHFA, by **9.75%** and **5.42%**, respectively. These results verify the effectiveness of STN again. **(3)** All the HDA methods can yield comparable or better

performance than the supervised learning method, *i.e.*, SVMt, which implies that these methods can produce positive transfer on all the text-to-text transfer tasks. **(4)** SHFA and G-JDA outperform MMDT and CDLS, respectively. We have the same observation as that on the image-to-image and text-to-image transfer tasks.

### 4.3 Analysis

**Variants Evaluation:** To go deeper with the efficacy of the iterative weighting mechanism, the distribution matching strategy, and the soft-label strategy, we evaluate several variants of STN: **(1)** $\text{STN}_{r=R}$, which removes the iterative weighting mechanism by setting $r = R$ in Eq. (6); **(2)** $\text{STN}_{r=0}$, which ignores the unlabeled target data by setting $r = 0$ in Eq. (6); **(3)** $\text{STN}_{\beta=0}$, which ablates the soft-MMD loss by setting $\beta = 0$ in Eq. (7); **(4)** $\text{STN}_{hard}$, which adopts the hard-label strategy of unlabeled target data and iteratively performs the hard-label assignment and objective optimizing in an end-to-end network; and **(5)** $\text{STN}_{Q_m}$, which neglects the divergence on the conditional distributions by ablating $Q_c$ in Eq. (2). We use the same experimental setting as above and plot the average classification accuracies on all the image-to-image transfer tasks in Figure 6. The results reveal several insightful observations. **(1)** As expected, STN delivers the best performance on all the tasks. **(2)** $\text{STN}_{r=R}$ is worse than STN, which suggests that the iterative weighting mechanism is useful to further improve the performance. $\text{STN}_{hard}$ is worse than STN and $\text{STN}_{r=R}$, which implies that the soft-label strategy is more suitable to match the conditional distributions between domains than the hard-label. $\text{STN}_{r=0}$ is worse than STN, $\text{STN}_{r=R}$, and $\text{STN}_{hard}$, which indicates that using the unlabeled target data can further improve the performance without resulting in negative transfer on these tasks. $\text{STN}_{Q_m}$ is worse than STN, $\text{STN}_{r=R}$, and $\text{STN}_{hard}$, which indicates that aligning the conditional distributions is useful. $\text{STN}_{Q_m}$ is similar to $\text{STN}_{r=0}$, which implies that using the unlabeled target data is important for reducing the divergence on conditional distributions. $\text{STN}_{\beta=0}$ yields the worst performance, which suggests that the distribution matching strategy is necessary and helpful for transferring knowledge across heterogeneous domains.

**Feature Visualization:** We use the t-SNE technique [38] to visualize the learned features of all the methods on the task of C→W except SHFA because it does not explicitly learn the feature projection matrices (please refer to [23] for details). We plot the visualization results in Figure 7. The results offer several interesting observations. **(1)** Comparing Figure 7(a) with Figure 7(b), we can see that the discriminative ability of $DeCAF_6$ is better than $SURF$, which is reasonable as $DeCAF_6$ is the deep feature. **(2)** For MMDT,
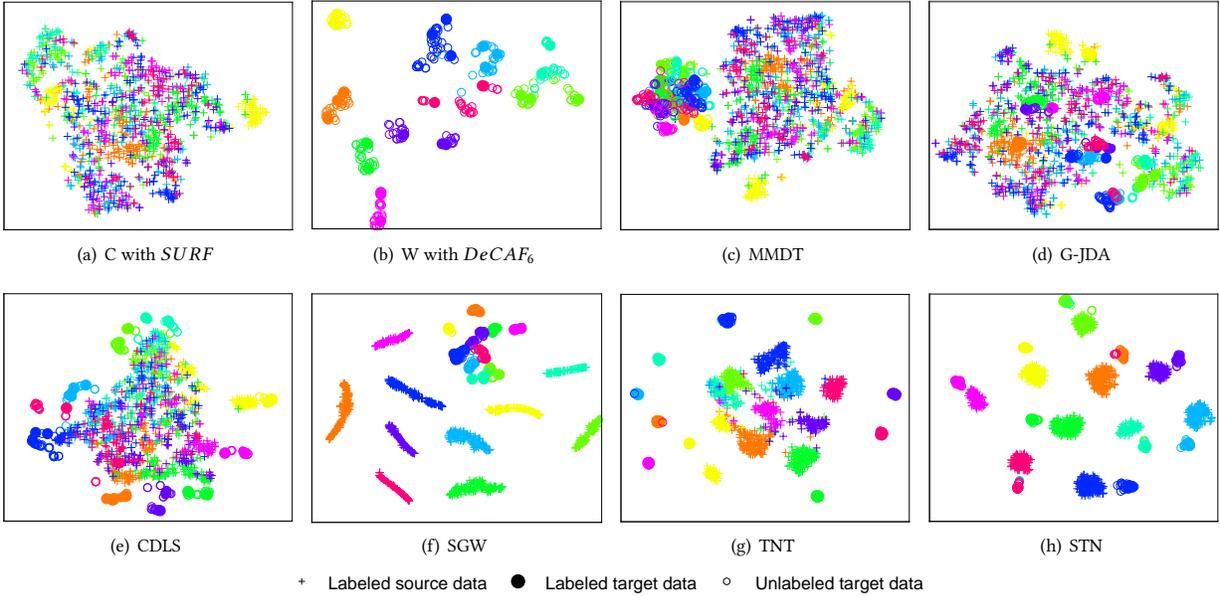
Figure 7: The t-SNE visualization on the C→W task. Here, (a) and (b) are the original feature representations of *SURF* and *DeCAF₆*, respectively. (c), (d), (e), (f), (g), and (h) are the learned feature representations of MMDT, G-JDA, CDLS, SGW, TNT, and STN, respectively.
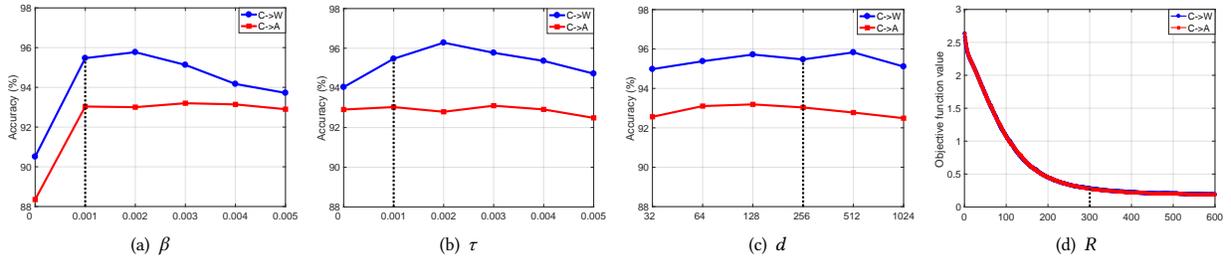


Figure 8: Parameter sensitivity and convergence analysis on the transfer tasks of C→W and C→A.

G-JDA, CDLS, SGW, and TNT, we can find that they do not align the distributions between domains very well, which explains their poor performance on this task. **(3)** The proposed STN matches the distributions between domains very well, which indicates that STN is powerful for transferring knowledge across heterogeneous domains.

**Parameter Sensitivity and Convergence:** We conduct experiments to analyze the parameter sensitivity and the convergence of STN on the transfer tasks of C→W and C→A. Figures 8(a)-8(c) show the accuracy *w.r.t.* $\beta$, $\tau$, and $d$, respectively. We can observe that the default values ($\beta = 0.001$, $\tau = 0.001$, $d = 256$) can achieve high accuracies. It is worth noting that STN yields the state-of-the-art performance on all the tasks by taking these default parameter values, which indicates STN is quite stable and effective. In addition, we plot the objective function value *w.r.t.* the number of iterations in Figure 8(d). We can see that the value of objective function first decreases and then tends to become steady as more iterations are executed, indicating the convergence of STN.

## 5 CONCLUSION

This paper proposes a STN to address the HDA problem, which jointly learns a domain-shared classifier and a domain-invariant subspace in an end-to-end way. Similar to many previous methods, STN aligns both the marginal and conditional distributions across domains. However, different from them, STN adopts the soft-label strategy of unlabeled target data to match the conditional distributions, which averts the hard assignment of each unlabeled target data. Furthermore, an adaptive coefficient is used to gradually increase the importance of the soft-labels. Experiments on three types of transfer tasks testify the effectiveness of STN. As a future direction, we plan to embed adversarial learning strategies into the STN by using, for instance, a domain discriminator. Another interesting direction is to apply the proposed Soft-MMD loss to other fields including homogeneous DA and few-shot learning.

# REFERENCES

[1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. TensorFlow: A System for Large-scale Machine Learning. In *OSDI*. 265–283.
[2] Massih Amini, Nicolas Usunier, and Cyril Goutte. 2009. Learning from Multiple Partially Observed Views - an Application to Multilingual Text Categorization. In *NIPS*. 28–36.
[3] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. 2006. Surf: Speeded up robust features. In *ECCV*. 404–417.
[4] Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. 2016. Domain Separation Networks. In *NIPS*. 343–351.
[5] Yue Cao, Mingsheng Long, and Jianmin Wang. 2018. Unsupervised Domain Adaptation With Distribution Matching Machines. In *AAAI*.
[6] Wei-Yu Chen, Tzu-Ming Harry Hsu, Yao-Hung Tsai, Yu-Chiang Frank Wang, and Ming-Syan Chen. 2016. Transfer Neural Trees for Heterogeneous Domain Adaptation. In *ECCV*.
[7] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. 2009. NUS-WIDE: A Real-world Web Image Database from National University of Singapore. In *CIVR*. 48:1–48:9.
[8] Gabriela Csurka. 2017. *A Comprehensive Survey on Domain Adaptation for Visual Applications*. Springer International Publishing, 1–35.
[9] Oscar Day and Taghi M. Khoshgoftaar. 2017. A survey on heterogeneous transfer learning. *Journal of Big Data* 4, 1 (26 Sep 2017), 29.
[10] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *CVPR*. 248–255.
[11] Zhengming Ding, Sheng Li, Ming Shao, and Yun Fu. 2018. Graph Adaptive Knowledge Transfer for Unsupervised Domain Adaptation. In *ECCV*. 36–52.
[12] Z. Ding, N. M. Nasrabadi, and Y. Fu. 2018. Semi-supervised Deep Domain Adaptation via Coupled Neural Networks. *TIP* 27, 11 (2018), 5214–5224.
[13] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. 2014. DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition. In *ICML*.
[14] Lixin Duan, Dong Xu, and Ivor W. Tsang. 2012. Learning with Augmented Features for Heterogeneous Domain Adaptation. In *ICML*. 711–718.
[15] Arthur Gretton, Karsten M. Borgwardt, Malte Rasch, Prof. Bernhard Schölkopf, and Alex J. Smola. 2007. A Kernel Method for the Two-Sample-Problem. In *NIPS*. 513–520.
[16] G. Griffin, A. Holub, and P. Perona. 2007. *Caltech-256 Object Category Dataset*. Technical Report 7694. California Institute of Technology.
[17] Judy Hoffman, Erik Rodner, Jeff Donahue, Brian Kulis, and Kate Saenko. 2014. Asymmetric and Category Invariant Feature Transformations for Domain Adaptation. *IJCV* 109, 1 (2014), 28–41.
[18] Judy Hoffman, Erik Rodner, Jeff Donahue, Kate Saenko, and Trevor Darrell. 2013. Efficient Learning of Domain-invariant Image Representations. In *ICLR*.
[19] Y. T. Hsieh, S. Y. Tao, Y. H. H. Tsai, Y. R. Yeh, and Y. C. F. Wang. 2016. Recognizing heterogeneous cross-domain data via generalized joint distribution adaptation. In *ICME*. 1–6.
[20] Diederick P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.
[21] B. Kulis, K. Saenko, and T. Darrell. 2011. What you saw is not what you get: Domain adaptation using asymmetric kernel transforms. In *CVPR*. 1785–1792.
[22] J. Li, K. Lu, Z. Huang, L. Zhu, and H. T. Shen. 2018. Heterogeneous Domain Adaptation Through Progressive Alignment. *TNNLS* (2018), 1–11.
[23] W. Li, L. Duan, D. Xu, and I. W. Tsang. 2014. Learning With Augmented Features for Supervised and Semi-Supervised Heterogeneous Domain Adaptation. *TPAMI* 36, 6 (2014), 1134–1148.
[24] M. Long, Y. Cao, Z. Cao, J. Wang, and M. I. Jordan. 2018. Transferable Representation Learning with Deep Adaptation Networks. *TPAMI* (2018), 1–1.
[25] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I. Jordan. 2015. Learning Transferable Features with Deep Adaptation Networks. In *ICML*. 97–105.

[26] M. Long, J. Wang, G. Ding, S. J. Pan, and P. S. Yu. 2014. Adaptation Regularization: A General Framework for Transfer Learning. *TKDE* 26, 5 (2014), 1076–1089.
[27] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu. 2013. Transfer Feature Learning with Joint Distribution Adaptation. In *ICCV*. 2200–2207.
[28] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I. Jordan. 2017. Deep Transfer Learning with Joint Adaptation Networks. In *ICML*. 2208–2217.
[29] Andrew L. Maas, Awni Y. Hannun, and Andrew Y. Ng. 2013. Rectifier nonlinearities improve neural network acoustic models. In *ICML Workshop on Deep Learning for Audio, Speech and Language Processing*.
[30] S. J. Pan and Q. Yang. 2010. A Survey on Transfer Learning. *TKDE* 22, 10 (2010), 1345–1359.
[31] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. 2010. Adapting Visual Category Models to New Domains. In *ECCV*. 213–226.
[32] X. Shi, Q. Liu, W. Fan, P. S. Yu, and R. Zhu. 2010. Transfer Learning on Heterogenous Feature Spaces via Spectral Transformation. In *ICDM*. 1049–1054.
[33] Xiangbo Shu, Guo-Jun Qi, Jinhui Tang, and Jingdong Wang. 2015. Weakly-Shared Deep Transfer Networks for Heterogeneous-Domain Knowledge Propagation. In *ACM MM*. New York, NY, USA, 35–44.
[34] Ben Tan, Yangqiu Song, Erheng Zhong, and Qiang Yang. 2015. Transitive Transfer Learning. In *KDD*. 1155–1164.
[35] Ben Tan, Yu Zhang, Sinno Jialin Pan, and Qiang Yang. 2017. Distant Domain Transfer Learning. In *AAAI*. 2604–2610.
[36] Y. H. H. Tsai, Y. R. Yeh, and Y. C. F. Wang. 2016. Heterogeneous domain adaptation with label and structure consistency. In *ICASSP*. 2842–2846.
[37] Y. H. H. Tsai, Y. R. Yeh, and Y. C. F. Wang. 2016. Learning Cross-Domain Landmarks for Heterogeneous Domain Adaptation. In *CVPR*. 5081–5090.
[38] Laurens van der Maaten. 2014. Accelerating t-SNE using Tree-Based Algorithms. *JMLR* (2014), 3221–3245.
[39] Vladimir N. Vapnik. 1998. *Statistical Learning Theory*. Wiley-Interscience.
[40] C. Villani. 2008. *Optimal Transport: Old and New*. Springer Berlin Heidelberg.
[41] Chang Wang and Sridhar Mahadevan. 2011. Heterogeneous Domain Adaptation Using Manifold Alignment. In *IJCAI*. 1541–1546.
[42] Jindong Wang, Wenjie Feng, Yiqiang Chen, Han Yu, Meiyu Huang, and Philip S. Yu. 2018. Visual Domain Adaptation with Manifold Embedded Distribution Alignment. In *ACM MM*. 402–410.
[43] Karl Weiss, Taghi M. Khoshgoftaar, and DingDing Wang. 2016. A survey of transfer learning. *Journal of Big Data* 3, 9 (2016).
[44] M. Xiao and Y. Guo. 2015. Feature Space Independent Semi-Supervised Domain Adaptation via Kernel Matching. *TPAMI* 37, 1 (2015), 54–66.
[45] Min Xiao and Yuhong Guo. 2015. Semi-supervised Subspace Co-Projection for Multi-class Heterogeneous Domain Adaptation. In *ECML PKDD*. 525–540.
[46] Hongliang Yan, Yukang Ding, Peihua Li, Qilong Wang, Yong Xu, and Wangmeng Zuo. 2017. Mind the Class Weight Bias: Weighted Maximum Mean Discrepancy for Unsupervised Domain Adaptation. In *CVPR*. 945–954.
[47] Yuguang Yan, Wen Li, Michael Ng, Mingkui Tan, Hanrui Wu, Huaqing Min, and Qingyao Wu. 2017. Learning Discriminative Correlation Subspace for Heterogeneous Domain Adaptation. In *IJCAI*. 3252–3258.
[48] Yuguang Yan, Wen Li, Hanrui Wu, Huaqing Min, Mingkui Tan, and Qingyao Wu. 2018. Semi-Supervised Optimal Transport for Heterogeneous Domain Adaptation. In *IJCAI*. 2969–2975.
[49] T. Yao, Yingwei Pan, C. W. Ngo, Houqiang Li, and Tao Mei. 2015. Semi-supervised Domain Adaptation with Subspace Learning for visual recognition. In *CVPR*. 2142–2150.
[50] Han-Jia Ye, Xiang-Rong Sheng, De-Chuan Zhan, and Peng He. 2018. Distance Metric Facilitated Transportation between Heterogeneous Domains. In *IJCAI*. 3012–3018.
[51] Jing Zhang, Wanqing Li, and Philip Ogunbona. 2017. Joint Geometrical and Statistical Alignment for Visual Domain Adaptation. In *CVPR*. 5150–5158.
[52] Joey Tianyi Zhou, Ivor W. Tsang, Sinno Jialin Pan, and Mingkui Tan. 2014. Heterogeneous Domain Adaptation for Multiple Classes. In *AISTATS*. 1095–1103.
[53] Yin Zhu, Yuqiang Chen, Zhongqi Lu, Sinno Jialin Pan, Gui-Rong Xue, Yong Yu, and Qiang Yang. 2011. Heterogeneous Transfer Learning for Image Classification. In *AAAI*. 1304–1309.