

Cycle Label-Consistent Networks for Unsupervised Domain Adaptation

Mei Wang, Weihong Deng

Abstract—Domain adaptation aims to leverage a labeled source domain to learn a classifier for the unlabeled target domain with a different distribution. Previous methods mostly match the distribution between two domains by global or class alignment. However, global alignment methods cannot achieve a fine-grained class-to-class overlap; class alignment methods supervised by pseudo-labels cannot guarantee their reliability. In this paper, we propose a simple yet efficient domain adaptation method, i.e. Cycle Label-Consistent Network (CLCN), by exploiting the cycle consistency of classification label, which applies dual cross-domain nearest centroid classification procedures to generate a reliable self-supervised signal for the discrimination in the target domain. The cycle label-consistent loss reinforces the consistency between ground-truth labels and pseudo-labels of source samples leading to statistically similar latent representations between source and target domains. This new loss can easily be added to any existing classification network with almost no computational overhead. We demonstrate the effectiveness of our approach on MNIST-USPS-SVHN, Office-31, Office-Home and Image CLEF-DA benchmarks. Results validate that the proposed method can alleviate the negative influence of falsely-labeled samples and learn more discriminative features, leading to the absolute improvement over source-only model by 9.4% on Office-31 and 6.3% on Image CLEF-DA.

Index Terms—Unsupervised domain adaptation, Cycle-consistency, Pseudo-label, Nearest centroid classification.

I. INTRODUCTION

Deep learning methods have propelled unprecedented advances in a wide variety of visual recognition tasks, such as image recognition [1] and object detection [2], demonstrating an excellent generalization ability. The recent success of deep learning heavily depends on large quantities of labeled data. However, collecting and annotating datasets for every new task and domain are extremely expensive and time-consuming processes, sufficient training data may not always be available. Training deep models on available labeled data from different but related source domains is a strong motivation to reduce the labeling consumption. But, analogously to other statistical machine learning techniques, this learning paradigm suffers from the domain shift problem [3], [4], which degrades the performance of pre-trained models when testing on target domains. In this regard, the research community has proposed different mechanisms, unsupervised domain adaptation (UDA) [5] is one of the promising methodologies to address this problem, which learns a good predictive model for the target domain using labeled examples from the source domain but only unlabeled examples from the target domain.

Mei Wang and Weihong Deng are with the Pattern Recognition and Intelligent System Laboratory, School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, 100876, China. E-mail: {wangmei1, whdeng}@bupt.edu.cn. (Corresponding Author: Weihong Deng)

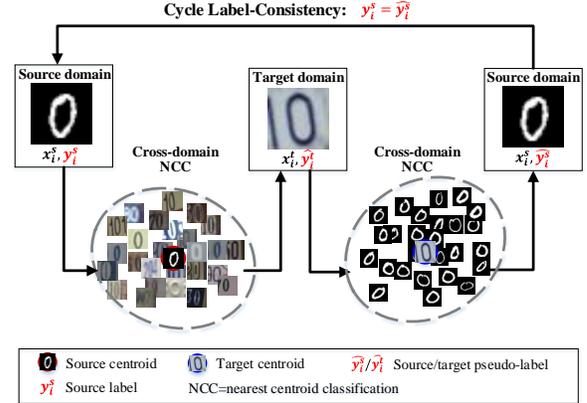


Fig. 1

ILLUSTRATION OF THE MAIN IDEA IN THIS WORK. IF WE CLASSIFY THE TARGET SAMPLES ACCORDING TO THE GROUND-TRUTH LABELS OF SOURCE SAMPLES, AND THEN CLASSIFY THE SOURCE SAMPLES ACCORDING TO PSEUDO-LABELS OF TARGET SAMPLES, WE SHOULD ASSIGN SOURCE SAMPLES BACK TO THE ORIGINAL LABELS.

Many deep UDA methods [5] try to match global distributions of source and target data to learn domain-invariant features, and then directly apply the classifier learned from only source labels to target domain [6], [7], [8], [9]. Other studies [10], [11], [12] recently propose to utilize pseudo-labels to take category information into consideration and learn target discriminative representations. These target pseudo-labels are generated by maximum posterior probability of source softmax classifier, and then are utilized by following two approaches. 1) These pseudo labels are directly used to finetune the network by supervised losses [13], [14]. However, due to large domain shift, the learned source classifiers might be incapable of precisely labeling target samples with an expected accuracy requirement. Supervising with these target pseudo-labels may lead to the error accumulation caused by some falsely pseudo-labeled samples. 2) The generated pseudo labels are used to align the centroids of source and target samples belonged to the same category [15], [10]. Although the target centroids computed by the mean vectors of embeddings of their constituent samples can alleviate the error generated by some mislabelled samples, these centroids correspondingly miss some backpropagation information provided by samples.

In this paper, we propose a simple yet efficient domain adaptation method, i.e. Cycle Label-Consistent Network (CLCN), exploiting cycle label-consistency and cross-domain nearest

centroid classification (NCC) algorithm to learn aligned and discriminative presentations for the target domain, as shown in Fig. 1. Different from other domain adaptation methods, CLCN is supervised with ground-truth labels and optimizes networks based on backpropagation information provided by each sample. During training, in addition to minimizing the cross-entropy loss on labeled source data, we reinforce the consistency between ground-truth labels and pseudo-labels of source samples leading to statistically similar latent representations between source and target domains. First, CLCN proposes to utilize cross-domain nearest centroid classification algorithm to estimate pseudo-labels for target domain and source domain in turn. The centroid of each source class is computed by the mean vector of embeddings of its constituent instances. With the help of source centroids, target pseudo-labels can be obtained by cross-domain classification which is performed for each target point by simply finding the nearest source centroids. Then, we update the target centroids in the same way as source centroids, and assign each source sample to its closest target centroid to get source pseudo-label. Such cross-domain classification enables to transport label information from the source domain to the target domain and back to the source domain. After obtaining source pseudo-labels, we reinforce the consistency between pseudo-labels and ground-truth labels of source samples by computing cross entropy between them so that statistically similar latent representations between source and target domains are learned at class level. Intuitively, with similar latent representations, the target pseudo-labels can be more accurately inferred from the labeled source samples.

Compared with other methods which directly use target pseudo-labels to finetune network, our CLCN just utilizes them to compute the target centroids and transport label information to source domain. The centroids computed by averaging the features of each class can alleviate the wrong information of several falsely-labeled samples. And cycle label-consistent loss is supervised with ground-truth labels of source samples, which makes our CLCN be more reliable and take full advantage of backpropagation information provided by each sample. Moreover, our source pseudo-labels here are “soft” (label probability) which are produced by computing a distribution over target pseudo-centroids for each source sample. Cycle consistency computed by “soft” pseudo-labels imposes a stricter restriction and obtains more compact clusters in latent space.

Following the standard evaluation protocol in the unsupervised domain adaptation community, we evaluate our method on the digit classification task using MNIST [16], SVHN [17] and USPS [18] as well as the object recognition task using the Office-31 [19], Office-Home [20] and ImageCLEF-DA dataset [21], and demonstrate the superiority of the proposed method. To summarize, the contributions of the paper are threefold:

- First, we propose a novel domain adaptation method, i.e. Cycle Label-Consistent Networks (CLCN), by exploring an intrinsic property that classification should be γ -cycle consistent, in the sense that if we classify the target samples according to the ground-truth labels of source samples, and then classify the source samples according

to pseudo-labels of target samples, we should assign source samples back to the original labels.

- Second, we formulate a cross-domain nearest centroid classification (NCC) algorithm to transport label information from the source domain to the target domain and back to the source domain. Especially, the label probability is computed by the cross-domain sample-to-centroid distances, and these “soft” pseudo-labels encourage stricter restriction and more compact clusters compared with “hard” pseudo-labels.
- Third, benefiting from being supervised with ground-truth labels, our CLCN can alleviate the negative influence of falsely-labeled samples and take full advantage of backpropagation information provided by each sample. We experimentally show that CLCN can achieve comparable performance with other complicated domain adaptation methods only by a source classification loss and a simple penalty item.

The rest of this paper is organized as follows. In the next section, we briefly review related work on deep unsupervised domain adaptation and cycle consistency. Section III describes the proposed method including cross-domain nearest centroid classification algorithm and cycle label-consistent loss. Additionally, experimental results are shown and analyzed in Section IV. Section V offers the concluding remarks.

II. RELATED WORK

A. Unsupervised domain adaptation

There is always a distribution change or domain gap between training and testing sets, which would degrade the performance. Yosinski et al. [22] comprehensively explored the transferability of deep neural networks and finetuned the network with sufficient target labeled data to improve performance on target domain. However, in practical scenario, labeled target data is usually limited or unacquirable. To address this issue, many UDA approaches [5] are proposed, as shown in Table II-A.

Global-alignment based method. Some papers explore domain-invariant feature spaces by minimizing global discrepancy measured by statistic loss [6], [23], [24] and adversarial loss [25], [7], [26]. Maximum mean discrepancy (MMD) [6], [23], [27], [28], Central Moment Discrepancy (CMD) [9] and Correlation alignment (CORAL) [24] are commonly-used statistic losses for UDA. For example, in Deep Domain Confusion (DDC) [6], the network is optimized by classification loss in the source domain, while domain difference is minimized by one adaptation layer with the MMD metric. Adversarial learning has shown great promise for use in domain adaptation. The domain-adversarial neural network (DANN) [7] integrated a gradient reversal layer (GRL) to train a feature extractor by maximizing the domain classifier loss and simultaneously minimizing the label predictor loss. Rahman et al. [29] proposed a correlation-aware adversarial domain adaptation framework where the correlation metric is used jointly with adversarial learning to minimize the domain disparity of the source and target data. Domain-symmetric networks (SymNets) overcomes the limitation in aligning the

TABLE I
COMPARISON OF EXISTING UNSUPERVISED DOMAIN ADAPTATION METHODS. THE AVERAGE ACCURACIES TESTED ON OFFICE-31 DATASET ARE REPORTED.

Category	Algorithm	Training method	Accuracy
Global alignment	DDC [6], DAN [23]	MMD	72.9 (Alexnet)
	JAN [35]	JMMD	76.0 (Alexnet)
	RTN [27]	MMD, residual learning	73.7 (Alexnet)
	WMMD [28]	weighted MMD	72.1 (Alexnet)
	CMD [9]	CMD	79.9 (Alexnet)
	Deep CORAL [24]	CORAL	72.1 (Alexnet)
	DANN [7]	adversarial learning	74.3 (Alexnet)
	CAADA [29]	adversarial learning, CORAL	78.3 (Alexnet)
	SymNet [36]	two-level adversarial learning	88.4 (Resnet50)
	Cicek et al. [37]	2K-way adversarial learning	-
	STA [30], UAN [31]	weighted adversarial learning	89.2 (Resnet50)
SAFN [38]	feature-norm alignment	87.6 (Resnet50)	
Local alignment	OT [8], SWD [39]	optimal transport	-
	Das et al. [32], [33], [34]	graph matching	-
Pseudo label	AsmTri [13]	tri-training, finetuning	-
	MSTN [15]	class alignment	79.1 (Alexnet)
	PFAN [10]	class alignment, progressive learning	80.4 (Alexnet)
	iCAN [14]	progressive learning, adversarial learning	87.2 (Resnet50)
	PACET [40]	progressive learning, weighted targets	76.6 (Alexnet)
	GCAN [41]	structure-, domain-, class-alignment	80.6 (Alexnet)

joint distributions of feature and category across domains via two-level domain confusion losses. Separate to Adapt (STA) [30] and Universal Adaptation Network (UAN) [31] use adversarial learning to align the shared classes of target domain and source domain, and reject samples of unknown classes to address open set DA problem.

Local-alignment based method. These methods [32], [33] place an emphasis on establishing a sample-to-sample correspondence between each source sample and each target sample to mitigate domain gap. Compared with global-alignment based methods, they consider the effect of each and every sample in the dataset explicitly. For example, Courty et al. [8] learned a transport plan for each source sample so that they are close to the target samples. Their transport plan is defined on a first-order, point-wise unary cost between each source sample and each target sample. Das et al. [32], [33], [34] developed a framework that exploit all the first-, second- and third-order relations to match the source and target samples along with a regularization using labels of the source data. Such higher order relations provide additional geometric and structural information about the data beyond the unary point-wise relations.

Pseudo-label based method. Target pseudo-labels [42], [40], [43] are utilized to finetune the network so that the lack of categorical information is compensated and discriminative representations are learned in the target domain. Saito et al. [13] introduced the idea of tri-training [44] into domain adaptation. Two different networks assign pseudo-labels to unlabeled samples through voting, another network is trained by these pseudo-labels to learn target discriminative representations. However, the generated pseudo-labels may be unre-

liable because of large domain shift. Supervising with these target pseudo-labels may lead to the error accumulation caused by falsely pseudo-labeled samples. Some methods propose to utilize progressive learning to tackle this error accumulation. Zhang et al. [14] iteratively selected confident pseudo-labels based on the classifier from the previous training epoch and re-trained the model by using the enlarged training set. Liang et al. [40] tackled the uncertainty in pseudo labels of the target domain from two aspects, progressive target instance selection and incorporating the learned class confidence scores to characterize both the within- and cross- domain relations. Other methods utilize pseudo labels to perform centroid alignment instead of finetuning to alleviate the negative effect of these falsely-labeled samples. Combining with adversarial loss, the moving semantic transfer network (MSTN) [15] took semantic information into consideration for unlabeled target samples by aligning labeled source centroids and pseudo-labeled target centroids. Chen et al. [10] combined easy-to-hard transfer strategy and centroid alignment, resulting in an improved target classification accuracy. Graph Convolutional Adversarial Network (GCAN) [41] utilized the concatenated CNN and GCN features to generate pseudo-labels and performed structure-aware alignment, domain alignment, and class centroid alignment to mitigate domain gap. However, these centroid alignment methods may correspondingly miss some backpropagation information provided by samples. In this paper, we exploit cycle label-consistency and cross-domain nearest centroid classification algorithm to address the error accumulation problem and take full advantage of back-propagation information, leading to improved performance on target domain.

B. Cycle consistency

The idea of using transitivity as a way to regularize structured data has a long history, for example, higher-order cycle consistency has been used in structure from motion [45], 3D shape matching [46], co-segmentation [47], dense semantic alignment [48], [49], and depth estimation [50].

In domain adaptation field, cycle-consistency was first introduced into deep domain adaptation by Zhu et al. [51]. Zhu et al. proposed Cycle-Consistent Adversarial Networks (CycleGAN) [51] which can translate one image from source domain X into target domain Y in the absence of any paired training examples. CycleGAN learns a mapping $G : X \rightarrow Y$ and an inverse mapping $F : Y \rightarrow X$. Two discriminators measure how realistic the generated image is ($G(X) \approx Y$ or $G(Y) \approx X$) by an adversarial loss and how well the original input is reconstructed after a sequence of two generations ($F(G(X)) \approx X$ or $G(F(Y)) \approx Y$) by a cycle-consistent loss. The dual-GAN [52] and the disco-GAN [53] were proposed at the same time, whose core idea is similar to CycleGAN. Cycle-Consistent Adversarial Domain Adaptation (CyCADA) [54] unifies prior feature-level and image-level adversarial domain adaptation methods together with cycle-consistent image-to-image translation techniques in which cycle-consistent loss encourages the cross-domain transformation to preserve local structural information. Associative domain adaptation [55] reinforces associations between source and target data directly in embedding space. It defines associative similarity as two-step transition probability of an imaginary random walker starting from an embedding of source domain and returning to another embedding via target domain embedding, and constraints the two-step probabilities to be similar to the uniform distribution over the class labels via a cross-entropy loss. Transduction with Domain Shift (TDS) [56] is most similar to our work, which combined metric learning and cycle consistency. It generates target pseudo-labels based on the k -nearest-neighbor rule, and the cycle consistency between ground-truth labels and pseudo-labels of source samples is enforced without explicitly computing source pseudo-labels but using the large-margin nearest neighbor (LMNN) [57] rule.

III. CYCLE LABEL-CONSISTENT NETWORK

A. Problem formulation

In unsupervised domain adaptation, we are given a set of labeled data from the source domain, and denote them as $\mathcal{D}_s = \{x_i^s, y_i^s\}_{i=1}^{N_s}$, where x_i^s is the i -th source sample, y_i^s is its category label, and N_s is the number of source images. A set of unlabeled data from the target domain is given as well and is denoted as $\mathcal{D}_t = \{x_i^t\}_{i=1}^{N_t}$, where x_i^t is the i -th target sample and N_t is the number of target images. The data distributions of two domains are different, $P(X_s, Y_s) \neq P(X_t, Y_t)$. We assume that the source and target domains contain the same object classes. Our goal is to learn a classifier $h : \mathcal{X} \rightarrow \mathcal{Y}$ (parameterized by Θ) that can decrease the domain discrepancy so as to minimize the target error by using the supervision information from source data. Generally, the classifier h is constructed as $h = g \circ f$ where f maps samples into features in the space \mathcal{F} and g outputs

TABLE II
SUMMARY OF MAJOR NOTATIONS USED IN THE PAPER.

Notations	Description
$\mathcal{D}_s/\mathcal{D}_t$	the source/target domain
x_i^s/x_i^t	sample in the source/target domain
N_s/N_t	number of source/target samples
y_i^s	true label of source sample x_i^s
\hat{y}_i^s/\hat{y}_i^t	pseudo label of source/target sample x_i^s/x_i^t
$x_{i,k}^s$	source data with true label k
$x_{i,\hat{k}}^s/x_{i,\hat{k}}^t$	source/target data with pseudo label \hat{k}
\mathcal{D}_k^s	the set of all source images with true label k
$\mathcal{D}_{\hat{k}}^s/\mathcal{D}_{\hat{k}}^t$	the set of source/target images with pseudo-label \hat{k}
c_k^s	centroid of source data with true label k
$c_{\hat{k}}^t$	centroid of target data with pseudo label \hat{k}
N_k^s	number of source samples with true label k
$N_{\hat{k}}^s/N_{\hat{k}}^t$	number of source/target samples with pseudo label \hat{k}
K	number of source or target categories
Θ	network parameters to be learned
α, θ	trade-off hyper-parameters

the predictions based on the extracted features. The learning process of our CLCN includes simultaneously optimizing the classifier h w.r.t. the labeled source data, and reinforcing the consistency between pseudo-labels and ground-truth labels of source samples. For the convenience of reading, we list some of the major notations that are used throughout this paper in Table II.

B. Cycle label-consistent network

In addition to minimizing the cross-entropy loss on labeled source data, our CLCN utilizes cross-domain nearest centroid classification algorithm to classify the target samples according to the ground-truth labels of source samples, and then classifies the source samples according to pseudo-labels of target samples, and finally reinforces consistency between ground-truth labels and pseudo-labels of source samples. By supervising with ground-truth labels, our cycle label-consistent loss alleviates the accumulation of errors and takes full advantage of backpropagation information provided by each sample, resulting in statistically similar latent representations between two domains and compact clusters in latent space. The overall architecture of CLCN is depicted in Fig. 2.

Source2target nearest centroid classification. This classification algorithm aims to estimate pseudo-labels in one domain with the help of centroids of classes in the other domain. First, for source domain, CLCN computes an M -dimensional representation c_k^s , or centroid, of k -th source class through the embedding function f with learnable parameters Θ . Each centroid is the mean vector of the embedded support points belonging to its class:

$$c_k^s = \frac{1}{N_k^s} \sum_{x_i^s \in \mathcal{D}_k^s} f(x_i^s) \quad (1)$$

where \mathcal{D}_k^s denotes the set of all source images with label k and N_k^s denotes the number of source samples in \mathcal{D}_k^s . We assume the existence of an embedding space in which the projections of samples in each class cluster around a single

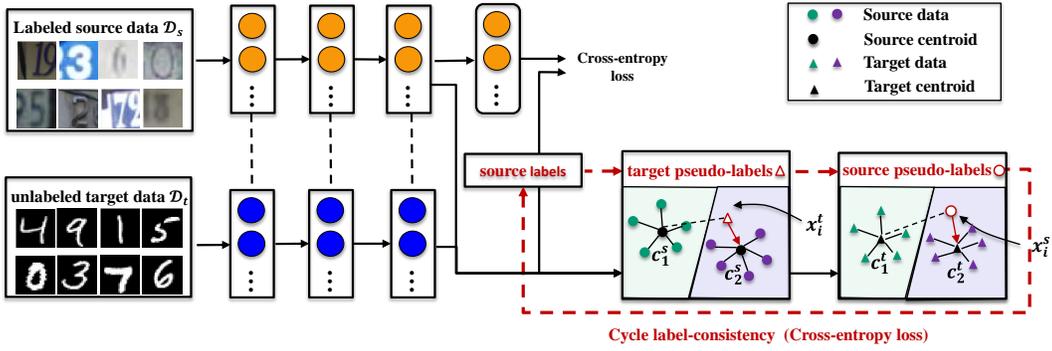


Fig. 2

THE OVERALL STRUCTURE OF OUR CLCN METHOD. THE INPUTS OF THE UPPER NETWORK ARE SOURCE LABELED IMAGES WHILE THOSE OF THE LOWER ARE TARGET UNLABELED DATA. IN ADDITION TO MINIMIZING THE CROSS-ENTROPY LOSS ON LABELED SOURCE DATA, CROSS-DOMAIN NEAREST CENTROID CLASSIFICATION ALGORITHM ESTIMATES TARGET PSEUDO-LABELS THROUGH SOURCE CENTROIDS AND ESTIMATES SOURCE PSEUDO-LABELS THROUGH TARGET PSEUDO-CENTROIDS IN TURN. THEN, CLCN REINFORCES THE CONSISTENCY BETWEEN THE GROUND-TRUTH LABELS AND THE PSEUDO-LABELS OF SOURCE SAMPLES RESULTING IN STATISTICALLY SIMILAR LATENT REPRESENTATIONS BETWEEN SOURCE AND TARGET DOMAINS AND COMPACT CLUSTERS IN LATENT SPACE.

centroid. Therefore, for each target sample x_i^t , we compute its cosine similarities with all source centroids, and assign it to corresponding centroid with the largest cosine similarity:

$$\hat{y}_i^t = \arg \max_k s_k \quad (2)$$

where, $s_k = \cos(f(x_i^t), c_k^s)$

where $\cos(\cdot, \cdot)$ denotes the cosine similarity function between two vectors. Finally, we can annotate each target node with pseudo label \hat{y}_i^t and add them into the \hat{k} -th pseudo-class of target domain $\mathcal{D}_{\hat{k}}^t$.

Target2source nearest centroid classification. After generating target pseudo-labels with the help of source centroids, we additionally generate source pseudo-labels with the help of target pseudo-centroids. The centroid of \hat{k} -th target pseudo-class, i.e. $c_{\hat{k}}^t$, are computed in the same way:

$$c_{\hat{k}}^t = \frac{1}{N_{\hat{k}}^t} \sum_{x_i^t \in \mathcal{D}_{\hat{k}}^t} f(x_i^t) \quad (3)$$

where $\mathcal{D}_{\hat{k}}^t$ denotes the set of all target images with pseudo-labels \hat{k} and $N_{\hat{k}}^t$ denotes the number of target samples in $\mathcal{D}_{\hat{k}}^t$. Then, the method of generating pseudo-labels for source domain is as follows. Given a source sample x_i^s , we compute the cross-domain sample-to-centroid distances between x_i^s and all target pseudo-centroids, and directly produce its score distribution $p_{score}(\hat{y}|x_i^s) \in \mathbb{R}^K$ over K target classes via a softmax function on these cross-domain sample-to-centroid distances:

$$p_{score}(\hat{y} = \hat{k}|x_i^s) = \frac{\exp(d(f(x_i^s), c_{\hat{k}}^t))}{\sum_{j=1}^K \exp(d(f(x_i^s), c_j^t))} \quad (4)$$

where $d(\cdot, \cdot)$ is the distance function between source sample and target centroid, and cosine distance is used in our paper; \hat{k} -th element of the score distribution, i.e. $p_s(\hat{y} = \hat{k}|x_i^s)$, is the probability of x_i^s belonging to pseudo-class \hat{k} . We treat

the score distribution $p_{score}(\hat{y}|x_i^s)$ as “soft” pseudo-labels \hat{y}_i^s of source samples.

Through such cross-domain nearest centroid classification algorithm, we finally transport label information from the source domain to the target domain and back to the source domain. Compared with generating pseudo-labels by softmax classifiers, our cross-domain nearest centroid classification algorithm can be performed easily with no extra structural and almost no training overhead, and can be more effective in cycle label-consistency, which has been proved in Section IV-D.

Cycle label-consistency. After classifying the target samples according to the ground-truth labels of source samples and then classifying the source samples according to pseudo-labels of target samples, we keep the consistency between predicted labels and ground-truth labels in source domain. Because such cycle label-consistency will be realized if the embedded features from both domains form well-aligned clusters. To establish this consistency, the model is optimized by minimizing the cross-entropy loss between pseudo-labels and ground-truth labels of source samples:

$$\mathcal{L}_{cyc} = -\frac{1}{N_s} \sum_{i=1}^{N_s} \sum_{k=1}^K \mathbf{1}[\hat{k} = y_i^s] \log p_{score}(\hat{y}|x_i^s) \quad (5)$$

where y_i^s and $p_{score}(\hat{y}|x_i^s)$ are ground-truth labels and “soft” pseudo-labels of source samples, respectively; $\mathbf{1}[\hat{k} = y_i^s]$ is 1 when $\hat{k} = y_i^s$, otherwise, it is 0. Compared with “hard” pseudo-labels of source samples which are generated in the same way as target pseudo-labels, minimizing \mathcal{L}_{cyc} computed by “soft” pseudo-labels imposes stricter constraint. It may encourage pseudo-labels and ground-truth labels to overlap rather than being nearest. Meanwhile, because “soft” pseudo-label indicates cross-domain sample-to-centroid distances, such cycle label consistency also clusters each source sample into its centroid leading to compact cluster per class.

C. Optimization

We can incorporate the cycle label-consistent loss into any popular deep convolutional neural networks (CNNs) architecture (e.g., AlexNet, VGG, ResNet, DenseNet, etc.) to learn robust features for unsupervised domain adaptation. Let $p_c(y = k|x_i^s)$ be the conditional probability that the CNN assigns x_i^s to k -th class. The classification loss on labeled source data, i.e. cross-entropy of ground-truth labels and network predictions, can be denoted as:

$$\mathcal{L}_C = -\frac{1}{N_s} \sum_{i=1}^{N_s} \sum_{k=1}^K \mathbf{1}[k = y_i^s] \log p_c(y|x_i^s) \quad (6)$$

where y_i^s is ground-truth labels of source samples. $\mathbf{1}[k = y_i^s]$ is 1 when $k = y_i^s$, otherwise, it is 0. In addition to minimizing the cross-entropy loss on labeled source data, we optimize cycle label-consistent loss to improve network performance on unlabeled target domain. Then, the final objective for our Cycle Label-Consistent Networks (CLCN) can be written as:

$$\mathcal{L} = \mathcal{L}_C + \alpha \mathcal{L}_{cyc} \quad (7)$$

where α is trade-off hyper-parameter. The cycle label-consistent loss learns aligned and compact embeddings for the source and target samples, while the source classification loss minimizes the prediction error of the source data. The entire procedure of computing total loss of CLCN is depicted in Algorithm 1.

Algorithm 1 Training episode loss computation for Cycle Label-Consistent Networks (CLCN).

Input:

Source domain labeled samples $\mathcal{D}_s = \{x_i^s, y_i^s\}_{i=1}^{N_s}$, and target domain unlabeled samples $\mathcal{D}_t = \{x_i^t\}_{i=1}^{N_t}$. Batch size N , the number of classes K , trade-off parameters α and θ .

Output:

The loss \mathcal{L} for a randomly generated training episode.

- 1: $\mathcal{S} = \text{RANDOMSAMPLE}(\mathcal{D}_s, N)$
 - 2: $\mathcal{T} = \text{RANDOMSAMPLE}(\mathcal{D}_t, N)$
 - 3: **Stage-1: // source2target nearest centroid classification**
 - 4: **for** $k = 1$ to K **do**
 - 5: $c_{k(t)}^s \leftarrow \frac{1}{|\mathcal{S}_k^s|} \sum_{x_i^s \in \mathcal{S}_k^s} f(x_i^s)$
 - 6: $c_k^s \leftarrow \theta c_{k(t)}^s + (1 - \theta) c_{k(t)}^s$
 - 7: **end for**
 - 8: Assign target samples to its closest source centroids according to Eqn. 2, and obtain target pseudo-labels
 - 9: **Stage-2: // target2source nearest centroid classification**
 - 10: **for** $k = 1$ to K **do**
 - 11: $c_k^t \leftarrow \frac{1}{|\mathcal{T}_k^t|} \sum_{x_i^t \in \mathcal{T}_k^t} f(x_i^t)$
 - 12: $c_k^t \leftarrow \theta c_k^t + (1 - \theta) c_{k(t)}^t$
 - 13: **end for**
 - 14: For each source sample, compute a distribution over target pseudo-centroids according to Eqn. 4, and obtain ‘‘soft’’ pseudo-labels
 - 15: **Stage-3: // cycle label-consistency**
 - 16: Reinforce the consistency \mathcal{L}_{cyc} between ground-truth labels and pseudo-labels of source samples according to Eqn. 5
 - 17: Compute softmax loss \mathcal{L}_C on source domain according to Eqn. 6
 - 18: $\mathcal{L} \leftarrow \mathcal{L}_C + \alpha \mathcal{L}_{cyc}$
-

In our CLCN, taking the entire training set into account and averaging the features of every class in each iteration are

inefficient even impractical. Therefore, we always use mini-batch SGD for optimization in practice. However, small batch size will lead to the huge deviation between the true centroid and local centroid (centroid of current-batch samples) caused by few mislabelled samples. To overcome this problem, global centroid [10], [15] is utilized to approach true centroid, so that more reliable pseudo-labels are obtained through our cross-domain nearest centroid classification. In each iteration, we first compute local centroid $c_{k(t)}^s$ for source domain using the mini-batch samples according to Eqn. 1, then update the global centroid c_k^s as follows:

$$c_k^s = \theta c_{k(t)}^s + (1 - \theta) c_k^s \quad (8)$$

where θ is the trade-off parameters. The global centroid c_k^t of target samples can be computed in the same way.

IV. EXPERIMENTS

In this section, we provide details about our implementation and training protocol, and report our experimental evaluation. Our algorithm is evaluated on various unsupervised domain adaptation tasks which focus on two different problems: handwritten digit classification and object recognition.

A. Datasets

MNIST-SVHN-USPS. We explore three digit datasets of varying difficulty as shown in Fig. 3: MNIST [16], SVHN [17] and USPS [18]. They contain digital images of 10 classes. MNIST and SVHN consist of grey images of size 28×28 and 16×16 , respectively; USPS consists of color images of size 32×32 , and has extra confusing digits around the centered digit of interest. Following previous works [25], we consider the three transfer tasks: MNIST→SVHN, SVHN→MNIST and MNIST→USPS. Digit images are cast to $28 \times 28 \times 1$ in all experiments for fair comparison.

Office-31 [19] is a standard dataset used for domain adaptation which contains 4110 images of 31 categories in total. This dataset contains three distinct domains, including Amazon (A) comprising 2817 images downloaded from online merchants, Webcam (W) involving 795 low resolution images acquired from webcams, and DSLR (D) containing 498 high resolution images of digital SLRs. We evaluate our method across all 6 transfer tasks. Data augmentation such as random flipping and cropping is used in training following JAN [35].

ImageCLEF-DA [21] is built for the ImageCLEF 2014 domain adaptation challenge with three domains, including Caltech-256 (C), ImageNet ILSVRC 2012 (I) and Pascal VOC 2012 (P). Each domain contains 12 classes and 50 images per class, which results in total 600 images for one domain. We also use data augmentation in training.

Office-Home [20] is a better organized but more difficult dataset than Office-31, which contains 15,500 images of 65 categories in total, forming 4 domains. Specifically, Art (Ar) denotes artistic depictions for object images, Clipart (Cl) means picture collection of clipart, Product (Pr) shows object images with a clear background and is similar to Amazon category in Office-31, and Real-World (Rw) represents object images collected with a regular camera. We use all domain combinations and build 12 transfer tasks.

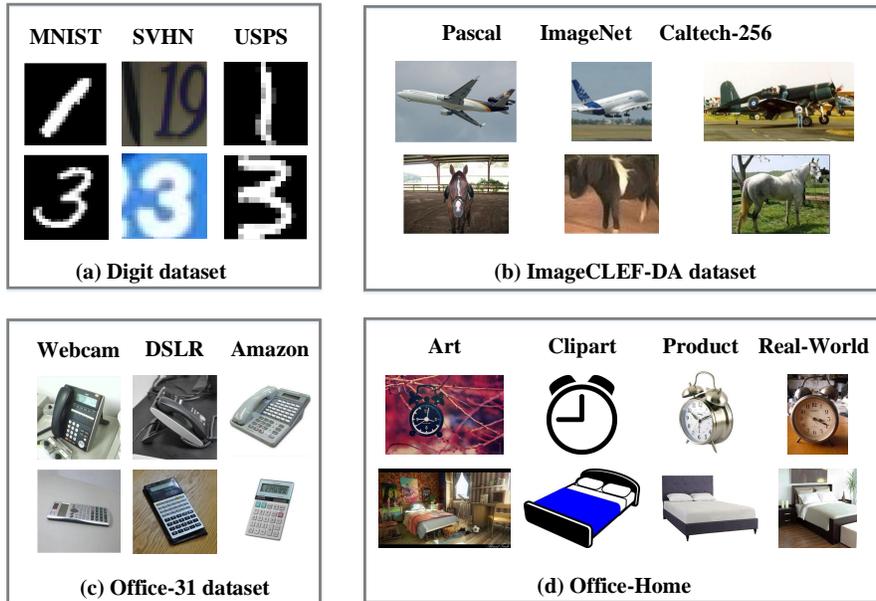


Fig. 3

DATASET SAMPLES OF (A) MNIST-SVHN-USPS [16], [17], [18], (B) IMAGECLEF-DA [21], (C) OFFICE-31 [19] AND (D) OFFICE-HOME [20].

B. Implementation detail

CNN architecture. For digit classification datasets, we use the same architecture with ADDA [25]: two convolution layers are followed by max pool layers and two fully connected layers are placed behind. Batch Normalization is inserted in convolutional layers. In experiments on Office-31, ImageCLEF-DA and Office-Home, we employ the AlexNet [1] pretrained on ImageNet [58] as our backbone networks. Following RTN [27] and RevGrad [59], a bottleneck layer $fc6$ with 256 units is added after the $fc7$ layer for safer transfer representation learning. We finetune the $conv1$, $conv2$, $conv3$, $conv4$, $conv5$, $fc6$, $fc7$ layers, and train the $fc6$ layer and last classifier from scratch.

Experimental setup. All the experiments were implemented using Tensorflow [60]. We use mini-batch stochastic gradient descent (SGD) with momentum of 0.9 to train the network and set weight decay to $5e-4$. We set the learning rate to 0.01 on most adaptation tasks except for MNIST→USPS, C→P and P→C adaptation tasks. The learning rate is set to be 0.001 on MNIST→USPS task, and is annealed by $\mu_p = \frac{\mu_0}{(1-\gamma^*p)^\beta}$ on C→P and P→C tasks, where $\mu_0 = 0.01$, $\gamma = 10$ and $\beta = 0.75$. In digit classification experiments, we set the batch size to 128. Following [61], we normalize the embedding features by 12 normalization and then re-scale them to 5. In experiments on Office-31, Office-Home and ImageCLEF-DA databases, we set the batch size to 400, 400 and 128, respectively. We also normalize the embedding features by 12 normalization and then re-scale them to 10.

Hyper-parameters. The hyper-parameter θ is set to be 0.7 in our method. In order to suppress noisy signal at the early stages of the training procedure, we change the parameter α from 0 to α_0 using the following strategy: $\alpha = \alpha_0 * \left(\frac{2}{1+\exp(-\gamma^*p)} - 1\right)$, where α_0 and γ are set to 2.5 and 10, and p is training progress changing from 0 to 1.

Evaluation protocols. We follow standard evaluation protocols for unsupervised domain adaptation as [23], [35], [59], [15]. All labeled source examples and all unlabeled target examples are used. We repeat each transfer task three times, and report the mean accuracy (number of correctly classified test samples divided by the total number of test samples) as well as the standard deviation of these three results.

C. Results

Results on MNIST-USPS-SVHN. We compare our method with Deep Domain Confusion (DDC) [6], Domain Separation Networks (DSN) [62], Gradient Reversal (RevGrad), Couple GAN (CoGAN) [63], Adversarial Discriminative Domain Adaptation (ADDA) [25], Label Efficient Learning (LEL) [64], Deep Reconstruction Classification Network (DRCN) [65], Domain Adversarial Adaptation neural network (DAA) [66], Asymmetric Tri-Training (AsmTri) [13], Moving Semantic Transfer Network (MSTN) [15] and TarGAN [67]. Table III shows the performance comparisons on three transfer directions among digit datasets.

From the results, we can see several important observations. 1) The performances of source-only model which is trained using only labeled source data could be regarded as a lower bound without domain adaptation. Source-only model gives the accuracies of 60.1%, 75.2% and 57.1% on the SVHN→MNIST, USPS→MNIST and MNIST→USPS adaptation tasks, respectively, showing deep network indeed suffers from the domain shift problem. 2) When matching global distribution of source and target domains through MMD or adversarial learning, DDC [6], RevGrad [59] and ADDA [25] outperform the source-only model on target domain, but this

TABLE III
CLASSIFICATION ACCURACIES (%) ON DIGIT RECOGNITIONS TASKS.

Methods	SVHN to MNIST	USPS to MNIST	MNIST to USPS
Source-only	60.1±1.1	75.2±1.6	57.1±1.7
DDC [6]	68.1±0.3	79.1±0.5	66.5±3.3
DSN w/ MMD [62]	72.2	-	-
RevGrad [59]	73.9	77.1±1.8	73.0±2.0
CoGAN [63]	-	91.2±0.8	-
ADDA [25]	76.0±1.8	89.4±0.2	90.1±0.8
DSN w/ DANN [62]	82.7	-	-
LEL [64]	81.0±0.3	-	-
DRCN [65]	82.0±0.1	91.8±0.09	73.7±0.04
DAA [66]	78.3±0.5	92.8±1.1	90.3±0.2
AsmTri [13]	86.0	-	-
MSTN [15]	91.7±1.5	92.9±1.1	-
TarGAN [67]	98.1	94.1	93.8
CLCN (ours)	97.5±0.1	98.5±0.1	94.4±0.3

TABLE IV
CLASSIFICATION ACCURACIES (%) ON OFFICE-31 DATASETS. (ALEXNET)

Methods	A to W	D to W	W to D	A to D	D to A	W to A	Average
AlexNet [1]	61.6±0.5	95.4±0.3	99.0±0.2	63.8±0.5	51.1±0.6	49.8±0.4	70.1
DDC [6]	61.8±0.4	95.0±0.5	98.5±0.4	64.4±0.3	52.1±0.6	52.2±0.4	70.6
DAN [23]	68.5±0.5	96.0±0.3	99.0±0.3	67.0±0.4	54.0±0.5	53.1±0.5	72.9
DRCN [65]	68.7±0.3	96.4±0.3	99.0±0.2	66.8±0.5	56.0±0.5	54.9±0.5	73.6
RevGrad [59]	73.0±0.5	96.4±0.3	99.2±0.3	72.3±0.3	53.4±0.4	51.2±0.5	74.3
RTN [27]	73.3±0.3	96.8±0.2	99.6±0.1	71.0±0.2	50.5±0.3	51.0±0.1	73.7
JAN [35]	74.9±0.3	96.6±0.2	99.5±0.2	71.8±0.2	58.3±0.3	55.0±0.4	76.0
AutoDIAL [68]	75.5	96.6	99.5	73.6	58.1	59.4	77.1
DAA [66]	74.3±0.3	97.1±0.2	99.6±0.2	72.5±0.2	52.5±0.3	53.2±0.1	74.8
PACET [40]	72.2	96.0	99.4	70.3	61.8	59.0	76.5
CAADA [29]	80.2	97.1	99.2	77.7	58.1	57.4	78.3
MSTN [15]	80.5±0.4	96.9±0.1	99.9±0.1	74.5±0.4	62.5±0.4	60.0±0.6	79.1
CLCN (ours)	78.4±0.4	97.6±0.1	99.9±0.2	73.9±0.2	64.3±0.2	62.8±0.1	79.5

improvement is limited. Because they are category agnostic leading to mismatch of label spaces across different domains. 3) When category information is incorporated, e.g. MSTN [15], the improvement becomes more significant. By aligning centroids of source and target classes, features in same class but different domains are mapped nearby, resulting in an improved target classification accuracy. 4) Our proposed CLCN achieves superior performances against other techniques. The accuracy of CLCN can achieve 97.5%, 98.5% and 94.4% on the SVHN→MNIST, USPS→MNIST and MNIST→USPS adaptation tasks, respectively, making the absolute improvement over MSTN [15] by 5.8% on SVHN→MNIST setting and 5.6% on MNIST→USPS setting. Compared with global alignment methods, CLCN takes consideration of category information leading to more similar latent representations between two domains; while compared with centroid alignment techniques, CLCN is supervised with ground-truth labels and takes full advantage of backpropagation information provided by each sample.

Results on Office-31. We also compare our method with DDC [6], DRCN [65], RevGrad [59], DAA [66], MSTN [15], Deep Adaptation Networks (DAN) [23], Residual Transfer Network (RTN) [27], Joint Adaptation Network (JAN) [35], Automatic Domain Alignment Layer (AutoDIAL) [68], Progressive learning with Confidence-weighted Targets (PACET) [40] and Correlation-aware adversarial domain adaptation (CAADA) [29] on Office-31 dataset. The results are reported in Table IV.

Observing the results in Table IV: 1) without adaptation, AlexNet [1] can not obtain perfect performance on target domain due to domain shift. 2) Existing domain-level alignment methods, such as RTN [27] and RevGrad [59], help to boost performance by reducing this discrepancy. 3) Our approach outperforms comparison methods on most adaptation tasks, which reveals that cycle label-consistency is effective and CLCN is scalable for different datasets. For example, we observe that CLCN is superior to JAN [35] significantly by about 3%. Compared with the best competitor, i.e. MSTN

TABLE V
CLASSIFICATION ACCURACIES (%) ON IMAGECLEF-DA DATASETS. (ALEXNET)

Methods	I to P	P to I	I to C	C to I	C to P	P to C	Average
AlexNet [1]	66.2±0.2	70.0±0.2	84.3±0.2	71.3±0.4	59.3±0.5	84.5±0.3	73.9
RTN [27]	67.4±0.3	81.3±0.3	89.5±0.4	78.0±0.2	62.0±0.2	89.1±0.1	77.9
RevGrad [59]	66.5±0.5	81.8±0.4	89.0±0.5	79.8±0.5	63.5±0.4	88.7±0.4	78.2
JAN [35]	67.2±0.5	82.8±0.4	91.3±0.5	80.0±0.5	63.5±0.4	91.0±0.4	79.3
MSTN [15]	67.3±0.3	82.8±0.2	91.5±0.1	81.7±0.3	65.3±0.2	91.2±0.2	80.0
CAADA [29]	67.8	84.5	91.7	81.3	63.9	91.8	80.2
CLCN (ours)	68.5±0.2	84.2±0.6	91.8±0.4	79.9±0.1	64.7±0.2	92.2±0.2	80.2

TABLE VI
CLASSIFICATION ACCURACIES (%) ON OFFICE-HOME DATASETS. (ALEXNET)

Source Target	Ar Cl	Ar Pr	Ar Rw	Cl Ar	Cl Pr	Cl Rw	Pr Ar	Pr Cl	Pr Rw	Rw Ar	Rw Cl	Rw Pr	Avg
GFK [69]	21.60	31.72	38.83	21.63	34.94	34.20	24.52	25.73	42.92	32.88	28.96	50.89	32.40
JDA [73]	25.34	35.98	42.94	24.52	40.19	40.90	25.96	32.72	49.25	35.10	35.35	55.35	36.97
CCSL [74]	23.51	34.12	40.02	22.54	35.69	36.04	24.84	27.09	46.36	34.61	31.75	52.89	34.12
LSC [11]	31.81	39.42	50.25	35.46	51.19	51.43	30.46	39.54	59.74	43.98	42.88	62.25	44.87
RTML [75]	27.57	36.20	46.09	29.49	44.69	44.66	28.21	36.12	52.99	38.54	40.62	57.80	40.25
JGSA [70]	28.81	37.57	48.92	31.67	46.30	46.76	28.72	35.90	54.47	40.61	40.83	59.16	41.64
PUnDA [71]	29.99	37.76	50.17	33.90	48.91	48.71	30.31	38.69	56.91	42.25	44.51	61.05	43.60
DAN [23]	30.66	42.17	54.13	32.83	47.59	49.78	29.07	34.05	56.70	43.58	38.25	62.73	43.46
DHN [20]	31.64	40.75	51.73	34.69	51.93	52.79	29.91	39.63	60.71	44.99	45.13	62.54	45.54
WDAN [28]	32.26	43.16	54.98	34.28	49.92	50.26	30.82	38.27	56.87	44.32	39.35	63.34	44.82
GAKT [72]	34.49	43.63	55.28	36.14	52.74	53.16	31.59	40.55	61.43	45.64	44.58	64.92	47.01
MSTN [15]	34.87	46.20	56.77	36.63	54.97	55.41	33.27	41.66	60.62	46.94	45.90	68.25	48.46
CLCN (ours)	37.58 ±0.2	49.39 ±0.3	57.70 ±0.2	37.08 ±0.1	55.31 ±0.4	56.24 ±0.1	34.80 ±0.1	39.85 ±0.4	61.03 ±0.2	48.63 ±0.1	46.14 ±0.2	68.93 ±0.1	49.39

[15], our CLCN obtains comparable results by a simpler and reliable way. MSTN performs domain-level alignment by adversarial learning and performs class-level alignment by minimizing the distances between source and target centroids with the same labels; while our CLCN only depends on a single penalty item, i.e. cycle label-consistent loss. The results demonstrate the effectiveness of our method which reinforces the consistency between ground-truth labels and pseudo-labels of source samples.

Results on ImageCLEF-DA. We compare our method with DAN [23], RTN [27], RevGrad [59], JAN [35], MSTN [15] and CAADA [29] on ImageCLEF-DA dataset. The results are reported in Table V. The improvement on ImageCLEF-DA is less than Office-31 since the difference in domain sizes will cause more serious shift [35]. The improvement of CLCN over Alexnet w.r.t. the average accuracy is 6.3%. Among the counterparts using domain adaptation algorithms, CAADA is the state-of-the-art approach. Our CLCN ranks within top two in both 5 out of 6 tasks, and it outperforms CAADA in most tasks. These convincing results on the challenging ImageCLEF-DA dataset indicate that our method has the potential to generalize to a variety of settings.

Results on Office-Home. On Office-Home dataset, we compare our CLCN with some shallow methods, e.g., Geodesic

Flow Kernel (GFK) [69], Joint Geometrical and Statistical Alignment (JGSA) [70], Probabilistic Unsupervised Domain Adaptation (PUnDA) [71], as well as deep methods, e.g. DAN [23], MSTN [15], Deep Hashing Network (DHN) [20], Weighted Domain Adaptation Network (WDAN) [28], Graph Adaptive Knowledge Transfer (GAKT) [72]. The results of Office-Home are reported in Table VI.

We can see that our CLCN outperforms the comparison methods on most transfer tasks, and exceeds the MSTN [15] about 1% by average. And we have the following observations. 1) Global alignment methods, i.e. DAN [23], can only obtain limited improvement on Office-Home datasets. Although benefiting from better representations of deep learning, deep methods, e.g. DAN [23], obtain a similar level of performance compared with shallow methods, e.g. PUnDA [71]. The reason may be that the four domains in Office-Home are with more categories, and are visually more dissimilar with each other. Since domain alignment is category agnostic in previous work, it is possible that the aligned domains are not classification friendly in the presence of large number of categories. This confirms that global alignment is not enough, and other constraints, e.g. compact clusters in latent space, are vital in domain adaptation problem. 2) Our CLCN yields larger boosts on Office-Home domain adaptation tasks. It takes into account

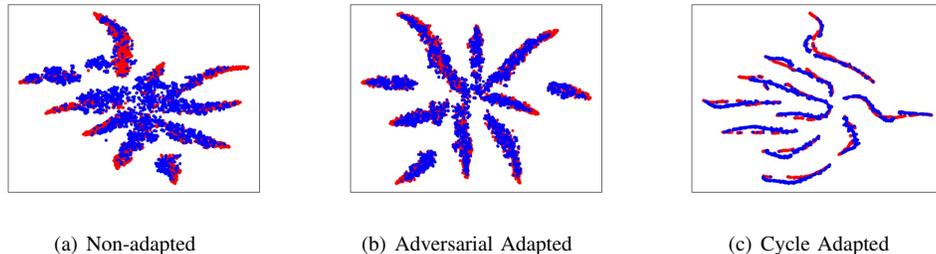


Fig. 4

FEATURE VISUALIZATION ON THE SVHN \rightarrow MNIST TASK. WE CONFIRM THE EFFECTS OF CLCN THROUGH A VISUALIZATION OF THE LEARNED REPRESENTATIONS USING T-DISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING (T-SNE) [76]. RED POINTS ARE SOURCE SAMPLES AND BLUE ARE TARGET SAMPLES. (A) IS TRAINED WITHOUT ANY ADAPTATION, (B) IS TRAINED WITH PREVIOUS ADVERSARIAL DOMAIN ADAPTATION METHOD, I.E. REVGRAD [59], (C) IS TRAINED WITH OUR CLCN METHOD. AS WE CAN SEE, COMPARED TO NON-ADAPTED METHOD, THE FEATURES GENERATED BY REVGRAD [59] ARE SUCCESSFULLY FUSED BUT ARE NOT DISCRIMINATED AND COMPACT. THE FEATURES NEAR CLASS BOUNDARY ARE OBVIOUSLY HARMFUL TO CLASSIFICATION TASKS. OUR CLCN METHOD ALIGNS FEATURE SPACE AND FORMS COMPACT CLUSTERS LEADING TO AN IMPROVED PERFORMANCE.

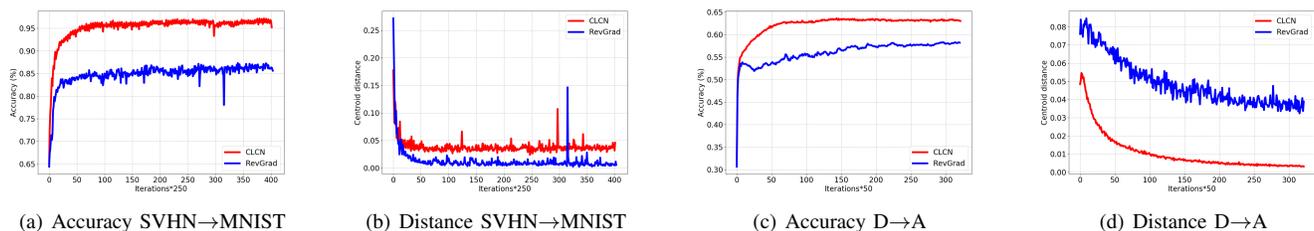


Fig. 5

REVGRAD [59] IN BLUE, OUR MODEL CLCN IN RED. (A)(C): COMPARISON OF TESTING ACCURACIES OF REVGRAD AND OUR PROPOSED METHOD CLCN ON SVHN \rightarrow MNIST AND D \rightarrow A ADAPTATION TASK. OUR MODEL HAS SIMILAR CONVERGENCE SPEED AS REVGRAD. (B)(D): THE DISTANCES BETWEEN SOURCE CENTROIDS AND TARGET CENTROIDS WITH THE SAME LABELS. FOR THE EASIER DIGIT TASKS, SUCH AS SVHN \rightarrow MNIST, THE CENTROID DISTANCES OF THE SAME CLASSES CAN BE BOTH REDUCED TO A SMALL VALUE IN REVGRAD AND CLCN MODEL; WHILE FOR THE HARDER TASKS, SUCH AS D \rightarrow A, THE CENTROID DISTANCES IN CLCN MODEL ARE MUCH SMALLER THAN THOSE IN REVGRAD MODEL.

the cycle consistency of classification across domain, resulting in higher similarity within the same class and better tightness of the clusters.

D. Experimental analysis

Feature visualization. To demonstrate the transferability of the CLCN learned features, the visualization comparisons are conducted at feature level. First, we randomly extract the deep features of source and target images in the SVHN \rightarrow MNIST task with source-only model, RevGrad [59] model and CLCN model, respectively. The features are visualized using t-distributed stochastic neighbor embedding (t-SNE) [76], as shown in Fig. 4. Fig. 4(a) shows the representations without any adaptation. As we can see, the distributions are separated between domains, which visually proves that there is domain shift between images of SVHN and those of MNIST. Fig. 4(b) shows the result of RevGrad method. Although features are successfully fused, the target points are not discriminated and compact. The ambiguous features are generated near class boundary which are obviously harmful to classification tasks.

Fig. 4(c) shows the representations of our CLCN method. We can see that the features with the same labels are concentrated and form tight clusters, and those from different classes are separated. Therefore, we conclude that the cycle label-consistency does help our CLCN to align feature space and form compact clusters at class level so that the target presentations are more discriminated and the performance of target domain is improved.

Convergence analysis. To inspect how CLCN converges, we show the test accuracy with respect to the number of iterations in Fig. 5(a) and 5(c). On SVHN \rightarrow MNIST and D \rightarrow A adaptation tasks, CLCN shows similar convergence rate with RevGrad [59] but better performance. Moreover, to verify that our cycle label-consistency loss can learn statistically similar latent representations between two domains, we further examine distances between source centroids and target centroids with the same labels. We compute global centroid of each source class and each target class according to Eqn. 8. The squared Euclidean distance $d(x, y) = \|x - y\|^2$ is utilized, so the total centroid distance can be formulated as: $d(c^s, c^t) = \sum_{k=\hat{k}=1}^K \|c_k^s - c_k^t\|^2$. The results of centroid

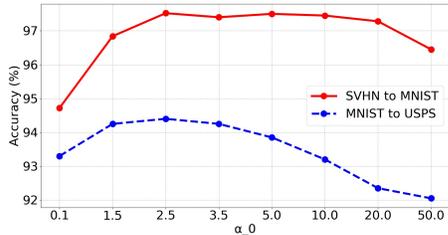


Fig. 6

PARAMETER SENSITIVITY STUDY ON THE SVHN→MMIST AND MNIST→USPS ADAPTATION TASKS.

distances with respect to the number of iterations are shown in Fig. 5(b) and 5(d). For the SVHN→MNIST adaptation task, the centroid distances can both be reduced to a small value in RevGrad [59] and CLCN model due to the simplicity of digit adaptation tasks. It is coincident with the observation in feature visualization in Fig. 4(b) and 4(c) where the source and target features are both successfully fused and aligned in RevGrad and CLCN model. However, our CLCN model is superior to RevGrad in the SVHN→MNIST task since CLCN additionally encourages each source sample to cluster into corresponding centroid so that more compact and discriminative presentations are learned. For the harder tasks, such as D→A, the centroid distances in CLCN model are much smaller than those in RevGrad model. The global alignment method, e.g. RevGrad model, can not effectively pass the class information to the adaptation network and can not guarantee that samples from different domains but with the same class label will map nearby in the feature space. Our CLCN model achieves stricter alignment resulting in better performance.

Parameter sensitivity. Our CLCN is optimized by source classification loss and cycle label-consistent loss, and achieves comparable performance with other complicated domain adaptation methods. As mentioned in Section IV-B, we set the trade-off hyper-parameter α as $\alpha_0 * \left(\frac{2}{1+\exp(-\gamma * p)} - 1 \right)$ for most experiments, where $\frac{2}{1+\exp(-\gamma * p)} - 1$ is the commonly-used strategy [15], [59] which gradually changes adaptation factor from 0 to 1 to suppress noisy signal at the early stages of the training procedure. To have a closer look at the parameter α_0 , we perform sensitivity analysis for it on the SVHN→MMIST and MNIST→USPS adaptation task by varying the parameter of interest in $\{0.1, 1.5, 2.5, 3.5, 5, 10, 20, 50\}$, and show the result in Fig. 6. From Fig. 6, we can see that the performance of our model remains largely stable across a wide range of parameter values on SVHN→MMIST tasks. And, the accuracy first increases and then decreases as α_0 varies and demonstrates a desirable bell-shaped curve. In principle, smaller values of α_0 could avoid CLCN learning perfect domain transfer features, and larger values of α_0 will weaken the effects of source classification loss and degrade the classification performance.

Ablation study. We perform an ablation study to investigate the effectiveness of two proposed techniques (i.e., cross-domain nearest centroid classification and cycle label-

TABLE VII

ABLATION STUDY ON SVHN→MNIST AND A→W ADAPTATION TASKS.

Methods	SVHN to MMIST	A to W
w/o NCC	70.6±0.2	65.5±0.3
w/o cycle	82.2±0.3	72.1±0.4
CLCN	97.5±0.1	78.4±0.4

consistency). Particularly, we introduce two different variants of CLCN, i.e. *softmax-based cycle (w/o NCC)* and *NCC-based finetuning (w/o cycle)*. In *softmax-based cycle (w/o NCC)* method, NCC is replaced by two softmax classifiers which are used to generate pseudo labels for two domains. After obtaining pseudo labels, it reinforces consistency between ground-truth labels and pseudo-labels of source samples. *NCC-based finetuning (w/o cycle)* method generates target pseudo-labels through source2target nearest centroid classification (NCC), and then directly utilizes these pseudo labels to finetune the network instead of performing cycle label-consistency. We show the comparison results in Table VII.

First, our CLCN is consistently superior to *NCC-based finetuning (w/o cycle)* method, which indicates the importance of such a cycle consistency technique. The *NCC-based finetuning (w/o cycle)* method only obtains 82.2% and 72.1% on the SVHN→MNIST and A→W tasks, respectively. Our CLCN achieves better performances and is superior to *NCC-based finetuning (w/o cycle)* by about 15.3% on the SVHN→MNIST task and 6.3% on the A→W task. Due to domain shift, the network trained with source data can not assign pseudo-labels for all target samples correctly. Some target samples lay far away from the source domain and they are ambiguous on the classification boundaries. These false-labeled samples introduce wrong information in *NCC-based finetuning (w/o cycle)* method and potentially result in the error accumulation. Rather than back-propagating the category loss based on each pseudo-labeled sample, our CLCN just utilizes the centroids of target pseudo-classes to transport label information to source domain. The centroids computed by averaging the features of each class can alleviate the wrong information of several falsely-labeled samples. When label information is transported back to source domain, cycle label-consistent loss is supervised with ground-truth labels of source samples, which makes our CLCN more reliable compared with *NCC-based finetuning (w/o cycle)* method.

Second, compared with our CLCN, *softmax-based cycle (w/o NCC)* method should additionally train a target softmax classifier to generate source pseudo labels, which is an extra burden. Moreover, we can see from the results that our CLCN outperforms *softmax-based cycle (w/o NCC)* by a large margin. Our hypothesis to explain this phenomenon is that the classification ability of softmax classifier is superior to NCC algorithm which makes it easier for target softmax classifier to overfit the target pseudo-labels generated by the source classifier. When target softmax classifier generates pseudo-labels for source domain, this overfitting results in high similarity between pseudo-labels and ground-truth labels of source samples so that cycle label-consistent loss can not

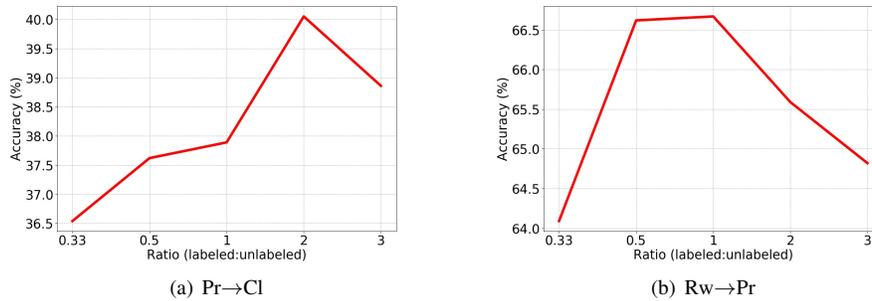


Fig. 7

ACCURACY ON PR→CL AND RW→PR TASKS WITH DIFFERENT RATIOS OF LABELED SOURCE SAMPLES TO UNLABELED TARGET SAMPLES.

TABLE VIII

CLASSIFICATION ACCURACIES (%) OF DIFFERENT FINETUNING METHODS ON THE SVHN→MNIST AND A→W ADAPTATION TASKS.

Methods	SVHN to MNIST	A to W
AsmTri ¹ [13]	86.0	-
softmax (finetuning)	86.6±2.7	72.5±1.2
NCC (finetuning)	82.2±0.3	72.1±0.4
NCC (cycle+finetuning)	98.5±0.1	75.4±0.6
CLCN (ours)	97.5±0.1	78.4±0.4

¹ AsmTri [13] utilizes two source softmax classifiers to generate target pseudo-labels through voting, then finetunes network with them.

obtain enough gradient to optimize the network.

Comparison with other finetuning methods. To verify the effectiveness of our CLCN, we additionally compare CLCN with other finetuning methods, e.g. *softmax-based finetuning*, on the SVHN→MNIST and A→W adaptation tasks, and show the results in Table VIII. Although *softmax-based finetuning* method performs better than *NCC-based finetuning* on these two tasks, it is still worse than our CLCN. Moreover, we additionally compare our CLCN with *NCC-based cycle and finetuning* method. It first generates pseudo labels by NCC, and then performs the cycle label-consistency and finetuning simultaneously. As we can see from Table VIII, *NCC-based cycle and finetuning* achieves about 1% gain over CLCN on the SVHN→MNIST adaptation task; but its result is unsatisfactory on the A→W task. This phenomenon is caused by the quality of target pseudo-labels. Benefiting from our CLCN, well aligned and compact presentations are learned in the easier task SVHN→MNIST. Therefore, most target samples are correctly labeled through target2source nearest centroid classification. Finetuning network with these reliable pseudo-labels is really helpful. When performing CLCN on the harder task A→W, the performance of target domain is still not perfect enough due to the large domain discrepancy. Even if finetuning method is combined with cycle label-consistency, more target samples will be wrongly labeled leading to the error accumulation and performance corrosion.

Sensitivity to the number of samples. In order to qualify how recognition performance changes as the number of sam-

ples is varied, we train our approach using different numbers of source and target samples on Pr→Cl and Rw→Pr tasks, respectively. We keep the total number of samples in two domains unchanged, and randomly select a part of data of each class from Office-home dataset. We vary the ratio of labeled source samples to unlabeled target samples ranging from 1/3 to 3, and observe the influence of different ratios on the recognition performance. The results are shown in Fig. 7. We can see that the accuracy first increases and then decreases as ratio varies and demonstrates a bell-shaped curve. First, without enough labeled data, the network cannot learn powerful representations and assign pseudo-labels correctly leading to poorer performance on target domain. Second, when the number of target data is small, less label information can be transported from target domain back to the source domain which weakens the effect of our cycle label-consistency.

V. CONCLUSION

In this paper, we propose a simple yet efficient method, i.e. Cycle Label-Consistent Networks (CLCN), for unsupervised domain adaptation. We exploit cycle label-consistency and cross-domain nearest centroid classification algorithm to learn statistically similar latent representations between source and target domains and regularize the latent space to form compact clusters at class level. Especially, “soft” pseudo-labels are generated to encourage stricter alignment and more compact clusters. Benefiting from being supervised with ground-truth labels, our CLCN can alleviate the negative influence of falsely-labeled samples without assistant of other technologies, and can take full advantage of backpropagation information provided by each sample. We experimentally show that CLCN optimized by classification loss and cycle label-consistent loss can achieve comparable performance with other complicated domain adaptation methods.

However, there are still some aspects to be improved. 1) Our method is somewhat heuristic and lacks some theoretical justification. We would like to discover some theoretical insights behind our method. 2) It highly relies on the closed set assumption where two domains share the same label space. In the future, we aim to extend it from the closed set setting to some challenging settings like open-set domain adaptation. 3) The quality of target pseudo-labels can be further improved.

We consider to use the easy-to-hard scheme which progressively selects reliable pseudo-labeled target samples in the future work.

VI. ACKNOWLEDGMENTS

This work was partially supported by National Key R&D Program of China (2019YFB1406504) and BUPT Excellent Ph.D. Students Foundation CX2020207.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [2] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [3] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, "Decaf: A deep convolutional activation feature for generic visual recognition," in *International conference on machine learning*, 2014, pp. 647–655.
- [4] A. Torralba, A. A. Efros *et al.*, "Unbiased look at dataset bias," in *CVPR*, vol. 1, no. 2. Citeseer, 2011, p. 7.
- [5] M. Wang and W. Deng, "Deep visual domain adaptation: A survey," *Neurocomputing*, vol. 312, pp. 135–153, 2018.
- [6] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell, "Deep domain confusion: Maximizing for domain invariance," *Computer Science*, 2014.
- [7] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *International Conference on Machine Learning*, 2015, pp. 1180–1189.
- [8] N. Courty, R. Flamary, D. Tuia, and A. Rakotomamonjy, "Optimal transport for domain adaptation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 9, pp. 1853–1865, 2017.
- [9] W. Zellinger, T. Grubinger, E. Lughofer, T. Natschläger, and S. Saminger-Platz, "Central moment discrepancy (cmd) for domain-invariant representation learning," *arXiv preprint arXiv:1702.08811*, 2017.
- [10] C. Chen, W. Xie, T. Xu, W. Huang, Y. Rong, X. Ding, Y. Huang, and J. Huang, "Progressive feature alignment for unsupervised domain adaptation," *arXiv preprint arXiv:1811.08585*, 2018.
- [11] C.-A. Hou, Y.-H. H. Tsai, Y.-R. Yeh, and Y.-C. F. Wang, "Unsupervised domain adaptation with label and structural consistency," *IEEE Transactions on Image Processing*, vol. 25, no. 12, pp. 5552–5562, 2016.
- [12] S. Li, S. Song, G. Huang, Z. Ding, and C. Wu, "Domain invariant and class discriminative feature learning for visual domain adaptation," *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4260–4273, 2018.
- [13] K. Saito, Y. Ushiku, and T. Harada, "Asymmetric tri-training for unsupervised domain adaptation," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 2988–2997.
- [14] W. Zhang, W. Ouyang, W. Li, and D. Xu, "Collaborative and adversarial network for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3801–3809.
- [15] S. Xie, Z. Zheng, L. Chen, and C. Chen, "Learning semantic representations for unsupervised domain adaptation," in *International Conference on Machine Learning*, 2018, pp. 5419–5428.
- [16] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner *et al.*, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [17] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *NIPS workshop on deep learning and unsupervised feature learning*, 2011, p. 5.
- [18] J. S. Denker, W. Gardner, H. P. Graf, D. Henderson, R. E. Howard, W. Hubbard, L. D. Jackel, H. S. Baird, and I. Guyon, "Neural network recognizer for hand-written zip code digits," in *Advances in neural information processing systems*, 1989, pp. 323–331.
- [19] K. Saenko, B. Kulis, M. Fritz, and T. Darrell, "Adapting visual category models to new domains," in *European conference on computer vision*. Springer, 2010, pp. 213–226.
- [20] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan, "Deep hashing network for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5018–5027.
- [21] B. Caputo, H. Müller, J. Martinez-Gomez, M. Villegas, B. Acar, N. Patricia, N. Marvasti, S. Üsküdarlı, R. Paredes, M. Cazorla *et al.*, "Imageclef 2014: Overview and analysis of the results," in *International Conference of the Cross-Language Evaluation Forum for European Languages*. Springer, 2014, pp. 192–211.
- [22] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" in *Advances in neural information processing systems*, 2014, pp. 3320–3328.
- [23] M. Long, Y. Cao, J. Wang, and M. I. Jordan, "Learning transferable features with deep adaptation networks," in *International Conference on Machine Learning*, 2015, pp. 97–105.
- [24] B. Sun and K. Saenko, "Deep coral: Correlation alignment for deep domain adaptation," in *European Conference on Computer Vision*. Springer, 2016, pp. 443–450.
- [25] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," *arXiv preprint arXiv:1702.05464*, 2017.
- [26] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko, "Simultaneous deep transfer across domains and tasks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 4068–4076.
- [27] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Unsupervised domain adaptation with residual transfer networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 136–144.
- [28] H. Yan, Y. Ding, P. Li, Q. Wang, Y. Xu, and W. Zuo, "Mind the class weight bias: Weighted maximum mean discrepancy for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2272–2281.
- [29] M. M. Rahman, C. Fookes, M. Baktashmotlagh, and S. Sridharan, "Correlation-aware adversarial domain adaptation and generalization," *Pattern Recognition*, p. 107124, 2019.
- [30] H. Liu, Z. Cao, M. Long, J. Wang, and Q. Yang, "Separate to adapt: Open set domain adaptation via progressive separation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2927–2936.
- [31] K. You, M. Long, Z. Cao, J. Wang, and M. I. Jordan, "Universal domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2720–2729.
- [32] D. Das and C. G. Lee, "Sample-to-sample correspondence for unsupervised domain adaptation," *Engineering Applications of Artificial Intelligence*, vol. 73, pp. 80–91, 2018.
- [33] —, "Graph matching and pseudo-label guided deep unsupervised domain adaptation," in *International Conference on Artificial Neural Networks*. Springer, 2018, pp. 342–352.
- [34] —, "Unsupervised domain adaptation using regularized hyper-graph matching," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 3758–3762.
- [35] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 2208–2217.
- [36] Y. Zhang, H. Tang, K. Jia, and M. Tan, "Domain-symmetric networks for adversarial domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5031–5040.
- [37] S. Cicek and S. Soatto, "Unsupervised domain adaptation via regularized conditional alignment," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1416–1425.
- [38] R. Xu, G. Li, J. Yang, and L. Lin, "Larger norm more transferable: An adaptive feature norm approach for unsupervised domain adaptation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1426–1435.
- [39] C.-Y. Lee, T. Batra, M. H. Baig, and D. Ulbricht, "Sliced wasserstein discrepancy for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 285–10 295.
- [40] J. Liang, R. He, Z. Sun, and T. Tan, "Exploring uncertainty in pseudo-label guided unsupervised domain adaptation," *Pattern Recognition*, vol. 96, p. 106996, 2019.
- [41] X. Ma, T. Zhang, and C. Xu, "Gcan: Graph convolutional adversarial network for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8266–8276.
- [42] M. Chen, K. Q. Weinberger, and J. Blitzer, "Co-training for domain adaptation," in *Advances in neural information processing systems*, 2011, pp. 2456–2464.

- [43] Y. Chen, C. Yang, Y. Zhang, and Y. Li, "Deep conditional adaptation networks and label correlation transfer for unsupervised domain adaptation," *Pattern Recognition*, vol. 98, p. 107072, 2020.
- [44] Z.-H. Zhou and M. Li, "Tri-training: Exploiting unlabeled data using three classifiers," *IEEE Transactions on Knowledge and Data Engineering*, vol. 17, no. 11, pp. 1529–1541, 2005.
- [45] C. Zach, M. Klopschitz, and M. Pollefeys, "Disambiguating visual relations using loop constraints," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, 2010, pp. 1426–1433.
- [46] Q.-X. Huang and L. Guibas, "Consistent shape maps via semidefinite programming," in *Computer Graphics Forum*, vol. 32, no. 5. Wiley Online Library, 2013, pp. 177–186.
- [47] F. Wang, Q. Huang, and L. J. Guibas, "Image co-segmentation via consistent functional maps," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 849–856.
- [48] T. Zhou, Y. Jae Lee, S. X. Yu, and A. A. Efros, "Flowweb: Joint image set alignment by weaving consistent, pixel-wise correspondences," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1191–1200.
- [49] T. Zhou, P. Krahenbuhl, M. Aubry, Q. Huang, and A. A. Efros, "Learning dense correspondence via 3d-guided cycle consistency," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 117–126.
- [50] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 270–279.
- [51] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [52] Z. Yi, H. Zhang, P. Tan, and M. Gong, "Dualgan: Unsupervised dual learning for image-to-image translation," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2849–2857.
- [53] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim, "Learning to discover cross-domain relations with generative adversarial networks," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 1857–1865.
- [54] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell, "Cycada: Cycle-consistent adversarial domain adaptation," *arXiv preprint arXiv:1711.03213*, 2017.
- [55] P. Haeusser, T. Frerix, A. Mordvintsev, and D. Cremers, "Associative domain adaptation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2765–2773.
- [56] O. Sener, H. O. Song, A. Saxena, and S. Savarese, "Learning transferable representations for unsupervised domain adaptation," in *Advances in Neural Information Processing Systems*, 2016, pp. 2110–2118.
- [57] K. Q. Weinberger and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," *Journal of Machine Learning Research*, vol. 10, no. Feb, pp. 207–244, 2009.
- [58] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [59] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," *arXiv preprint arXiv:1409.7495*, 2014.
- [60] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, "Tensorflow: A system for large-scale machine learning," in *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, 2016, pp. 265–283.
- [61] R. Ranjan, C. D. Castillo, and R. Chellappa, "L2-constrained softmax loss for discriminative face verification," *arXiv preprint arXiv:1703.09507*, 2017.
- [62] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan, "Domain separation networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 343–351.
- [63] M.-Y. Liu and O. Tuzel, "Coupled generative adversarial networks," in *Advances in neural information processing systems*, 2016, pp. 469–477.
- [64] Z. Luo, Y. Zou, J. Hoffman, and L. F. Fei-Fei, "Label efficient learning of transferable representations across domains and tasks," in *Advances in Neural Information Processing Systems*, 2017, pp. 165–177.
- [65] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, and W. Li, "Deep reconstruction-classification networks for unsupervised domain adaptation," in *European Conference on Computer Vision*. Springer, 2016, pp. 597–613.
- [66] X. Jia, Y. Jin, X. Su, and Y. Hu, "Domain-invariant representation learning using an unsupervised domain adversarial adaptation deep neural network," *Neurocomputing*, vol. 355, pp. 209–220, 2019.
- [67] F. Lv, J. Zhu, G. Yang, and L. Duan, "Targan: Generating target data with class labels for unsupervised domain adaptation," *Knowledge-Based Systems*, vol. 172, pp. 123–129, 2019.
- [68] F. M. Carlucci, L. Porzi, B. Caputo, E. Ricci, and S. R. Bulò, "Auto-dial: Automatic domain alignment layers," in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2017, pp. 5077–5085.
- [69] B. Gong, Y. Shi, F. Sha, and K. Grauman, "Geodesic flow kernel for unsupervised domain adaptation," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 2066–2073.
- [70] J. Zhang, W. Li, and P. Ogunbona, "Joint geometrical and statistical alignment for visual domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1859–1867.
- [71] B. Gholami, V. Pavlovic *et al.*, "Punda: Probabilistic unsupervised domain adaptation for knowledge transfer across visual categories," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3581–3590.
- [72] Z. Ding, S. Li, M. Shao, and Y. Fu, "Graph adaptive knowledge transfer for unsupervised domain adaptation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 37–52.
- [73] M. Long, G. Ding, J. Wang, J. Sun, Y. Guo, and P. S. Yu, "Transfer sparse coding for robust image representation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 407–414.
- [74] T. Ming Harry Hsu, W. Yu Chen, C.-A. Hou, Y.-H. Hubert Tsai, Y.-R. Yeh, and Y.-C. Frank Wang, "Unsupervised domain adaptation with imbalanced cross-domain data," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 4121–4129.
- [75] Z. Ding and Y. Fu, "Robust transfer metric learning for image classification," *IEEE Transactions on Image Processing*, vol. 26, no. 2, pp. 660–670, 2016.
- [76] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.