

Deep Reconstruction-Classification Networks for Unsupervised Domain Adaptation

Muhammad Ghifary^{1*}, W. Bastiaan Kleijn¹, Mengjie Zhang¹, David Balduzzi¹, Wen Li²

¹Victoria University of Wellington, ²ETH Zürich

Abstract. In this paper, we propose a novel unsupervised domain adaptation algorithm based on deep learning for visual object recognition. Specifically, we design a new model called Deep Reconstruction-Classification Network (DRCN), which jointly learns a shared encoding representation for two tasks: i) supervised classification of labeled source data, and ii) unsupervised reconstruction of unlabeled target data. In this way, the learnt representation not only preserves discriminability, but also encodes useful information from the target domain. Our new DRCN model can be optimized by using backpropagation similarly as the standard neural networks.

We evaluate the performance of DRCN on a series of cross-domain object recognition tasks, where DRCN provides a considerable improvement (up to $\sim 8\%$ in accuracy) over the prior state-of-the-art algorithms. Interestingly, we also observe that the reconstruction pipeline of DRCN transforms images from the source domain into images whose appearance resembles the target dataset. This suggests that DRCN's performance is due to constructing a single composite representation that encodes information about both the structure of target images and the classification of source images. Finally, we provide a formal analysis to justify the algorithm's objective in domain adaptation context.

Keywords: domain adaptation, object recognition, deep learning, convolutional networks, transfer learning

1 Introduction

An important task in visual object recognition is to design algorithms that are robust to *dataset bias* [1]. Dataset bias arises when labeled training instances are available from a source domain and test instances are sampled from a related, but different, target domain. For example, consider a person identification application in *unmanned aerial vehicles* (UAV), which is essential for a variety of tasks, such as surveillance, people search, and remote monitoring [2]. One of the critical tasks is to identify people from a bird's-eye view; however collecting labeled data from that viewpoint can be very challenging. It is more desirable that a UAV can be trained on some already available *on-the-ground* labeled

* Current affiliation: Weta Digital

images (source), e.g., people photographs from social media, and then successfully applied to the actual UAV view (target). Traditional supervised learning algorithms typically perform poorly in this setting, since they assume that the training and test data are drawn from the same domain.

Domain adaptation attempts to deal with dataset bias using unlabeled data from the target domain so that the task of manual labeling the target data can be reduced. Unlabeled target data provides auxiliary training information that should help algorithms generalize better on the target domain than using source data only. Successful domain adaptation algorithms have large practical value, since acquiring a huge amount of labels from the target domain is often expensive or impossible. Although domain adaptation has gained increasing attention in object recognition, see [3] for a recent overview, the problem remains essentially unsolved since model accuracy has yet to reach a level that is satisfactory for real-world applications. Another issue is that many existing algorithms require optimization procedures that do not scale well as the size of datasets increases [4,5,6,7,8,9,10]. Earlier algorithms were typically designed for relatively small datasets, e.g., the Office dataset [11].

We consider a solution based on learning representations or features from raw data. Ideally, the learned feature should model the label distribution as well as reduce the discrepancy between the source and target domains. We hypothesize that a possible way to approximate such a feature is by (supervised) learning the *source label* distribution and (unsupervised) learning of the *target data distribution*. This is in the same spirit as *multi-task learning* in that learning auxiliary tasks can help the main task be learned better [12,13]. The goal of this paper is to develop an accurate, scalable multi-task feature learning algorithm in the context of domain adaptation.

Contribution: To achieve the goal stated above, we propose a new deep learning model for unsupervised domain adaptation. Deep learning algorithms are highly scalable since they run in linear time, can handle streaming data, and can be parallelized on GPUs. Indeed, deep learning has come to dominate object recognition in recent years [14,15].

We propose *Deep Reconstruction-Classification Network* (DRCN), a convolutional network that jointly learns two tasks: i) supervised source label prediction and ii) unsupervised target data reconstruction. The encoding parameters of the DRCN are shared across both tasks, while the decoding parameters are separated. The aim is that the learned label prediction function can perform well on classifying images in the target domain – the data reconstruction can thus be viewed as an auxiliary task to support the adaptation of the label prediction. Learning in DRCN alternates between unsupervised and supervised training, which is different from the standard *pretraining-finetuning* strategy [16,17].

From experiments over a variety of cross-domain object recognition tasks, DRCN performs better than the state-of-the-art domain adaptation algorithm [18], with up to $\sim 8\%$ accuracy gap. The DRCN learning strategy also provides a considerable improvement over the pretraining-finetuning strategy, indicating that it is more suitable for the unsupervised domain adaptation setting. We

furthermore perform a visual analysis by reconstructing source images through the learned reconstruction function. It is found that *the reconstructed outputs resemble the appearances of the target images* suggesting that the encoding representations are successfully adapted. Finally, we present a probabilistic analysis to show the relationship between the DRCN’s learning objective and a semi-supervised learning framework [19], and also the soundness of considering only data from a target domain for the data reconstruction training.

2 Related Work

Domain adaptation is a large field of research, with related work under several names such as class imbalance [20], covariate shift [21], and sample selection bias [22]. In [23], it is considered as a special case of transfer learning. Earlier work on domain adaptation focused on text document analysis and NLP [24,25]. In recent years, it has gained a lot of attention in the computer vision community, mainly for object recognition application, see [3] and references therein. The domain adaptation problem is often referred to as *dataset bias* in computer vision [1].

This paper is concerned with *unsupervised domain adaptation* in which labeled data from the target domain is not available [26]. A range of approaches along this line of research in object recognition have been proposed [4,5,27,28,29,30,9], most were designed specifically for small datasets such as the Office dataset [11]. Furthermore, they usually operated on the SURF-based features [31] extracted from the raw pixels. In essence, the unsupervised domain adaptation problem remains open and needs more powerful solutions that are useful for practical situations.

Deep learning now plays a major role in the advancement of domain adaptation. An early attempt addressed large-scale sentiment classification [32], where the concatenated features from fully connected layers of stacked denoising autoencoders have been found to be domain-adaptive [33]. In visual recognition, a fully connected, shallow network pretrained by denoising autoencoders has shown a certain level of effectiveness [34]. It is widely known that deep convolutional networks (ConvNets) [35] are a more natural choice for visual recognition tasks and have achieved significant successes [36,14,15]. More recently, ConvNets pretrained on a large-scale dataset, ImageNet, have been shown to be reasonably effective for domain adaptation [14]. They provide significantly better performances than the SURF-based features on the Office dataset [37,38]. An earlier approach on using a convolutional architecture without pretraining on ImageNet, DLID, has also been explored [39] and performs better than the SURF-based features.

To further improve the domain adaptation performance, the pretrained ConvNets can be *fine-tuned* under a particular constraint related to minimizing a domain discrepancy measure [18,40,41,42]. Deep Domain Confusion (DDC) [41] utilizes the maximum mean discrepancy (MMD) measure [43] as an additional loss function for the fine-tuning to adapt the last fully connected layer. Deep Adaptation Network (DAN) [40] fine-tunes not only the last fully connected

layer, but also some convolutional and fully connected layers underneath, and outperforms DDC. Recently, the deep model proposed in [42] extends the idea of DDC by adding a criterion to guarantee the class alignment between different domains. However, it is limited only to the *semi-supervised* adaptation setting, where a small number of target labels can be acquired.

The algorithm proposed in [18], which we refer to as ReverseGrad, handles the domain invariance as a binary classification problem. It thus optimizes two contradictory objectives: i) minimizing label prediction loss and ii) maximizing domain classification loss via a simple *gradient reversal* strategy. ReverseGrad can be effectively applied both in the pretrained and randomly initialized deep networks. The randomly initialized model is also shown to perform well on cross-domain recognition tasks other than the Office benchmark, i.e., large-scale handwritten digit recognition tasks. Our work in this paper is in a similar spirit to ReverseGrad in that it does not necessarily require pretrained deep networks to perform well on some tasks. However, our proposed method undertakes a fundamentally different learning algorithm: finding a good label classifier while simultaneously learning the structure of the target images.

3 Deep Reconstruction-Classification Networks

This section describes our proposed deep learning algorithm for unsupervised domain adaptation, which we refer to as *Deep Reconstruction-Classification Networks* (DRCN). We first briefly discuss the unsupervised domain adaptation problem. We then present the DRCN architecture, learning algorithm, and other useful aspects.

Let us define a *domain* as a probability distribution \mathbb{D}_{XY} (or just \mathbb{D}) on $\mathcal{X} \times \mathcal{Y}$, where \mathcal{X} is the input space and \mathcal{Y} is the output space. Denote the source domain by \mathbb{P} and the target domain by \mathbb{Q} , where $\mathbb{P} \neq \mathbb{Q}$. The aim in *unsupervised domain adaptation* is as follows: given a labeled i.i.d. sample from a source domain $S^s = \{(x_i^s, y_i^s)\}_{i=1}^{n_s} \sim \mathbb{P}$ and an unlabeled sample from a target domain $S_u^t = \{(x_i^t)\}_{i=1}^{n_t} \sim \mathbb{Q}_X$, find a good labeling function $f : \mathcal{X} \rightarrow \mathcal{Y}$ on S_u^t . We consider a feature learning approach: finding a function $g : \mathcal{X} \rightarrow \mathcal{F}$ such that the discrepancy between distribution \mathbb{P} and \mathbb{Q} is minimized in \mathcal{F} .

Ideally, a discriminative representation should model both the label and the structure of the data. Based on that intuition, we hypothesize that a domain-adaptive representation should satisfy two criteria: i) classify well the source domain labeled data and ii) reconstruct well the target domain unlabeled data, which can be viewed as an approximate of the ideal discriminative representation. Our model is based on a convolutional architecture that has two pipelines with a shared encoding representation. The first pipeline is a standard convolutional network for *source label prediction* [35], while the second one is a convolutional autoencoder for *target data reconstruction* [44,45]. Convolutional architectures are a natural choice for object recognition to capture spatial correlation of images. The model is optimized through multitask learning [12], that is, jointly learns the (supervised) source label prediction and the (unsupervised) target

data reconstruction tasks.¹ The aim is that the encoding shared representation should learn the commonality between those tasks that provides useful information for cross-domain object recognition. Figure 1 illustrates the architecture of DRCN.

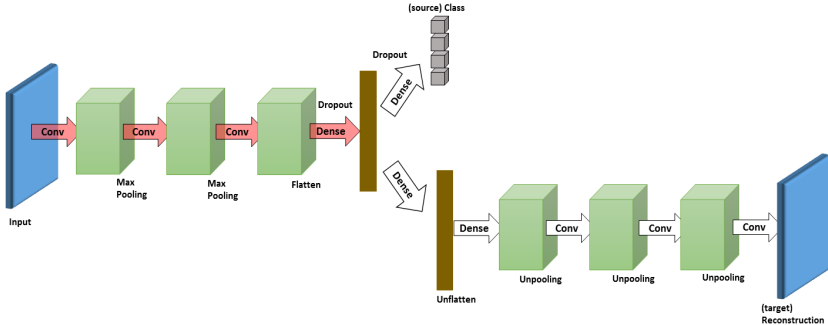


Fig. 1. Illustration of the DRCN’s architecture. It consists of two pipelines: i) label prediction and ii) data reconstruction pipelines. The shared parameters between those two pipelines are indicated by the red color.

We now describe DRCN more formally. Let $f_c : \mathcal{X} \rightarrow \mathcal{Y}$ be the (supervised) label prediction pipeline and $f_r : \mathcal{X} \rightarrow \mathcal{X}$ be the (unsupervised) data reconstruction pipeline of DRCN. Define three additional functions: 1) an encoder / feature mapping $g_{\text{enc}} : \mathcal{X} \rightarrow \mathcal{F}$, 2) a decoder $g_{\text{dec}} : \mathcal{F} \rightarrow \mathcal{X}$, and 3) a feature labeling $g_{\text{lab}} : \mathcal{F} \rightarrow \mathcal{Y}$. For m -class classification problems, the output of g_{lab} usually forms an m -dimensional vector of real values in the range $[0, 1]$ that add up to 1, i.e., *softmax* output. Given an input $x \in \mathcal{X}$, one can decompose f_c and f_r such that

$$f_c(x) = (g_{\text{lab}} \circ g_{\text{enc}})(x), \quad (1)$$

$$f_r(x) = (g_{\text{dec}} \circ g_{\text{enc}})(x). \quad (2)$$

Let $\Theta_c = \{\Theta_{\text{enc}}, \Theta_{\text{lab}}\}$ and $\Theta_r = \{\Theta_{\text{enc}}, \Theta_{\text{dec}}\}$ denote the parameters of the supervised and unsupervised model. Θ_{enc} are shared parameters for the feature mapping g_{enc} . Note that $\Theta_{\text{enc}}, \Theta_{\text{dec}}, \Theta_{\text{lab}}$ may encode parameters of multiple layers. The goal is to seek a single feature mapping g_{enc} model that supports both f_c and f_r .

Learning algorithm: The learning objective is as follows. Suppose the inputs lie in $\mathcal{X} \subseteq \mathbb{R}^d$ and their labels lie in $\mathcal{Y} \subseteq \mathbb{R}^m$. Let $\ell_c : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ and $\ell_r : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be the classification and reconstruction loss respectively. Given

¹ The unsupervised convolutional autoencoder is not trained via the greedy layer-wise fashion, but only with the standard back-propagation over the whole pipeline.

labeled source sample $S^s = \{(\mathbf{x}_i^s, \mathbf{y}_i^s)\}_{i=1}^{n_s} \sim \mathbb{P}$, where $\mathbf{y}_i \in \{0, 1\}^m$ is a *one-hot* vector, and unlabeled target sample $S_u^t = \{(\mathbf{x}_j^t)\}_{j=1}^{n_t} \sim \mathbb{Q}$, we define the empirical losses as:

$$\mathcal{L}_c^{n_s}(\{\theta_{\text{enc}}, \theta_{\text{lab}}\}) := \sum_{i=1}^{n_s} \ell_c(f_c(\mathbf{x}_i^s; \{\theta_{\text{enc}}, \theta_{\text{lab}}\}), \mathbf{y}_i^s), \quad (3)$$

$$\mathcal{L}_r^{n_t}(\{\theta_{\text{enc}}, \theta_{\text{dec}}\}) := \sum_{j=1}^{n_t} \ell_r(f_r(\mathbf{x}_j^t; \{\theta_{\text{enc}}, \theta_{\text{dec}}\}), \mathbf{x}_j^t). \quad (4)$$

Typically, ℓ_c is of the form *cross-entropy loss* $\sum_{k=1}^m y_k \log[f_c(\mathbf{x})]_k$ (recall that $f_c(\mathbf{x})$ is the softmax output) and ℓ_r is of the form *squared loss* $\|\mathbf{x} - f_r(\mathbf{x})\|_2^2$.

Our aim is to solve the following objective:

$$\min \lambda \mathcal{L}_c^{n_s}(\{\theta_{\text{enc}}, \theta_{\text{lab}}\}) + (1 - \lambda) \mathcal{L}_r^{n_t}(\{\theta_{\text{enc}}, \theta_{\text{dec}}\}), \quad (5)$$

where $0 \leq \lambda \leq 1$ is a hyper-parameter controlling the trade-off between classification and reconstruction. The objective is a convex combination of supervised and unsupervised loss functions. We justify the approach in Section 5.

Objective (5) can be achieved by alternately minimizing $\mathcal{L}_c^{n_s}$ and $\mathcal{L}_r^{n_t}$ using *stochastic gradient descent* (SGD). In the implementation, we used RM-Sprop [46], the variant of SGD with a gradient normalization – the current gradient is divided by a moving average over the previous root mean squared gradients. We utilize dropout regularization [47] during $\mathcal{L}_c^{n_s}$ minimization, which is effective to reduce overfitting. Note that dropout regularization is applied in the fully-connected/dense layers only, see Figure 1.

The stopping criterion for the algorithm is determined by monitoring the average reconstruction loss of the unsupervised model during training – the process is stopped when the average reconstruction loss stabilizes. Once the training is completed, the optimal parameters $\hat{\theta}_{\text{enc}}$ and $\hat{\theta}_{\text{lab}}$ are used to form a classification model $f_c(\mathbf{x}^t; \{\hat{\theta}_{\text{enc}}, \hat{\theta}_{\text{lab}}\})$ that is expected to perform well on the target domain. The DRCN learning algorithm is summarized in Algorithm 1 and implemented using Theano [48].

Data augmentation and denoising: We use two well-known strategies to improve DRCN’s performance: data augmentation and denoising. Data augmentation generates additional training data during the supervised training with respect to some plausible transformations over the original data, which improves generalization, see e.g. [49]. Denoising involves reconstructing *clean* inputs given their *noisy* counterparts. It is used to improve the feature invariance of denoising autoencoders (DAE) [33]. Generalization and feature invariance are two properties needed to improve domain adaptation. Since DRCN has both classification and reconstruction aspects, we can naturally apply these two tricks simultaneously in the training stage.

Let $\mathbb{Q}_{\tilde{X}|X}$ denote the noise distribution given the original data from which the noisy data are sampled from. The classification pipeline of DRCN f_c thus

Algorithm 1 The Deep Reconstruction-Classification Network (DRCN) learning algorithm.

Input:

- Labeled source data: $S^s = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^{n_s}$;
- Unlabeled target data: $S_u^t = \{\mathbf{x}_j^t\}_{i=j}^{n_t}$;
- Learning rates: α_c and α_r ;

- 1: Initialize parameters $\Theta_{\text{enc}}, \Theta_{\text{dec}}, \Theta_{\text{lab}}$
- 2: **while** not stop **do**
- 3: **for each** source batch of size m_s **do**
- 4: Do a forward pass according to (1);
- 5: Let $\Theta_c = \{\Theta_{\text{enc}}, \Theta_{\text{lab}}\}$. Update Θ_c :

$$\Theta_c \leftarrow \Theta_c - \alpha_c \lambda \nabla_{\Theta_c} \mathcal{L}^{m_s}(\Theta_c);$$

- 6: **end for**
- 7: **for each** target batch of size m_t **do**
- 8: Do a forward pass according to (2);
- 9: Let $\Theta_r = \{\Theta_{\text{enc}}, \Theta_{\text{dec}}\}$. Update Θ_r :

$$\Theta_r \leftarrow \Theta_r - \alpha_r (1 - \lambda) \nabla_{\Theta_r} \mathcal{L}^{m_t}(\Theta_r).$$

- 10: **end for**
- 11: **end while**

Output:

- DRCN learnt parameters: $\hat{\Theta} = \{\hat{\Theta}_{\text{enc}}, \hat{\Theta}_{\text{dec}}, \hat{\Theta}_{\text{lab}}\}$;
-

actually observes additional pairs $\{(\tilde{\mathbf{x}}_i^s, y_i^s)\}_{i=1}^{n_s}$ and the reconstruction pipeline f_r observes $\{(\tilde{\mathbf{x}}_i^t, \mathbf{x}_i^t)\}_{i=1}^{n_t}$. The noise distribution $\mathbb{Q}_{\tilde{\mathbf{X}}|X}$ are typically geometric transformations (translation, rotation, skewing, and scaling) in data augmentation, while either zero-masked noise or Gaussian noise is used in the denoising strategy. In this work, we combine all the fore-mentioned types of noise for denoising and use only the geometric transformations for data augmentation.

4 Experiments and Results

This section reports the evaluation results of DRCN. It is divided into two parts. The first part focuses on the evaluation on large-scale datasets popular with deep learning methods, while the second part summarizes the results on the Office dataset [11].

4.1 Experiment I: SVHN, MNIST, USPS, CIFAR, and STL

The first set of experiments investigates the empirical performance of DRCN on five widely used benchmarks: MNIST [35], USPS [50], Street View House Numbers (SVHN) [51], CIFAR [52], and STL [53], see the corresponding references for more detailed configurations. The task is to perform cross-domain recognition:

taking the training set from one dataset as the source domain and the test set from another dataset as the target domain. We evaluate our algorithm’s recognition accuracy over three cross-domain pairs: 1) MNIST vs USPS, 2) SVHN vs MNIST, and 3) CIFAR vs STL.

MNIST (MN) vs USPS (US) contains 2D grayscale handwritten digit images of 10 classes. We preprocessed them as follows. USPS images were rescaled into 28×28 and pixels were normalized to $[0, 1]$ values. From this pair, two cross-domain recognition tasks were performed: $MN \rightarrow US$ and $US \rightarrow MN$.

In SVHN (SV) vs MNIST (MN) pair, MNIST images were rescaled to 32×32 and SVHN images were grayscaled. The $[0, 1]$ normalization was then applied to all images. Note that we did not preprocess SVHN images using local contrast normalization as in [54]. We evaluated our algorithm on $SV \rightarrow MN$ and $MN \rightarrow SV$ cross-domain recognition tasks.

STL (ST) vs CIFAR (CI) consists of RGB images that share eight object classes: *airplane*, *bird*, *cat*, *deer*, *dog*, *horse*, *ship*, and *truck*, which forms 4,000 (train) and 6,400 (test) images for STL, and 40,000 (train) and 8,000 (test) images for CIFAR. STL images were rescaled to 32×32 and pixels were standardized into zero-mean and unit-variance. Our algorithm was evaluated on two cross-domain tasks, that is, $ST \rightarrow CI$ and $CI \rightarrow ST$.

The architecture and learning setup: The DRCN architecture used in the experiments is adopted from [44]. The label prediction pipeline has three convolutional layers: 100 5×5 filters (CONV1), 150 5×5 filters (CONV2), and 200 3×3 filters (CONV3) respectively, two max-pooling layers of size 2×2 after the first and the second convolutional layers (POOL1 and POOL2), and three fully-connected layers (FC4, FC5, and FC_OUT) – FC_OUT is the output layer. The number of neurons in FC4 or FC5 was treated as a tunable hyper-parameter in the range of $[300, 350, \dots, 1000]$, chosen according to the best performance on the validation set. The shared encoder g_{enc} has thus a configuration of CONV1-POOL1-CONV2-POOL2-CONV3-FC4-FC5. Furthermore, the configuration of the decoder g_{dec} is the inverse of that of g_{enc} . Note that the unpooling operation in g_{dec} performs by upsampling-by-duplication: inserting the pooled values in the appropriate locations in the feature maps, with the remaining elements being the same as the pooled values.

We employ ReLU activations [55] in all hidden layers and linear activations in the output layer of the reconstruction pipeline. Updates in both classification and reconstruction tasks were computed via RMSprop with learning rate of 10^{-4} and moving average decay of 0.9. The control penalty λ was selected according to accuracy on the source validation data – typically, the optimal value was in the range $[0.4, 0.7]$.

Benchmark algorithms: We compare DRCN with the following methods. 1) ConvNet_{src}: a supervised convolutional network trained on the labeled source domain only, with the same network configuration as that of DRCN’s label prediction pipeline, 2) SCAE: ConvNet preceded by the layer-wise pretraining of stacked convolutional autoencoders on all unlabeled data [44], 3) SCAE_t:

similar to SCAE, but only unlabeled data from the target domain are used during pretraining, 4) SDA_{sh} [32]: the deep network with three fully connected layers, which is a successful domain adaptation model for sentiment classification, 5) Subspace Alignment (SA) [27],² and 6) ReverseGrad [18]: a recently published domain adaptation model based on deep convolutional networks that provides the state-of-the-art performance.

All deep learning based models above have the same architecture as DRCN for the label predictor. For ReverseGrad, we also evaluated the “original architecture” devised in [18] and chose whichever performed better of the original architecture or our architecture. Finally, we applied the data augmentation to all models similarly to DRCN. The ground-truth model is also evaluated, that is, a convolutional network trained from and tested on images from the target domain only ($ConvNet_{tgt}$), to measure the difference between the cross-domain performance and the ideal performance.

Classification accuracy: Table 1 summarizes the cross-domain recognition accuracy (*mean \pm std*) of all algorithms over ten independent runs. DRCN performs best in all but one cross-domain tasks, better than the prior state-of-the-art ReverseGrad. Notably on the $SV \rightarrow MN$ task, DRCN outperforms ReverseGrad with $\sim 8\%$ accuracy gap. DRCN also provides a considerable improvement over ReverseGrad ($\sim 5\%$) on the reverse task, $MN \rightarrow SV$, but the gap to the groundtruth is still large – this case was also mentioned in previous work as a failed case [18]. In the case of $CI \rightarrow ST$, the performance of DRCN almost matches the performance of the target baseline.

DRCN also convincingly outperforms the greedy-layer pretraining-based algorithms (SDA_{sh} , SCAE, and $SCAE_t$). This indicates the effectiveness of the simultaneous reconstruction-classification training strategy over the standard pretraining-finetuning in the context of domain adaptation.

Table 1. Accuracy (*mean \pm std* %) on five cross-domain recognition tasks over ten independent runs. Bold and underline indicate the best and second best domain adaptation performance. $ConvNet_{tgt}$ denotes the ground-truth model: training and testing on the target domain only.

Methods	MN \rightarrow US	US \rightarrow MN	SV \rightarrow MN	MN \rightarrow SV	ST \rightarrow CI	CI \rightarrow ST
$ConvNet_{src}$	85.55 \pm 0.12	65.77 \pm 0.06	62.33 \pm 0.09	25.95 \pm 0.04	54.17 \pm 0.21	63.61 \pm 0.17
SDA_{sh} [32]	43.14 \pm 0.16	37.30 \pm 0.12	55.15 \pm 0.08	8.23 \pm 0.11	35.82 \pm 0.07	42.27 \pm 0.12
SA [27]	85.89 \pm 0.13	51.54 \pm 0.06	63.17 \pm 0.07	28.52 \pm 0.10	54.04 \pm 0.19	62.88 \pm 0.15
SCAE [44]	85.78 \pm 0.08	63.11 \pm 0.04	60.02 \pm 0.16	27.12 \pm 0.08	54.25 \pm 0.13	62.18 \pm 0.04
$SCAE_t$ [44]	86.24 \pm 0.11	65.37 \pm 0.03	65.57 \pm 0.09	27.57 \pm 0.13	54.68 \pm 0.08	61.94 \pm 0.06
ReverseGrad [18]	91.11 \pm 0.07	74.01 \pm 0.05	73.91 \pm 0.07	35.67 \pm 0.04	56.91 \pm 0.05	66.12 \pm 0.08
DRCN	91.80 \pm 0.09	73.67 \pm 0.04	81.97 \pm 0.16	40.05 \pm 0.07	58.86 \pm 0.07	66.37 \pm 0.10
$ConvNet_{tgt}$	96.12 \pm 0.07	98.67 \pm 0.04	98.67 \pm 0.04	91.52 \pm 0.05	78.81 \pm 0.11	66.50 \pm 0.07

² The setup follows one in [18]: the inputs to SA are the last hidden layer activation values of $ConvNet_{src}$.

Comparison of different DRCN flavors: Recall that DRCN uses only the unlabeled target images for the unsupervised reconstruction training. To verify the importance of this strategy, we further compare different flavors of DRCN: DRCN_s and DRCN_{st} . Those algorithms are conceptually the same but different only in utilizing the unlabeled images during the unsupervised training. DRCN_s uses only unlabeled source images, whereas DRCN_{st} combines both unlabeled source and target images.

The experimental results in Table 2 confirm that DRCN always performs better than DRCN_s and DRCN_{st} . While DRCN_{st} occasionally outperforms ReverseGrad, its overall performance does not compete with that of DRCN. The only case where DRCN_s and DRCN_{st} flavors can closely match DRCN is on $\text{MN} \rightarrow \text{US}$. This suggests that the use of *unlabeled source data* during the reconstruction training do not contribute much to the cross-domain generalization, which verifies the DRCN strategy in using the unlabeled target data only.

Table 2. Accuracy (%) of DRCN_s and DRCN_{st} .

Methods	$\text{MN} \rightarrow \text{US}$	$\text{US} \rightarrow \text{MN}$	$\text{SV} \rightarrow \text{MN}$	$\text{MN} \rightarrow \text{SV}$	$\text{ST} \rightarrow \text{CI}$	$\text{CI} \rightarrow \text{ST}$
DRCN_s	89.92 ± 0.12	65.96 ± 0.07	73.66 ± 0.04	34.29 ± 0.09	55.12 ± 0.12	63.02 ± 0.06
DRCN_{st}	91.15 ± 0.05	68.64 ± 0.05	75.88 ± 0.09	37.77 ± 0.06	55.26 ± 0.06	64.55 ± 0.13
DRCN	91.80 ± 0.09	73.67 ± 0.04	81.97 ± 0.16	40.05 ± 0.07	58.86 ± 0.07	66.37 ± 0.10

Data reconstruction: A useful insight was found when reconstructing source images through the reconstruction pipeline of DRCN. Specifically, we observe the visual appearance of $f_r(x_1^s), \dots, f_r(x_m^s)$, where x_1^s, \dots, x_m^s are some images from the source domain. Note that x_1^s, \dots, x_m^s are unseen during the unsupervised reconstruction training in DRCN. We visualize such a reconstruction in the case of $\text{SV} \rightarrow \text{MN}$ training in Figure 3. Figure 2(a) and 3(a) display the original source (SVHN) and target (MNIST) images.

The main finding of this observation is depicted in Figure 3(c): the reconstructed images produced by DRCN given some SVHN images as the source inputs. We found that *the reconstructed SVHN images resemble MNIST-like digit appearances, with white stroke and black background*, see Figure 3(a). Remarkably, DRCN still can produce “correct” reconstructions of some noisy SVHN images. For example, all SVHN digits 3 displayed in Figure 2(a) are clearly reconstructed by DRCN, see the fourth row of Figure 3(c). DRCN tends to pick only the digit in the middle and ignore the remaining digits. This may explain the superior cross-domain recognition performance of DRCN on this task. However, such a cross-reconstruction appearance does not happen in the reverse task, $\text{MN} \rightarrow \text{SV}$, which may be an indicator for the low accuracy relative to the groundtruth performance.

We also conduct such a diagnostic reconstruction on other algorithms that have the reconstruction pipeline. Figure 3(d) depicts the reconstructions of the SVHN images produced by ConvAE trained on the MNIST images only. They do

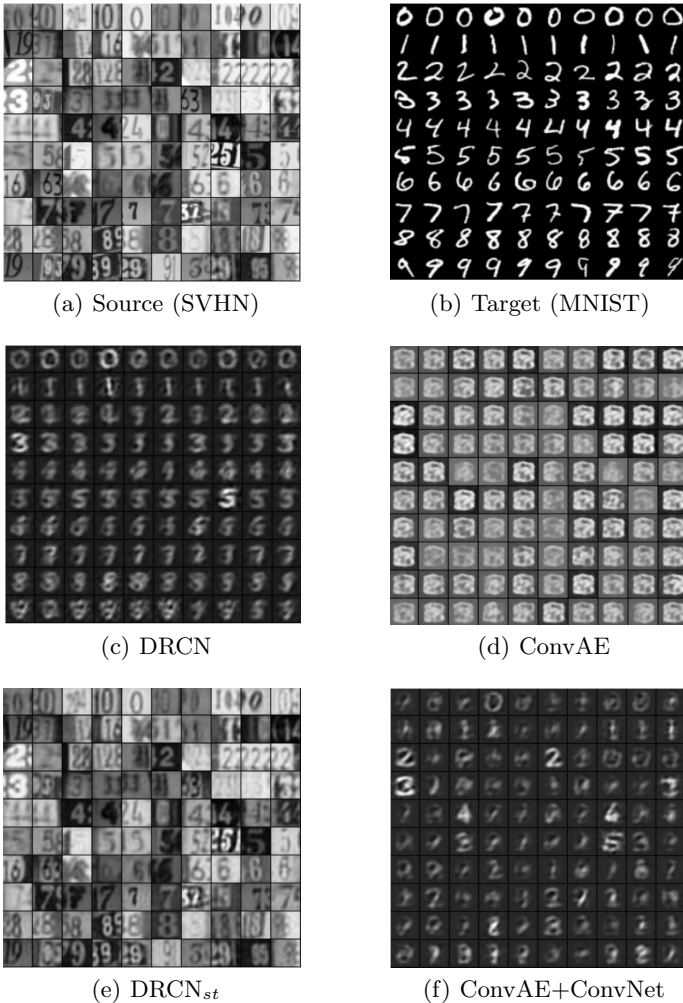


Fig. 2. Data reconstruction after training from SVHN \rightarrow MNIST. Fig. (a)-(b) show the original input pixels, and (c)-(f) depict the reconstructed source images (SVHN). The reconstruction of DRCN appears to be MNIST-like digits, see the main text for a detailed explanation.

not appear to be digits, suggesting that ConvAE recognizes the SVHN images as noise. Figure 3(e) shows the reconstructed SVHN images produced by DRCN_{st}. We can see that they look almost identical to the source images shown in Figure 2(a), which is not surprising since the source images are included during the reconstruction training.

Finally, we evaluated the reconstruction induced by ConvNet_{src} to observe the difference with the reconstruction of DRCN. Specifically, we trained ConvAE

on the MNIST images in which the encoding parameters were initialized from those of ConvNet_{src} and not updated during training. We refer to the model as ConvAE+ConvNet_{src}. The reconstructed images are visualized in Figure 3(f). Although they resemble the style of MNIST images as in the DRCN’s case, only a few source images are correctly reconstructed.

To summarize, the results from this diagnostic data reconstruction correlate with the cross-domain recognition performance. More visualization on other cross-domain cases can be found in the Supplemental materials.

4.2 Experiments II: Office dataset

In the second experiment, we evaluated DRCN on the standard domain adaptation benchmark for visual object recognition, OFFICE [11], which consists of three different domains: AMAZON (A), DSLR (D), and WEBCAM (W). OFFICE has 2817 labeled images in total distributed across 31 object categories. The number of images is thus relatively small compared to the previously used datasets.

We applied the DRCN algorithm to *finetune* AlexNet [14], as was done with different methods in previous work [18,40,41].³ The fine-tuning was performed only on the fully connected layers of AlexNet, *fc6* and *fc7*, and the last convolutional layer, *conv5*. Specifically, the label prediction pipeline of DRCN contains *conv4-conv5-fc6-fc7-label* and the data reconstruction pipeline has *conv4-conv5-fc6-fc7-fc6'-conv5'-conv4'* (the ' denotes the the inverse layer) – it thus does not reconstruct the original input pixels. The learning rate was selected following the strategy devised in [40]: cross-validating the base learning rate between 10^{-5} and 10^{-2} with a multiplicative step-size $10^{1/2}$.

We followed the standard unsupervised domain adaptation training protocol used in previous work [39,7,40], that is, using *all* labeled source data and unlabeled target data. Table 3 summarizes the performance accuracy of DRCN based on that protocol in comparison to the state-of-the-art algorithms. We found that DRCN is competitive against DAN and ReverseGrad – the performance is either the best or the second best except for one case. In particular, DRCN performs best with a convincing gap in situations when the target domain has relatively many data, i.e., AMAZON as the target dataset.

Table 3. Accuracy (*mean ± std %*) on the Office dataset with the standard unsupervised domain adaptation protocol used in [7,39].

Method	A → W	W → A	A → D	D → A	W → D	D → W
DDC [41]	61.8 ± 0.4	52.2 ± 0.4	64.4 ± 0.3	52.1 ± 0.8	98.5 ± 0.4	95.0 ± 0.5
DAN [40]	68.5 ± 0.4	53.1 ± 0.3	67.0 ± 0.4	54.0 ± 0.4	99.0 ± 0.2	96.0 ± 0.3
ReverseGrad [18]	72.6 ± 0.3	52.7 ± 0.2	67.1 ± 0.3	<u>54.5</u> ± 0.4	99.2 ± 0.3	96.4 ± 0.1
DRCN	68.7 ± 0.3	54.9 ± 0.5	66.8 ± 0.5	56.0 ± 0.5	99.0 ± 0.2	96.4 ± 0.3

³ Recall that AlexNet consists of five convolutional layers: *conv1*, ..., *conv5* and three fully connected layers: *fc6*, *fc7*, and *fc8/output*.

5 Analysis

This section provides a first step towards a formal analysis of the DRCN algorithm. We demonstrate that optimizing (5) in DRCN relates to solving a semi-supervised learning problem on the target domain according to a framework proposed in [19]. The analysis suggests that unsupervised training using only unlabeled target data is sufficient. That is, adding unlabeled source data might not further improve domain adaptation.

Denote the labeled and unlabeled distributions as $\mathbb{D}_{XY} =: \mathbb{D}$ and \mathbb{D}_X respectively. Let $P^\theta(\cdot)$ refer to a family of models, parameterized by $\theta \in \Theta$, that is used to learn a maximum likelihood estimator. The DRCN learning algorithm for domain adaptation tasks can be interpreted probabilistically by assuming that $P^\theta(x)$ is Gaussian and $P^\theta(y|x)$ is a multinomial distribution, fit by logistic regression.

The objective in Eq.(5) is equivalent to the following maximum likelihood estimate:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \lambda \sum_{i=1}^{n_s} \log P_{Y|X}^\theta(y_i^s | x_i^s) + (1 - \lambda) \sum_{j=1}^{n_t} \log P_{X|\tilde{X}}^\theta(x_j^t | \tilde{x}_j^t), \quad (6)$$

where \tilde{x} is the noisy input generated from $\mathbb{Q}_{\tilde{X}|X}$. The first term represents the model learned by the supervised convolutional network and the second term represents the model learned by the unsupervised convolutional autoencoder. Note that the discriminative model only observes labeled data from the source distribution \mathbb{P}_X in objectives (5) and (6).

We now recall a semi-supervised learning problem formulated in [19]. Suppose that labeled and unlabeled samples are taken from the *target domain* \mathbb{Q} with probabilities λ and $(1 - \lambda)$ respectively. By Theorem 5.1 in [19], the maximum likelihood estimate ζ is

$$\zeta = \operatorname{argmax}_{\zeta} \lambda \mathbb{E}_{\mathbb{Q}}[\log P^\zeta(x, y)] + (1 - \lambda) \mathbb{E}_{\mathbb{Q}_X}[\log P_X^\zeta(x)] \quad (7)$$

The theorem holds if it satisfies the following assumptions: *consistency*, the model contains true distribution, so the MLE is consistent; and *smoothness and measurability* [56]. Given target data $(x_1^t, y_1^t), \dots, (x_{n_t}^t, y_{n_t}^t) \sim \mathbb{Q}$, the parameter ζ can be estimated as follows:

$$\hat{\zeta} = \operatorname{argmax}_{\zeta} \lambda \sum_{i=1}^{n_t} [\log P^\zeta(x_i^t, y_i^t)] + (1 - \lambda) \sum_{i=1}^{n_t} [\log P_X^\zeta(x_i^t)] \quad (8)$$

Unfortunately, $\hat{\zeta}$ cannot be computed in the unsupervised domain adaptation setting since we do not have access to target labels.

Next we inspect a certain condition where $\hat{\theta}$ and $\hat{\zeta}$ are closely related. Firstly, by the *covariate shift* assumption [21]: $\mathbb{P} \neq \mathbb{Q}$ and $\mathbb{P}_{Y|X} = \mathbb{Q}_{Y|X}$, the first term in (7) can be switched from an expectation over target samples to source samples:

$$\mathbb{E}_{\mathbb{Q}} \left[\log P^\zeta(x, y) \right] = \mathbb{E}_{\mathbb{P}} \left[\frac{\mathbb{Q}_X(x)}{\mathbb{P}_X(x)} \cdot \log P^\zeta(x, y) \right]. \quad (9)$$

Secondly, it was shown in [57] that $P_{X|\tilde{X}}^\theta(x|\tilde{x})$, see the second term in (6), defines an ergodic Markov chain whose asymptotic marginal distribution of X converges to the data-generating distribution \mathbb{P}_X . Hence, Eq. (8) can be rewritten as

$$\hat{\zeta} \approx \underset{\zeta}{\operatorname{argmax}} \lambda \sum_{i=1}^{n_s} \frac{\mathbb{Q}_X(x_i^s)}{\mathbb{P}_X(x_i^s)} \log P^\zeta(x_i^s, y_i^s) + (1 - \lambda) \sum_{j=1}^{n_t} [\log P_{X|\tilde{X}}^\zeta(x_j^t|\tilde{x}_j^t)]. \quad (10)$$

The above objective differs from objective (6) only in the first term. Notice that $\hat{\zeta}$ would be approximately equal $\hat{\theta}$ if the ratio $\frac{\mathbb{Q}_X(x_i^s)}{\mathbb{P}_X(x_i^s)}$ is constant for all x^s . In fact, it becomes the objective of DRCN_{st}. Although the constant ratio assumption is too strong to hold in practice, comparing (6) and (10) suggests that $\hat{\zeta}$ can be a reasonable approximation to $\hat{\theta}$.

Finally, we argue that using unlabeled source samples during the unsupervised training may not further contribute to domain adaptation. To see this, we expand the first term of (10) as follows

$$\lambda \sum_{i=1}^{n_s} \frac{\mathbb{Q}_X(x_i^s)}{\mathbb{P}_X(x_i^s)} \log P_{Y|X}^\zeta(y_i^s|x_i^s) + \lambda \sum_{i=1}^{n_s} \frac{\mathbb{Q}_X(x_i^s)}{\mathbb{P}_X(x_i^s)} \log P_X^\zeta(x_i^s).$$

Observe the second term above. As $n_s \rightarrow \infty$, P_X^θ will converge to \mathbb{P}_X . Hence, since $\int_{x \sim \mathbb{P}_X} \frac{\mathbb{Q}_X(x)}{\mathbb{P}_X(x)} \log \mathbb{P}_X(x) \leq \int_{x \sim \mathbb{P}_X} \mathbb{P}_X^t(x)$, adding more unlabeled source data will only result in a constant. This implies an optimization procedure equivalent to (6), which may explain the *uselessness* of unlabeled source data in the context of domain adaptation.

Note that the latter analysis does not necessarily imply that incorporating unlabeled source data degrades the performance. The fact that DRCN_{st} performs worse than DRCN could be due to, e.g., the model capacity, which depends on the choice of the architecture.

6 Conclusions

We have proposed Deep Reconstruction-Classification Network (DRCN), a novel model for unsupervised domain adaptation in object recognition. The model performs multitask learning, i.e., alternately learning (source) label prediction and (target) data reconstruction using a shared encoding representation. We have shown that DRCN provides a considerable improvement for some cross-domain recognition tasks over the state-of-the-art model. It also performs better than deep models trained using the standard *pretraining-finetuning* approach. A useful insight into the effectiveness of the learned DRCN can be obtained from its data reconstruction. The appearance of DRCN's reconstructed source images resemble that of the target images, which indicates that DRCN learns the domain correspondence. We also provided a theoretical analysis relating the DRCN algorithm to semi-supervised learning. The analysis was used to support the strategy in involving only the target unlabeled data during learning the reconstruction task.

References

1. Torralba, A., Efros, A.A.: Unbiased Look at Dataset Bias. In: CVPR. (2011) 1521–1528
2. Hsu, H.J., Chen, K.T.: Face recognition on drones: Issues and limitations. In: Proceedings of ACM DroNet 2015. (2015)
3. Patel, V.M., Gopalan, R., Li, R., Chellapa, R.: Visual domain adaptation: A survey of recent advances. *IEEE Signal Processing Magazine* **32**(3) (2015) 53–69
4. Aljundi, R., Emonet, R., Muselet, D., Sebban, M.: Landmarks-Based Kernelized Subspace Alignment for Unsupervised Domain Adaptation. In: CVPR. (2015)
5. Baktashmotlagh, M., Harandi, M.T., Lovell, B.C., Salzmann, M.: Unsupervised Domain Adaptation by Domain Invariant Projection. In: ICCV. (2013) 769–776
6. Bruzzone, L., Marconcini, M.: Domain Adaptation Problems: A DASVM Classification Technique and a Circular Validation Strategy. *IEEE TPAMI* **32**(5) (2010) 770–787
7. Gong, B., Grauman, K., Sha, F.: Connecting the Dots with Landmarks: Discriminatively Learning Domain-Invariant Features for Unsupervised Domain Adaptation. In: ICML. (2013)
8. Long, M., Ding, G., Wang, J., Sun, J., Guo, Y., Yu, P.S.: Transfer Sparse Coding for Robust Image Representation. In: CVPR. (2013) 404–414
9. Long, M., Wang, J., Ding, G., Sun, J., Yu, P.S.: Transfer Joint Matching for Unsupervised Domain Adaptation. In: CVPR. (2014) 1410–1417
10. Pan, S.J., Tsang, I.W.H., Kwok, J.T., Yang, Q.: Domain adaptation via transfer component analysis. *IEEE Trans. Neural Networks* **22**(2) (2011) 199–210
11. Saenko, K., Kulis, B., Fritz, M., Darrell, T.: Adapting Visual Category Models to New Domains. In: ECCV. (2010) 213–226
12. Caruana, R.: Multitask Learning. *Machine Learning* **28** (1997) 41–75
13. Argyriou, A., Evgeniou, T., Pontil, M.: Multi-task Feature Learning. In: Advances in Neural Information Processing Systems 19. (2006) 41–48
14. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Classification with Deep Convolutional Neural Networks. In: NIPS. Volume 25. (2012) 1106–1114
15. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR. (2015)
16. Hinton, G.E., Osindero, S.: A fast learning algorithm for deep belief nets. *Neural Computation* **18**(7) (2006) 1527–1554
17. Bengio, Y., Lamblin, P., Popovici, D., Larochelle, H.: Greedy Layer-Wise Training of Deep Networks. In: NIPS. Volume 19. (2007) 153–160
18. Ganin, Y., Lempitsky, V.S.: Unsupervised domain adaptation by backpropagation. In: ICML. (2015) 1180–1189
19. Cohen, I., Cozman, F.G.: Risks of semi-supervised learning: how unlabeled data can degrade performance of generative classifiers. In: *Semi-Supervised Learning*. MIT Press (2006)
20. Japkowicz, N., Stephen, S.: The class imbalance problem: A systematic study. *Intelligent Data Analysis* **6**(5) (2002) 429–450
21. Shimodaira, H.: Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference* **90**(2) (2000) 227–244
22. Zadrozny, B.: Learning and evaluating classifiers under sample selection bias. In: Proceedings of the 21th Annual International Conference on Machine Learning. (2004) 114–121

23. Pan, S.J., Yang, Q.: A Survey on Transfer Learning. *IEEE Transactions on Knowledge and Data Engineering* **22**(10) (2010) 1345–1359
24. Blitzer, J., McDonald, R., Pereira, F.: Domain Adaptation with Structural Correspondence Learning. In: *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing*. (2006) 120–128
25. Daumé-III, H.: Frustratingly Easy Domain Adaptation. In: *Proceedings of ACL*. (2007)
26. Margolis, A.: A literature review of domain adaptation with unlabeled data. Technical report, University of Washington (2011)
27. Fernando, B., Habrard, A., Sebban, M., Tuytelaars, T.: Unsupervised Visual Domain Adaptation Using Subspace Alignment. In: *ICCV*. (2013) 2960–2967
28. Ghifary, M., Balduzzi, D., Kleijn, W.B., Zhang, M.: Scatter component analysis: A unified framework for domain adaptation and domain generalization. *CoRR abs/1510.04373* (2015)
29. Gopalan, R., Li, R., Chellapa, R.: Domain Adaptation for Object Recognition: An Unsupervised Approach. In: *ICCV*. (2011) 999–1006
30. Gong, B., Shi, Y., Sha, F., Grauman, K.: Geodesic Flow Kernel for Unsupervised Domain Adaptation. In: *CVPR*. (2012) 2066–2073
31. Bay, H., Tuytelaars, T., Gool, L.V.: SURF: Speeded Up Robust Features. *CVIU* **110**(3) (2008) 346–359
32. Glorot, X., Bordes, A., Bengio, Y.: Domain Adaptation for Large-Scale Sentiment Classification: A Deep Learning Approach. In: *ICML*. (2011) 513–520
33. Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., Manzagol, P.A.: Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion. *Journal of Machine Learning Research* **11** (2010) 3371–3408
34. Ghifary, M., Kleijn, W.B., Zhang, M.: Domain adaptive neural networks for object recognition. In: *PRICAI: Trends in AI. Volume 8862*. (2014) 898–904
35. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. In: *Proceedings of the IEEE. Volume 86*. (1998) 2278–2324
36. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: *CVPR*. (2014)
37. Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., Darrell, T.: DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition. In: *ICML*. (2014)
38. Hoffman, J., Tzeng, E., Donahue, J., Jia, Y., Saenko, K., Darrell, T.: One-Shot Adaptation of Supervised Deep Convolutional Models. *CoRR abs/1312.6204* (2013)
39. Chopra, S., Balakrishnan, S., Gopalan, R.: DLID: Deep Learning for Domain Adaptation by Interpolating between Domains. In: *ICML Workshop on Challenges in Representation Learning*. (2013)
40. Long, M., Cao, Y., Wang, J., Jordan, M.I.: Learning transferable features with deep adaptation networks. In: *ICML*. (2015)
41. Tzeng, E., Hoffman, J., Zhang, N., Saenko, K., Darrell, T.: Deep domain confusion: Maximizing for domain invariance. *CoRR abs/1412.3474* (2014)
42. Tzeng, E., Hoffman, J., Darrell, T., Saenko, K.: Simultaneous deep transfer across domains and tasks. In: *ICCV*. (2015)
43. Borgwardt, K.M., Gretton, A., Rasch, M.J., Kriegel, H.P., Schölkopf, B., Smola, A.J.: Integrating structured biological data by Kernel Maximum Mean Discrepancy. *Bioinformatics* **22**(14) (2006) e49–e57

44. Masci, J., Meier, U., Ciresan, D., Schmidhuber, J.e.: Stacked Convolutional Auto-Encoders for Hierarchical Feature Extraction. In: ICANN. (2011) 52–59
45. Zeiler, M.D., Krishnan, D., Taylor, G.W., Fergus, R.: Deconvolutional networks. In: CVPR. (2010) 2528–2535
46. Tieleman, T., Hinton, G.: Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning (2012)
47. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A Simple Way to Prevent Neural Networks from Overfitting. JMLR (2014)
48. Bastien, F., Lamblin, P., Pascanu, R., Bergstra, J., Goodfellow, I.J., Bergeron, A., Bouchard, N., Bengio, Y.: Theano: new features and speed improvements. In: Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop. (2012)
49. Simard, P.Y., Steinkraus, D., Platt, J.C.: Best practices for convolutional neural networks applied to visual document analysis. In: ICDAR. Volume 2. (2003) 958–962
50. Hull, J.J.: A database for handwritten text recognition research. IEEE TPAMI **16**(5) (1994) 550–554
51. Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A.Y.: Reading Digits in Natural Images with Unsupervised Feature Learning. In: NIPS Workshop on Deep Learning and Unsupervised Feature Learning. (2011)
52. Krizhevsky, A.: Learning Multiple Layers of Features from Tiny Images. Master’s thesis, Department of Computer Science, University of Toronto (April 2009)
53. Coates, A., Lee, H., Ng, A.Y.: An Analysis of Single-Layer Networks in Unsupervised Feature Learning. In: AISTATS. (2011) 215–223
54. Sermanet, P., Chintala, S., LeCun, Y.: Convolutional neural networks applied to house number digit classification. In: ICPR. (2012) 3288–3291
55. Nair, V., Hinton, G.E.: Rectified Linear Units Improve Restricted Boltzmann Machines. In: ICML. (2010)
56. White, H.: Maximum likelihood estimation of misspecified models. *Econometrica* **50**(1) (1982) 1–25
57. Bengio, Y., Yao, L., Guillaume, A., Vincent, P.: Generalized denoising auto-encoders as generative models. In: NIPS. (2013) 899–907
58. van der Maaten, L., Hinton, G.: Visualizing High-Dimensional Data Using t-SNE. *Journal of Machine Learning Research* **9** (2008) 2579–2605

Supplemental Material

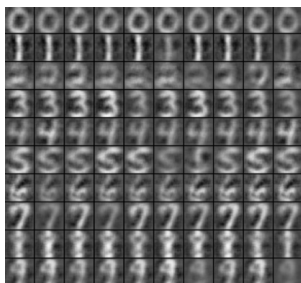
This document is the supplemental material for the paper *Deep Reconstruction-Classification for Unsupervised Domain Adaptation*. It contains some more experimental results that cannot be included in the main manuscript due to a lack of space.



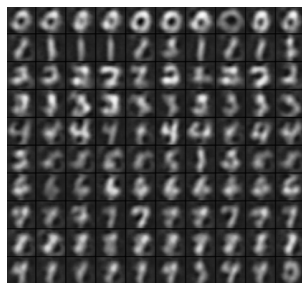
(a) Source (MNIST)



(b) Target (USPS)



(c) DRCN



(d) ConvAE

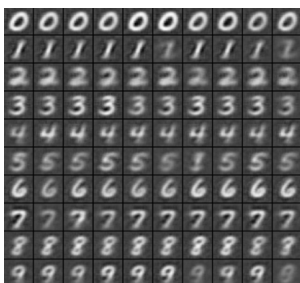
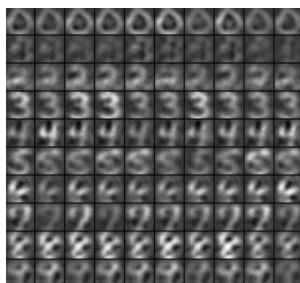
(e) DRCN_{st}(f) ConvAE+ConvNet_{src}

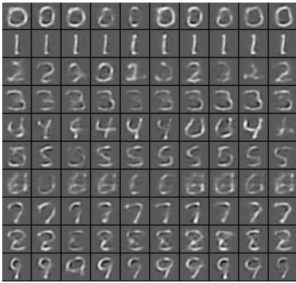
Fig. 3. Data reconstruction after training from MNIST \rightarrow USPS. Fig. (a)-(b) show the original input pixels, and (c)-(f) depict the reconstructed source images (MNIST).



(a) Source (USPS)



(b) Target (MNIST)



(c) DRCN



(d) ConvAE

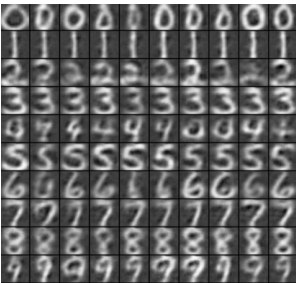
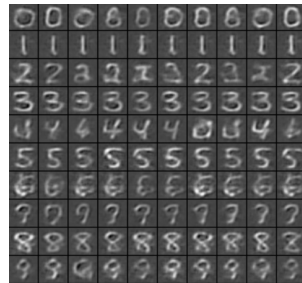
(e) DRCN_{st}(f) ConvAE+ConvNet_{src}

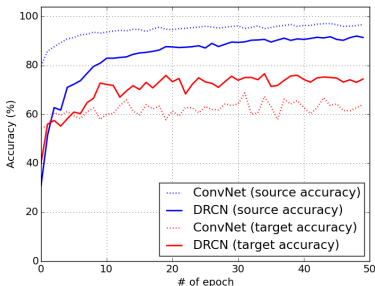
Fig. 4. Data reconstruction after training from USPS \rightarrow MNIST. Fig. (a)-(b) show the original input pixels, and (c)-(f) depict the reconstructed source images (USPS).

Data Reconstruction

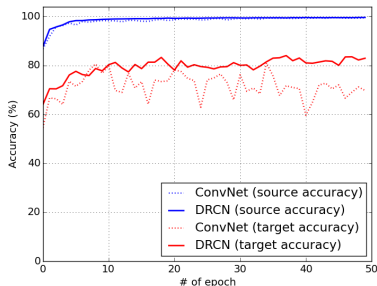
Figures 3 and 4 depict the reconstruction of the source images in cases of MNIST \rightarrow USPS and USPS \rightarrow MNIST, respectively. The trend of the outcome is similar to that of SVHN \rightarrow MNIST, see Figure 2 in the main manuscript. That is, the reconstructed images produced by DRCN resemble the *style* of the target images.

Training Progress

Recall that DRCN has two pipelines with a shared encoding representation; each corresponds to the classification and reconstruction task, respectively. One can consider that the unsupervised reconstruction learning acts as a regularization for the supervised classification to reduce overfitting onto the source domain. Figure 5 compares the source and target accuracy of DRCN with that of the standard ConvNet during training. The most prominent results indicating the overfitting reduction can be seen in SVHN \rightarrow MNIST case, i.e., DRCN produces higher target accuracy, but with lower source accuracy, than ConvNet.



(a) SVHN \rightarrow MNIST training



(b) MNIST \rightarrow USPS training

Fig. 5. The source accuracy (blue lines) and target accuracy (red lines) comparison between ConvNet and DRCN during training stage on SVHN \rightarrow MNIST cross-domain task. DRCN induces lower source accuracy, but higher target accuracy than ConvNet.

t-SNE visualization.

For completeness, we also visualize the 2D point cloud of the last hidden layer of DRCN using t-SNE [58] and compare it with that of the standard ConvNet. Figure 6 depicts the feature-point clouds extracted from the target images in the case of MNIST \rightarrow USPS and SVHN \rightarrow MNIST. Red points indicate the source feature-point cloud, while gray points indicate the target feature-point cloud. Domain invariance should be indicated by the degree of overlap between the source and target feature clouds. We can see that the overlap is more prominent in the case of DRCN than ConvNet.

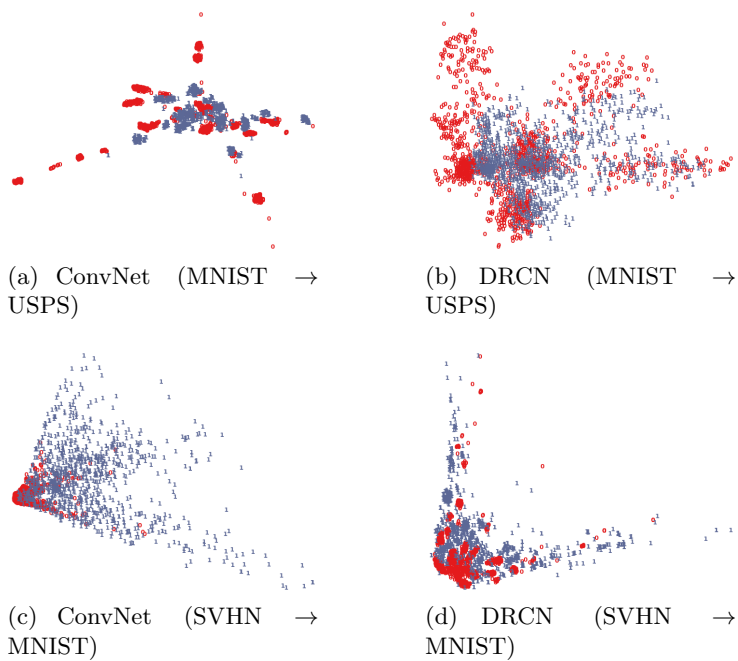


Fig. 6. The t-SNE visualizations of the last layer's activations. Red and gray points indicate the source and target domain examples, respectively.