

Unsupervised Domain Adaptation using Deep Networks with Cross-Grafted Stacks

Jinyong Hou¹, Xuejie Ding², Jeremiah D. Deng¹, Stephen Cranfield¹

Department of Information Science, University of Otago¹

Institute of Information Engineering, Chinese Academy of Sciences²

robert.hou@postgrad.otago.ac.nz, dingxuejie@iie.ac.cn,

{jeremiah.deng, stephen.cranfield}@otago.ac.nz

Abstract

Current deep domain adaptation methods used in computer vision have mainly focused on learning discriminative and domain-invariant features across different domains. In this paper, we present a novel approach that bridges the domain gap by projecting the source and target domains into a common association space through an unsupervised “cross-grafted representation stacking” (CGRS) mechanism. Specifically, we construct variational auto-encoders (VAE) for the two domains, and form bidirectional associations by cross-grafting the VAEs’ decoder stacks. Furthermore, generative adversarial networks (GAN) are employed for label alignment (LA), mapping the target domain data to the known label space of the source domain. The overall adaptation process hence consists of three phases: feature representation learning by VAEs, association generation, and association label alignment by GANs. Experimental results demonstrate that our CGRS-LA approach outperforms the state-of-the-art on a number of unsupervised domain adaptation benchmarks.

1. Introduction

In machine learning, domain adaptation aims to transfer knowledge learned previously from one or more “source” tasks to a new but related “target” domain. As a special form of transfer learning, it helps to overcome the lack of labelled data in computer vision tasks by utilizing labelled data of the source domain and trying to automatically annotate unlabelled data in the target domain [28]. It may also be used to recognize unfamiliar objects in a dynamically changing environment in robotics. Therefore, in recent years domain adaptation, especially unsupervised domain adaptation, has become an appealing research topic [3, 2, 25, 12, 39, 32, 14].

For domain adaptation to occur, it is assumed that the

source and target domains are located in the same label space, but there is a domain bias. The challenge is to extract the domain-invariant representations from the data, and find an effective mechanism to overcome the domain bias and map the unlabelled targets to the label space.

To address the challenge, we propose to recruit different levels of deep unsupervised receptive fields from both the source and target domains and construct grafted representations for domain adaptation. Our approach is inspired by UNIT [21], but we generate the cross-domain association differently, employing grafted deep network layers. Specifically, we construct two parallel variational auto-encoders (VAEs) [17] to extract the latent encodings of the source and target. Then we recruit the different parts of the decoders to construct some cross-grafted representation stacks (CGRS), which produces bi-directional cross association between the two domains. Furthermore, generative adversarial networks (GANs) [11] are employed to carry out label alignment (LA), so that associations between the source and target contribute to accurate classification.

Due to these treatments our proposed CGRS-LA framework gives a promising direction for domain adaptation. Building cross associations between the domains, feature learning is hence achieved across domains, owning reduced domain dependency and increased domain-invariance, while adversarial networks further push feature representations away from the differences between domains, contributing to robust domain adaptation performance. Also, the cross-grafting process is entirely symmetric, leading to similar performance regardless of the adaptation direction, as revealed by our experiment results. Another advantage revealed by our experiments is that the CGRS is rather transferable across different tasks, which is an attractive trait for developing practical applications.

The rest of paper is organized as follows. In Section 2, we will briefly review some related work. In Section 3, we outline the overall structure of our proposed model, intro-

duce the CGRS scheme, and present the learning metrics used by the model. Finally the experimental results are presented in Section 4. We conclude the paper in Section 5, indicating our plan of future work.

2. Related Work

There are existing works that utilize intermediate feature representations to transfer previously learned knowledge to the target tasks. Self-taught learning [29] uses unsupervised learning trained on natural images to construct a sparse coding space, to which targets are projected to complete the recognition. In geodesic flow kernel [10, 13], the source and target datasets are embedded in a Grassman manifold, and a geodesic flow is constructed between the domains. A number of feature subspaces are sampled along the geodesic flow, and a kernel can be defined on the incremental feature vector, allowing a classifier to be built for the target dataset. DLID [7] uses deep sparse learning to extract the interpolated representation from a set of intermediate datasets constructed by combining the source and target datasets using progressively varying proportions, and the features from these intermediate datasets are concatenated to train a classifier.

Recent works have shown that deep networks involved in domain adaptation have achieved impressive performance due to their strong feature learning capacity. This provides a considerable improvement for some cross-domain recognition tasks [37, 23, 34, 26, 30, 21, 6, 9]. Specifically, a number of deep domain adaptation models have applied the adversarial training strategy [35, 36, 8, 21, 5, 20, 22]. DANN [8] employs a gradient reversal layer between the feature layer and the domain discriminator, causing feature representation to anti-learn the domain difference and hence adapt well to the target domain. ADDA [35] firstly trains a convolutional neural network (CNN) using the source dataset. An adversarial phase then follows, with the CNN assigned to the target for domain discriminator training, and the new target encoder CNN is finally combined with source classifier to achieve the adaptation.

Using generative adversarial networks (GAN), the PixeIDA framework [5] generates synthetic images from source-domain images that are mapped to the target domain. A task classifier then is trained by the source and synthetic images using the source labels. UNIT [21] introduces an unsupervised image-to-image translation framework based on couple of variational auto-encoders (VAEs) and GANs. To achieve this, a pair of corresponding images in different domains are mapped to a shared latent representation space.

Inspired by these previous works, our proposed CGRS-LA framework combines two ideas: constructing cross-domain feature representations, and employing adversarial networks for label alignment. Specifically, it incorporates VAEs to learn feature representations, a cross-grafting step

to generate bidirectional cross-domain associations, and a generative adversarial approach that carries out classification on source-target associations. A detailed descriptions of our framework are given next.

3. The CGRS-LA Framework

3.1. Model Description

We consider two domains: one is a source domain \mathcal{D}_s , which is constructed by n_s images $\mathbf{X}_s = \{\mathbf{x}_i^s\}_{i=1}^{n_s}$ and their correspond labels $\mathbf{y}_s = \{y_i^s\}_{i=1}^{n_s}$; the other is a target domain $\mathcal{D}_t = \{\mathbf{X}_t, \mathbf{y}_t\}$, where $\mathbf{X}_t = \{\mathbf{x}_i^t\}_{i=1}^{n_t}$ and their labels $\mathbf{y}_t = \{y_i^t\}_{i=1}^{n_t}$ are not available during adaptation. The source and target domain are drawn from joint distributions $\mathcal{P}(\mathbf{X}_s, \mathbf{y}_s)$ and $\mathcal{Q}(\mathbf{X}_t, \mathbf{y}_t)$, with a domain bias making \mathcal{P} and \mathcal{Q} different. Our goal is to learn some representations bearing similarity to both domains, i.e. some joint distribution between \mathcal{P} and \mathcal{Q} as a bridge for knowledge transfer, based on which the target images can be successfully classified.

Our framework is shown in Figure 1, split into five modular sub-tasks based on the ideas outlined as above. Firstly, in module *A*, the VAEs couple are implemented by CNNs. Both the encoders and decoders are divided into high and low level stacks. The high-level layers of the encoders are shared between domains. The source and target data are encoded to latent representation \mathbf{z}_s and \mathbf{z}_t , and then decoded to the reconstruction images $\hat{\mathbf{x}}_s$ and $\hat{\mathbf{x}}_t$ respectively. We assume that they have the same latent space, and the prior distribution is a normal one, $\mathcal{N}(0, I)$.

Secondly, the latent encodings pass through the cross-grafted stacks, forming cross-domain associations that are aligned to the label space. In module *B*, we construct two parallel CGRS by grafting the decoder stacks of the source and the target. Therefore, the cross-domain association images ($\mathbf{X}_s^{st}, \mathbf{X}_t^{st}, \mathbf{X}_s^{ts}, \mathbf{X}_t^{ts}$) are generated when the latent encodings from different domains (indicated by subscripts) pass through the CGRS (order indicated by superscripts). The detailed generation of associations is described in the next section. In the domain alignment module *C*, G_1 and G_2 are two adversarial generators for associations. They are used to generate the target association adversarial to the source’s association, and vice versa. The situation when the source associations works as the “real player” for the adversarial generation is shown in Figure 1¹. Here the adversarials of the corresponding target associations are $\tilde{\mathbf{X}}_t^{st}$, and $\tilde{\mathbf{X}}_t^{ts}$. The discriminators D_1, D_2 are used to distinguish associations of \mathbf{X}_s^{st} from $\tilde{\mathbf{X}}_t^{st}$, and \mathbf{X}_s^{ts} from $\tilde{\mathbf{X}}_t^{ts}$ respectively.

Finally, L_G and L_T in module *D* and *E* are the learning metrics for domain confusion and task classification. Mod-

¹The arrangement can be flexible, i.e. it also works if the target association is used as the real player.

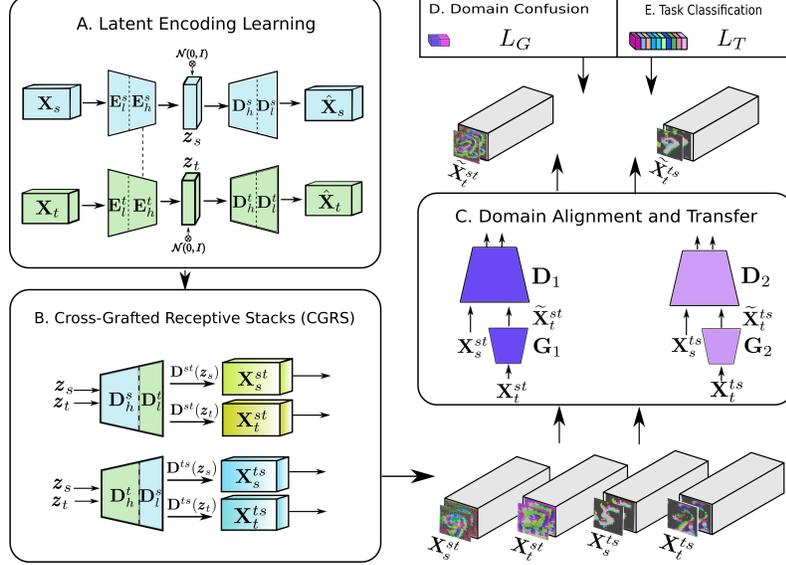


Figure 1: Overview of the the proposed model. There are 5 modules in it. In module A, the high-level layers of encoders E_h^s , E_h^t are shared (demonstrated by the dashed line). The outputs of D_h^s and D_h^t are the high-level representation of the source and target, whereas D_l^s , D_l^t are the low-level ones. The \mathbf{X}_s^{st} , \mathbf{X}_t^{st} , \mathbf{X}_s^{ts} , \mathbf{X}_t^{ts} in module B are the adversarial images reproduced by CGRS ($D^{st} \equiv [D_h^s \circ D_l^t]$ and $D^{ts} \equiv [D_h^t \circ D_l^s]$) from latent encodings. In module C, G_1 and G_2 are adversarial generators, D_1 , D_2 are discriminators. L_G and L_T are learning metric for the domain and task respectively. Best viewed in color.

ule C combines the learning metric modules to align the label space of the source and target images, and complete the adaptation. The training process adopts standard back-propagation. In contrast to the conventional domain adaptation framework in which the classifier input is $\{\mathbf{X}_s, \mathbf{y}_s\}$ and output is $\{\mathbf{X}_t, \hat{\mathbf{y}}_t\}$, our model’s classifier is trained by $\{\mathbf{X}_s^{st}, \mathbf{y}_s\}$, $\{\mathbf{X}_s^{ts}, \mathbf{y}_s\}$ and tested by $\{\tilde{\mathbf{X}}_t^{st}, \mathbf{y}_t\}$, $\{\tilde{\mathbf{X}}_t^{ts}, \mathbf{y}_t\}$. In short, the associations of the source data are used for training, and the adversarial generation of the target data are used in testing.

3.2. Generation of Cross-Grafted Associations

In module A, we obtain the latent encodings of source and target domains using VAEs [17], assuming they have a normal prior distribution. They encode a data sample \mathbf{x} to a latent space \mathbf{z} and decode the latent representation back to data space image $\hat{\mathbf{x}}$. We get all the latent encodings \mathbf{z}_s and \mathbf{z}_t , which are conceptually sampled from conditional probability densities $q(\mathbf{z}_s|\mathbf{X}_s)$ and $q(\mathbf{z}_t|\mathbf{X}_t)$ respectively. In module B, the cross-grafted receptive stacks are constructed to map the encodings to the cross-domain association spaces, which are later aligned to the source domain’s label space in module C.

Here CGRS recruits the high level (i) and low level (j) of the decoders of source (s) and target (t). It maps the latent space \mathbf{z}_k to the common association distributions \mathcal{P}_{ij} ,

which the associations \mathbf{X}_k^{ij} are sampled from:

$$D_{CGRS}(\mathbf{z}_k) \mapsto \mathbf{X}_k^{ij} \in \mathcal{P}_{ij}, \quad (1)$$

where $i, j, k \in \{s, t\}, i \neq j$. In detail, when the latent encoding \mathbf{z}_k passes through CGRS, the generation of associations can be expressed in a generative approach [4] as follows:

$$\mathcal{P}_i = (\prod_{l'=1}^N p_i(\mathbf{m}^{l'+1}|\mathbf{m}^{l'}, \theta_{D_{ih}}^{l'})) p_i(\mathbf{m}^1|\mathbf{z}_k, \theta_{D_{ih}}^1), \quad (2)$$

where N is the number of high-level decoder layers, $\theta_{D_{ih}}^{l'}$ is the map parameters of i in l' layer, $\mathbf{m}^{l'}$ is the output space of high-level decoder of layer l' . Then \mathbf{m}^N is transferred to final association space:

$$\mathcal{P}_{ij} = (\prod_{l''=1}^M p_j(\mathbf{n}^{l''+1}|\mathbf{n}^{l''}, \theta_{D_{jl}}^{l''})) p_j(\mathbf{n}^1|\mathbf{m}^N, \theta_{D_{jl}}^1), \quad (3)$$

where M is the number of low-level decoder layers, $\mathbf{n}^{l''}$ is the output space of low-level decoder of layer l'' , and $\theta_{D_{jl}}^{l''}$ is the map parameters of j in l'' layer. We assume the grafted parts p_j can be regarded as the corresponding reconstruction decoder injected with noise $\epsilon_{jl}^{l''}, \epsilon_{jl}^{l''} \in \mathcal{P}(\theta_{D_{jl}}^{l''}|\mathbf{X}_j)$ in a normal distribution (more details in the supplementary). This bridges the gap between the source and target domains, and also enhances the generalization of the model. Figure 2 gives the schematic illustration for the generation of associations \mathbf{X}_s^{st} and \mathbf{X}_t^{st} , when $i = s, j = t$. Another case is

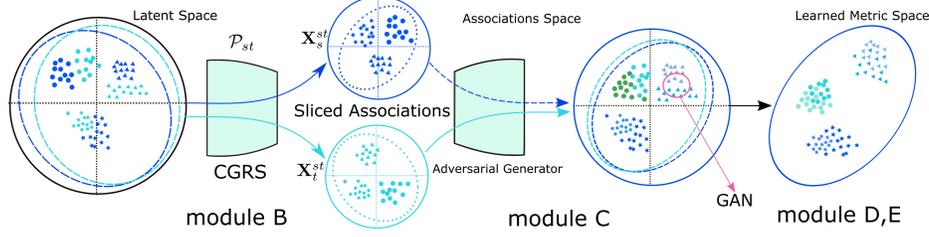


Figure 2: Generation of associations for channel \mathbf{X}^{st} . The encodings of source s and target t are formed in the latent space first. Then, they are projected to association spaces by CGRS. Finally, the latent and association spaces are aligned by the adversarial training combined with learning metrics. The adversarial process is flexible, and can be from the target to the source, and vice versa. The former scenario is shown here. The dashed line means the source associations are the real player in adversarial generation.

$i = t, j = s$, which corresponds to the associations \mathbf{X}_s^{ts} and \mathbf{X}_t^{ts} .

The associations are constructed, but they are yet to be aligned to the same label space of the source domain. To get the label distributions aligned, we use discriminator D to confuse the associations generated by different encodings. The discriminators make the distributions of associations more similar by minimizing the Jensen-Shannon divergence [11] (JSD):

$$\begin{aligned} \tilde{\mathbf{X}}_t^{ij} &\in p(\tilde{\mathbf{X}}_t^{ij} | \mathbf{X}_t^{ij}, \theta_D, \theta_{G_k}) \\ \text{w.r.t } \min JSD(p(\mathbf{X}_s^{ij}) || p(\tilde{\mathbf{X}}_t^{ij})), \end{aligned} \quad (4)$$

In the model, θ_{G_k} (G_1 when $i = s, j = t$ and G_2 for $i = t, j = s$) are used as generators for $\tilde{\mathbf{X}}_t^{st}$ and $\tilde{\mathbf{X}}_t^{ts}$ during the alignment, as shown in Figure 1. The encoders in module A and adversarial generators of module C are updated during training to minimum the Jensen-Shannon divergence of associations.

3.3. Learning

To train our model, we jointly solve the learning problems of the subnetworks. There are four loss functions, for the within-domain VAEs [17], cross-domain adversarial networks, content constancy and classifier training loss respectively.

First, we need to learn the representations of the source and target domains from encoders and decoders. Here, we minimize the within-domain VAEs loss functions. The loss function of our VAEs consists of both reconstruction error and prior regularization:

$$L_{VAEs} = L_{like}^{pixel} + L_{prior}. \quad (5)$$

The L_{like}^{pixel} and L_{prior} are given by

$$\begin{aligned} L_{like}^{pixel} &= -\lambda_1 \{ \mathbb{E}_{q_s(\mathbf{z}_s | \mathbf{x}_s)} [\log p_s(\mathbf{X}_s | \mathbf{z}_s)] \\ &\quad + \mathbb{E}_{q_t(\mathbf{z}_t | \mathbf{x}_t)} [\log p_t(\mathbf{X}_t | \mathbf{z}_t)] \}, \end{aligned} \quad (6)$$

$$\begin{aligned} L_{prior} &= \lambda_2 \{ D_{KL}(q_s(\mathbf{z}_s | x_s) || p(z)) \\ &\quad + D_{KL}(q_t(\mathbf{z}_t | x_t) || p(z)) \}, \end{aligned} \quad (7)$$

where D_{KL} is the Kullback-Leibler divergence. λ_1 and λ_2 are the trade-off hyper-parameters to control the priority of variational encoding and reconstruction.

To align the source and target domains, we use the adversarial training for the two association spaces \mathcal{P}_{st} and \mathcal{P}_{ts} . Their adversarial objectives L_G^{st} and L_G^{ts} are:

$$\begin{aligned} L_G^{st}(E_s, D^{st}, D_1) &= \lambda_0 \{ \mathbb{E}_{x_s} [\log D_1(D^{st}(\mathbf{z}_s))] \\ &\quad + \mathbb{E}_{x_s, z_s} [\log(1 - D_1(G_1(D^{st}(\mathbf{z}_t))))] \}, \end{aligned} \quad (8)$$

$$\begin{aligned} L_G^{ts}(E_t, D^{ts}, D_2) &= \lambda_0 \{ \mathbb{E}_{x_t} [\log D_2(D^{ts}(\mathbf{z}_s))] \\ &\quad + \mathbb{E}_{x_t, z_t} [\log(1 - D_2(G_2(D^{ts}(\mathbf{z}_t))))] \}, \end{aligned} \quad (9)$$

where $D^{st} \equiv D_h^s \circ D_l^t$ and $D^{ts} \equiv D_h^t \circ D_l^s$. $D(x)$ is the probability function assigned by the discriminator network, which tries to distinguish the generated source-based associations from the target-based ones. At last, the overall adversarial generative cost function is:

$$L_G = L_G^{st} + L_G^{ts}. \quad (10)$$

For the training stability, we introduce a content constancy loss function for the associations. Both the $L1$ and $L2$ penalty can be used to regularize the associations. Here we render a masked pairwise mean squared error [5]. Formally, when a binary mask \mathbf{m} is given ($\mathbf{m} \in \mathcal{R}^k$), the masked PMSE loss for associations \mathbf{X}^{st} and \mathbf{X}^{ts} is given as follows:

$$\begin{aligned} L_s^{st} &= \mathbb{E}_{\mathbf{X}_s^{st}, z} \left(\frac{1}{k} \| D^{st}(\mathbf{z}_s) - G_1(D^{st}(\mathbf{z}_t)) \circ \mathbf{m} \|_2^2 \right. \\ &\quad \left. - \frac{1}{k^2} ((D^{st}(\mathbf{z}_s) - G_1(D^{st}(\mathbf{z}_t)))^T \mathbf{m})^2 \right), \end{aligned} \quad (11)$$

and

$$\begin{aligned} L_s^{ts} &= \mathbb{E}_{\mathbf{X}_s^{ts}, z} \left(\frac{1}{k} \| D^{ts}(\mathbf{z}_s) - G_2(D^{ts}(\mathbf{z}_t)) \circ \mathbf{m} \|_2^2 \right. \\ &\quad \left. - \frac{1}{k^2} ((D^{ts}(\mathbf{z}_s) - G_2(D^{ts}(\mathbf{z}_t)))^T \mathbf{m})^2 \right). \end{aligned} \quad (12)$$

So the overall content objective for associations is:

$$L_s = \lambda_3 (L_s^{st} + L_s^{ts}). \quad (13)$$



Figure 3: Examples of the Datasets used for Experiments.

At last, for classification we use a typical soft-max cross-entropy loss:

$$L_T = \mathbb{E}[-y_s^T \log T(\mathbf{X}_s^{st}) - y_s^T \log T(\mathbf{X}_s^{ts})], \quad (14)$$

where y_s is the class label for source \mathbf{X}_s , and T is the task classifier. Finally, the overall loss function of our model is:

$$L^* = \min_{E, D, G} \max_{D_1, D_2} (L_{VAE_s} + L_G + L_s + L_T). \quad (15)$$

We solve this minimax problem of the loss function optimization by three alternating steps. First, the latent encodings are learned by the self-mapped process, which updates (E_s, E_t, D_s, D_t) , but keeps CGRS (D^{st}, D^{ts}) , (D_1, D_2) and (G_1, G_2) fixed. Then, we apply a gradient ascent step to update two discriminators D_1, D_2 and the classifier T , while keeping two VAEs channels (E_s, E_t, D_s, D_t) and CGRS (D^{st}, D^{ts}) , (G_1, G_2) fixed. Finally, a gradient descent step is applied to update (E_1, E_2, G_1, G_2) , while (D^{st}, D^{ts}) , D_1, D_2 and T are fixed.

4. Experiments and Results

We have evaluated our model on some benchmark datasets used commonly in the domain adaptation literature, including MNIST [19], MNIST-M [8], and USPS [18]. Also included is a multi-digit dataset ‘‘M-Digits’’, which we developed based on MNIST. The Fashion dataset [38] and its polluted version ‘‘Fashion-M’’ are also used in the experiments. Example images of these datasets are shown in Figure 3.

We compare our CGRS-LA method with the state-of-the-art domain adaptation methods: Pixel-level domain adaptation (PixelDA) [5], Domain Adversarial Neural Network (DANN) [8], Unsupervised Image-to-Image translation (UNIT) [21], Cycle-Consistent Adversarial Domain Adaption (CyCADA) [15], Generate to Adapt (GtA) [31] and Conditional Domain Adversarial Network (CDAN) [24]. In addition, we also use the source-only and target-only training as the lower and upper bound respectively, following the practice in [5, 8].

4.1. Datasets and Adaptation Scenarios

We use six popular datasets to construct four domain adaptation scenarios:

MNIST \rightleftharpoons MNIST-M: This is a scenario when the image content is the same, but the target data are polluted by noise. MNIST handwritten dataset [19] is a very popular machine learning dataset. It has a training set of 60,000 binary images, and a test set of 10,000. There are 10 classes in the dataset. MNIST-M [8] is a modified version for the MNIST, with random RGB background cropped from the Berkeley Segmentation Dataset². In our experiments, we use the standard split of the dataset.

MNIST \rightleftharpoons USPS: For this scenario, source and target domains have different contents but the same background. USPS is a handwritten zip digits datasets [18]. It is collected by the U.S Postal Service from envelopes processed at the Buffalo, N.Y Post Office. It contains 9298 binary images (16×16), 7291 of which are used as the training set, while the remaining 2007 are used as the test set. The USPS samples are resized to 28×28 , the same as MNIST.

Fashion \rightleftharpoons Fashion-M: Fashion-MNIST [38] contains 60,000 images for training, and 10,000 for testing. All the images are grayscale, 28×28 in size space. The samples are collected from 10 fashion categories: T-shirt/Top, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag and Ankle Boot. There are some complex textures in the images. In addition, following the protocol in [8], we add random noise to the Fashion images to generate the Fashion-M dataset.

MNIST \rightleftharpoons M-Digits In this scenario, we design a multi-digits dataset to evaluate the proposed model, noted as M-Digits. The MNIST digits are cropped first, and then are randomly selected, combined and randomly aligned in a new image, limited to 3 digits in maximum. The label for the new image is decided by the central digit. Finally, the new dataset is resized to 28×28 .

4.2. Implementation Details

All the models are implemented using the TensorFlow³ [1] and are trained with Mini-Batch Gradient Descent using the Adam optimizer [16]. The initial learning rate is 0.0002. Then it adopts an annealing method, with a decay of 0.95 after every 20,000 mini-batch steps. The mini-batch size for both the source and target domains are 64 samples, and the input images are rescaled to $[-1, 1]$. The hyper-parameters are $\lambda_0 = 1$, $\lambda_1 = 10$, $\lambda_2 = 0.01$, $\lambda_3 = 1$.

In our implementation, the latent space is sampled from a normal distribution $\mathcal{N}(0, I)$, and is achieved by the convolution encoders. The transpose convolution [40] is used in the decoder to build the reconstruction image space. This follows a similar structure protocol of [21], but we modify the padding strategy to ‘same’ for convolution layers. For sake of convenience in experiments, we add another 32-kernel layer before the last layer in the decoders. The stride is 2 for down-sampling in the encoders, and their counter-

²URL <https://www2.eecs.berkeley.edu/Research/Projects/CS/vision/bsds/>

³Our code will be available on github after the double blind review.

Table 1: Mean classification accuracy comparison. The "source only" row is the accuracy for target without domain adaptation training only on the source. And the "target only" is the accuracy of the full adaptation training on the target. For each source-target task the best performance is in bold.

Source Target	MNIST MNIST-M	MNIST-M MNIST	MNIST USPS	USPS MNIST	MNIST M-Digits	M-Digits MNIST	Fashion Fashion-M	Fashion-M Fashion
Source Only	0.561	0.633	0.634	0.625	0.603	0.651	0.527	0.612
CORAL [33]	0.817	-	0.577	-	-	-	-	-
MMD [23]	0.811	-	0.769	-	-	-	-	-
CyCADA [15]	-	-	0.956	0.965	-	-	-	-
GtA [31]	-	-	0.953	0.908	-	-	-	-
CDAV [24]	-	-	0.956	0.980	-	-	-	-
DANN [8]	0.766	0.851	0.774	0.833	0.864	0.920	0.604	0.822
PixelDA [5]	0.982	0.922	0.959	0.942	0.734	0.913	0.805	0.762
UNIT [21]	0.920	0.932	0.960	0.951	0.903	0.910	0.796	0.805
CGRS-LA (\mathbf{X}^{st})	0.821	0.935	0.946	0.938	0.895	0.902	0.735	0.805
CGRS-LA (\mathbf{X}^{ts})	0.923	0.840	0.902	0.930	0.853	0.851	0.792	0.760
CGRS-LA-C (\mathbf{X}^{st})	0.890	0.983	0.961	0.956	0.916	0.923	0.766	0.825
CGRS-LA-C (\mathbf{X}^{ts})	0.983	0.871	0.943	0.953	0.883	0.892	0.813	0.811
Target Only	0.983	0.985	0.980	0.985	0.982	0.985	0.920	0.942

part in decoders is also 2 so as to get the same dimensionality of the original image. The encoders for source and target domains share their high-level layers. We add the batch normalization between each layer in the encoders and the decoders. The CGRS of associations is the composition of different levels of the source and target’s representation. The stride step keeps 1 for all the dimensions in the adversarial generator, and the kernel is 3×3 . This adopts the structure of PixelDA [5], which uses a ResNet architecture. The discriminator fuses the domains, and also plays as a task classifier for the label space learning. It follows the design as in [21]. However, we do not share the layers of discriminators of \mathbf{X}^{st} and \mathbf{X}^{ts} channels. Also, we replace the max-pooling with a stride of 2×2 steps.

4.3. Results

4.3.1 Quantitative Results

Now we report the classification performance of our proposed model. During the experiments, associations \mathbf{X}_s^{st} and \mathbf{X}_t^{ts} are used to train the classifier, and the adversarial generation of \mathbf{X}_t^{st} and \mathbf{X}_s^{ts} are used for testing. The accuracy of the target domain classification after domain adaptation is listed in Table 1, presenting the result of 12 methods (4 versions of our model CGRS-LA, and 8 state-of-the-art methods) across 4 tasks (each in two directions). Our proposed model outperforms the state-of-the-art in most of the scenarios, especially when content constancy is considered. Also, it can be seen that the adaptation performance is usually asymmetric for the methods in comparison, e.g. the accuracies for MNIST→M-Digits and M-Digits→MNIST

are quite different for DANN and PixelDA. The CGRS-LA models, however, perform almost equally well on both directions for these adaptation tasks.

For MNIST⇌MNIST-M and MNIST⇌USPS, the mean classification accuracy nearly reaches the upper bound, suggesting these are easier tasks. On the other hand, we can see the adaptation task between Fashion and Fashion-M is more difficult than others. For this task, our method again not only achieves the best performance but also demonstrates balanced performance in two directions.

4.3.2 Qualitative Results

Since our model adopts a generative approach, we can have direct visual evaluation of the associations generated by the CGRS. The generative associations obtained by CGRS are shown in Figure 4, obtained after 100k mini-batch steps for the Fashion scenario and 50k for other three scenarios. The CGRS generate the associations with very similar appearance for the source and target domains. Then the GANs is employed to move them closer. During association generation, the CGRS eliminate the strong noise of MNIST-M and Fashion-M. Though there are more complex textures in the Fashion task, the proposed model still performs well to produce reasonable visualizations of the associations. The associations of the Fashion scenario seem to suffer some information loss, possibly due to the complex textures and strongly polluted images. However, they still look reasonable upon visualization. The MNIST→M-Digits scenario maintains the original content style, while the associations display some style variation in the MNIST→USPS sce-

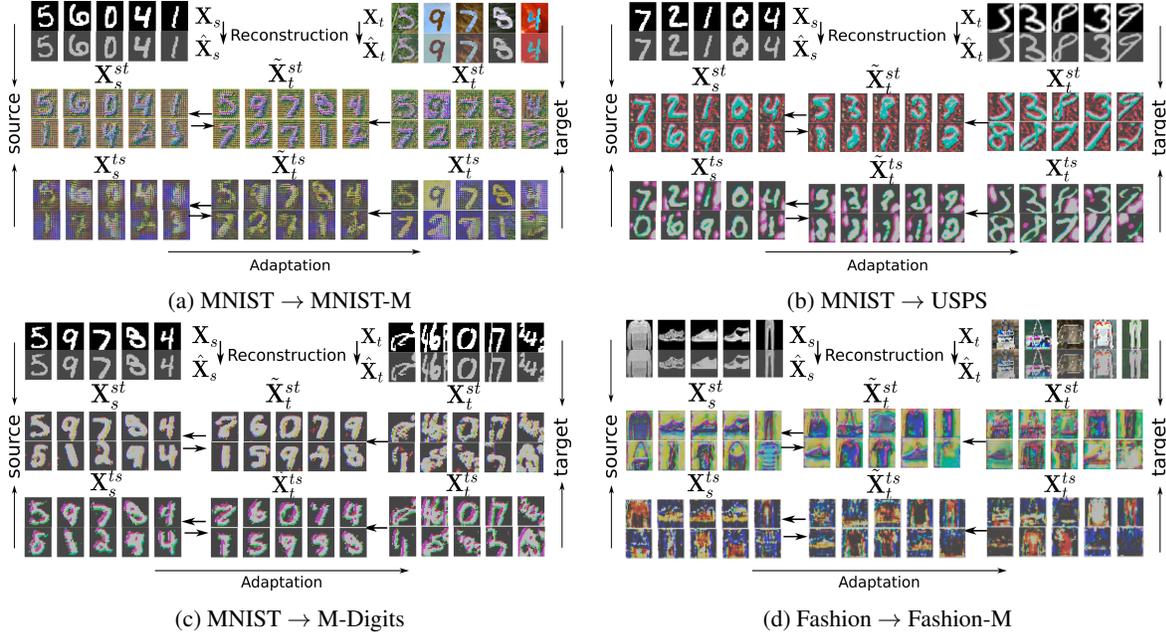


Figure 4: The visualization of association generations. For each scenario, the leftmost column is the source and its association, and the rightmost is for target. During the experiments, the associations of source are real player. The adversarial generations for target associations are in the middle column.

nario.

4.3.3 Model Analysis

Some further experiments are done to evaluate our model.

Ablation Study: We evaluate the potential effect of employing the content constancy strategy in our model. From the Table 1, we can see that the model with content constancy (denoted by CRGS-LA-C) outperforms its CRGS-LA counterpart. The constancy loss encourages the adversarial generation in a consistent way.

Sensitivity of CGRS: CGRS plays a critical role in the proposed model. In this section, we evaluate the performance of diverse structures of CGRS. During the experiments, we use a fix depth of network (6 layers) for the generation process. We apply various settings for splitting the high-level and low-level decoder stacks. For example, H5L1 denotes the scheme using 5 layers for high-level and 1 layer as low-level. The results of changing the CGRS setup for different scenarios are shown in Figure 5. It can be seen that for the channel X^{st} in MNIST \rightarrow MNIST-M and Fashion \rightarrow Fashion-M tasks, the highest accuracies are at the point H5L1, and for MNIST \rightarrow USPS and MNIST \rightarrow M-Digits tasks, there is a peak value at the point H2L4. The X^{ts} channel somehow seems more sensitive to varying CGRS settings.

Generalization of CGRS: Can we utilize the trained CGRS in one scenario to another adaptation task? In this evalua-

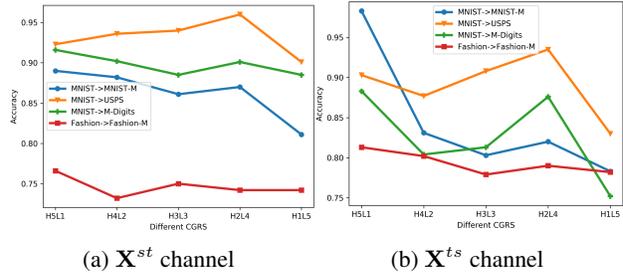


Figure 5: The Adaptation Accuracy of Different CGRS.

tion, we use our pre-trained CGRS from one scenario to adapt to a different task. These models are trained with a trade-off H4L2 CGRS according to the sensitivity analysis. During the experiments, we keep the CGRS fixed, then fine-tune the adversarial and label alignment parts. The results are shown in Table 2. Although there is a slight reduction to the accuracies reported before, the performance of adaptation to other tasks remains reasonable. Specifically, the CGRS of MNIST \rightarrow MNIST-M and Fashion \rightarrow Fashion-M adapts to other three scenarios pretty well, while the CGRS of the MNIST \rightarrow USPS and MNIST \rightarrow M-Digits get a lower accuracy for Fashion \rightarrow Fashion-M.

Visualization of Extracted Features: We also evaluate the features of top, fully connected layers in the discriminator for task MNIST \rightarrow USPS. The features are embedded by the t-SNE algorithm [27]. Figure 6 shows that the two domains

Table 2: Mean classification accuracy for Generalization Evaluation. The results of \mathbf{X}^{ts} channel is shown in the parentheses.

Source→Target	MNIST→MNIST-M	MNIST→USPS	MNIST→M-Digits	Fashion→Fashion-M
MNIST→MNIST-M	0.890(0.983)	0.958(0.945)	0.915(0.853)	0.762(0.730)
MNIST→USPS	0.915(0.859)	0.961(0.943)	0.882(0.914)	0.605(0.587)
MNIST→M-Digits	0.843(0.928)	0.944(0.958)	0.916(0.883)	0.613(0.593)
Fashion→Fashion-M	0.925(0.881)	0.932(0.935)	0.825(0.913)	0.766(0.813)

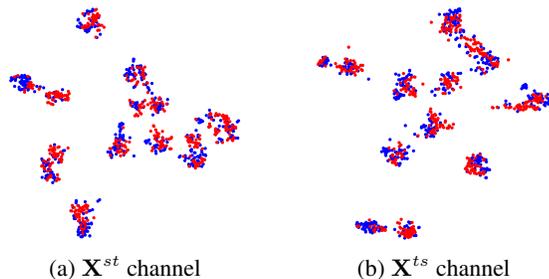


Figure 6: The visualization of top associations features embedded by t-SNE w.r.t source and target. The Blue dots are for source and red ones for target.

Table 3: Mean classification accuracy for semi-supervised evaluation.

Source Extra	MNIST MNIST-M	MNIST USPS	MNIST M-Digits	Fashion Fashion-M
1000	0.988	0.966	0.925	0.846
2000	0.990	0.970	0.932	0.855

can be aligned well on both channels after adaptation.

4.3.4 Semi-supervised Evaluation

Finally, we evaluate the performance of our model for semi-supervised learning. Under this scenario, it is assumed that we can get a small number of labeled target samples. Similar to the approach in [5], we choose 1000 samples from every category in the target domain as the baseline. These are added as extra to the source domain for training. The results are shown in Table 3. The adaptation performance is better when some target data are added into the source to train the classifier. It outperforms the unsupervised scenario when only 1000 target samples are fed to the classifier, whereas having 2000 target samples will further improve the performance.

4.3.5 Discussion

To sum up, our method can maintain stable performance when we vary the settings of CGRS for stack splitting. There seems to be a tendency to favour a higher ratio of high-level to low-level layers when the domains contain

similar contents but different background, while adaptation tasks with similar background but different content favour more low-level layers.

Another interesting observation is that CGRS have very good generalization ability. The CGRS trained by one task can be employed for domain adaptation in another task. This demonstrates a merit of our method for practical applications, that is the CGRS are transferable.

Finally, while the both association channels are well aligned, from our experiment it seems \mathbf{X}^{ts} claims better classification performance more often. In practical applications, it may be possible to design a classification combination method so that an optimal final decision can be developed from both association channels.

5. Conclusion

In this paper, we have proposed a novel unsupervised domain adaptation model based on cross-domain association generation, and label alignment using adversarial networks. In particular, cross-grafted representation stacks between different domains are constructed for bi-directional associations. The domain adaptation task hence transforms to constructing an effective mapping of the cross-domain associations onto the label space of the original source domain, a methodology we believe contributes to its robust performance in domain adaptation tasks. This is verified by the empirical results we have obtained from a number of tasks involving 6 benchmark tasks, which also demonstrate that the proposed CGRS also have strong cross-task generalization abilities. For future work, we would like to explore the extension of the framework for continual learning with cross-task adaptation.

References

- [1] M. Abadi et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *CoRR*, 2016. 1603.04467. 5
- [2] T. Adel and A. Wong. A probabilistic covariate shift assumption for domain adaptation. In *Proceeding of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 2476–2482, 2015. 1
- [3] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning from different domains. *Machine Learning*, 79(1-2):151–175, 2010. 1

- [4] Y. Bengio, A. C. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 35(8):1798–1828, 2013. [3](#)
- [5] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 95–104, 2017. [2](#), [4](#), [5](#), [6](#), [8](#)
- [6] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 343–351, 2016. [2](#)
- [7] S. Chopra, S. Balakrishnan, and R. Gopalan. Dlid: Deep learning for domain adaptation by interpolating between domains. *ICML Workshop on Challenges in Representation Learning (WREPL)*, 2013. [2](#)
- [8] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research (JMLR)*, 17(1):2096–2030, 2016. [2](#), [5](#), [6](#)
- [9] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, and W. Li. Deep reconstruction-classification networks for unsupervised domain adaptation. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 597–613. Springer, 2016. [2](#)
- [10] B. Gong, K. Grauman, and F. Sha. Geodesic flow kernel and landmarks: Kernel methods for unsupervised domain adaptation. In G. Csurka, editor, *Domain Adaptation in Computer Vision Applications.*, Advances in Computer Vision and Pattern Recognition, pages 59–79. Springer, 2017. [2](#)
- [11] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems (NeurIPS)*, pages 2672–2680, 2014. [1](#), [4](#)
- [12] R. Gopalan, R. Li, and R. Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *Proceeding of the IEEE International Conference on Computer Vision (ICCV)*, pages 999–1006, 2011. [1](#)
- [13] R. Gopalan, R. Li, and R. Chellappa. Unsupervised adaptation across domain shifts by generating intermediate data representations. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 36(11):2288–2302, 2014. [2](#)
- [14] S. Herath, M. T. Harandi, and F. Porikli. Learning an invariant hilbert space for domain adaptation. In *Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3956–3965, 2017. [1](#)
- [15] J. Hoffman, E. Tzeng, T. Park, J. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In J. G. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning (ICML)*, volume 80, pages 1994–2003, 2018. [5](#), [6](#)
- [16] D. P. Kingma and J. L. Ba. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. [5](#)
- [17] D. P. Kingma and M. Welling. Auto-encoding Variational-Bayes. *arXiv preprint arXiv:1312.6114*, 2013. [1](#), [3](#), [4](#)
- [18] Y. Le Cun, L. Jackel, B. Boser, J. Denker, H. Graf, I. Guyon, D. Henderson, R. Howard, and W. Hubbard. Handwritten digit recognition: Applications of neural network chips and automatic learning. *IEEE Communications Magazine*, 27(11):41–46, 1989. [5](#)
- [19] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. [5](#)
- [20] M. Liu and O. Tuzel. Coupled generative adversarial networks. In D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems (NeurIPS)*, pages 469–477, 2016. [2](#)
- [21] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 700–708, 2017. [1](#), [2](#), [5](#), [6](#)
- [22] Y.-C. Liu, Y.-Y. Yeh, T.-C. Fu, S.-D. Wang, W.-C. Chiu, and Y.-C. F. Wang. Detach and adapt: Learning cross-domain disentangled deep representation. *arXiv preprint arXiv:1705.01314*, 2017. [2](#)
- [23] M. Long, Y. Cao, J. Wang, and M. I. Jordan. Learning transferable features with deep adaptation networks. *arXiv preprint arXiv:1502.02791*, 2015. [2](#), [6](#)
- [24] M. Long, Z. Cao, J. Wang, and M. I. Jordan. Conditional adversarial domain adaptation. In S. Bengio, H. M. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31 (NeurIPS)*, pages 1647–1657, 2018. [5](#), [6](#)
- [25] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu. Transfer joint matching for unsupervised domain adaptation. In *Proceeding of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 1410–1417, 2014. [1](#)
- [26] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Unsupervised domain adaptation with residual transfer networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 136–144, 2016. [2](#)
- [27] L. v. d. Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of machine learning research (JMLR)*, 9:2579–2605, Nov 2008. [7](#)
- [28] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering (TKDE)*, 22(10):1345–1359, 2010. [1](#)
- [29] R. Raina, A. Battle, H. Lee, B. Packer, and A. Y. Ng. Self-taught learning: transfer learning from unlabeled data. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 227, pages 759–766, 2007. [2](#)
- [30] A. Rozantsev, M. Salzmann, and P. Fua. Beyond sharing weights for deep domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018. [2](#)
- [31] S. Sankaranarayanan, Y. Balaji, C. D. Castillo, and R. Chellappa. Generate to adapt: Aligning domains using generative adversarial networks. In *2018 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8503–8512, 2018. [5](#), [6](#)

- [32] O. Sener, H. O. Song, A. Saxena, and S. Savarese. Learning transferrable representations for unsupervised domain adaptation. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2110–2118, 2016. 1
- [33] B. Sun, J. Feng, and K. Saenko. Return of frustratingly easy domain adaptation. In *Proceeding of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 2058–2065, 2016. 6
- [34] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko. Simultaneous deep transfer across domains and tasks. In *Proceeding of the IEEE International Conference on Computer Vision, (ICCV)*, pages 4068–4076, 2015. 2
- [35] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2962–2971, 2017. 2
- [36] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014. 2
- [37] 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 (CVPR)*, pages 5018–5027, 2017. 2
- [38] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017. 5
- [39] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems (NeurIPS)*, pages 3320–3328, 2014. 1
- [40] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus. Deconvolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, pages 2528–2535, 2010. 5