

Received October 1, 2019, accepted October 21, 2019, date of publication October 25, 2019, date of current version November 7, 2019.

Digital Object Identifier 10.1109/ACCESS.2019.2949697

Self-Supervised Domain Adaptation for Computer Vision Tasks

JIAOLONG XU¹, LIANG XIAO¹, AND ANTONIO M. LÓPEZ², (Member, IEEE)

¹Unmanned System Research Center, NIIT, Beijing 100071, China

²Computer Vision Center, Universitat Autònoma de Barcelona, 08193 Bellaterra, Spain

Corresponding author: Jiaolong Xu (jiaolong_xu@126.com)

This work was supported in part by the National Natural Science Foundation of China (NSFC) under Grant 61601508, Grant 61790565, and Grant 61803380, in part by the Spanish project under Grant TIN2017-88709-R (MINECO/AEI/FEDER, UE), and in part by the ICREA under the ICREA Academia Program.

ABSTRACT Recent progress of self-supervised visual representation learning has achieved remarkable success on many challenging computer vision benchmarks. However, whether these techniques can be used for domain adaptation has not been explored. In this work, we propose a generic method for self-supervised domain adaptation, using object recognition and semantic segmentation of urban scenes as use cases. Focusing on simple pretext/auxiliary tasks (e.g. image rotation prediction), we assess different learning strategies to improve domain adaptation effectiveness by self-supervision. Additionally, we propose two complementary strategies to further boost the domain adaptation accuracy on semantic segmentation within our method, consisting of prediction layer alignment and batch normalization calibration. The experimental results show adaptation levels comparable to most studied domain adaptation methods, thus, bringing self-supervision as a new alternative for reaching domain adaptation. The code is available at this link.¹

INDEX TERMS Domain adaptation, semantic segmentation, object recognition.

I. INTRODUCTION

Since supervised (deep) machine learning became the key to solve computer vision tasks, the availability of task ground truth (*i.e.* supervision information) associated to the raw data (*i.e.* images and videos) has been a major practical problem. Training an image or video classifier requires to associate some class or attributes to the whole image/video [11], [27], [28], [53], training an object detector requires manual drawing of object bounding boxes [14], [32], training a CNN for semantic segmentation requires the delineation of the borders between the considered classes [10], [41], etc. This kind of ground truth (bounding boxes, class borders) is usually provided by human labeling, which is a costly process prone to errors due to subjectivity and fatigue. Therefore, procedures aiming at reducing human labeling became a research topic in itself too; or alternatively obtaining the most from a fixed budget for new labels. This underlying aim appears under different names depending on the practical situation at hand, *i.e.* the learning conditions. Under this umbrella we

The associate editor coordinating the review of this manuscript and approving it for publication was Kang Li¹.

¹<https://github.com/Jiaolong/self-supervised-da>

find concepts such as active learning, self-labeling, transfer learning, domain adaptation, and self-supervision.

In *active learning* [1], [47], [52], the learner receives a set of unlabeled data (images/videos) for training a visual accurate model, which must be done minimizing the labeling effort by choosing the best training data out of the total amount. This turns out into an iterative process where a human worker labels new automatically selected data in each cycle for model refinement. This contrasts with passive learning, where the training data is selected at random, eventually requiring more labeling budget.

In *self-labeling* [62], [63], [69], an initial visual model is trained on labeled data, after, the model is applied on unlabeled data to self-collect samples which are used then for refining the model by assuming that their label corresponds to the prediction of the model; turning out in an iterative process that must avoid drifting to systematic errors or easy samples.

In *transfer learning* [18], [43], a model is trained to perform a visual task (*e.g.* image classification) but aiming at reusing it to perform a new task (*e.g.* object detection) in a way that we minimize the amount of labeled data required to train for the new task (*e.g.* fine-tuning CNNs across tasks is a basic form of transfer learning).

In *domain adaptation* [8], [21], [60], [67], a model is trained to perform a visual task in a specific domain (e.g. semantic segmentation in synthetic images), however, we need to apply it to perform the same task in a correlated, but significantly different, domain (e.g. semantic segmentation in real-world images); which is done by reusing the previous knowledge (in the form of model or labeled data) for minimizing the labeling effort in the new domain.

Finally, *self-supervised learning* [17], [25], [26] focuses on learning visual models without manual labeling; more specifically, auxiliary relatively simple tasks, known as *pretext tasks* in this context, are created for training a generic visual model in the form of CNN. The supervision consists in modifying the original visual data (e.g. a set of images) according to known transforms (e.g. image rotations [17]), training the pretext CNN to predict such transforms; thus, the transforms are the labels/supervision for the pretext task. This pretext CNN is then concatenated with another task-specific CNN. The former acting as generic feature extractor, and the later leveraging such features to create new ones specific for the *main task* of interest. Sometimes, both CNN blocks are fine-tuned [4], and sometimes the pretext CNN block is frozen and only the task-specific CNN block is fine-tuned [26]. Overall, the idea is that we can have a high number of supervised samples for the pretext task and this should compensate for a lower number of manually labeled samples for the main task.

Active learning can be naturally combined with transfer learning or domain adaptation [58]. Self-labeling can also be combined with transfer learning or domain adaptation [62]. Self-supervised learning, as usually performed, can be seen as a type of transfer learning (from the pretext task to the main task). What has not been explored, up to the best of our knowledge, is how self-supervised learning can support domain adaptation. This is the main focus of this paper, *i.e.* can we incorporate self-supervision to learn domain invariant feature representation? The goal of this work is not to propose new self-supervised learning methods but investigate how existing self-supervised representation learning methods can be used to address domain adaptation problems. With this aim, we design a multi-task learning method to jointly train pretext and main tasks (Figure 1). The pretext task acts as nexus between source and target domains for learning a domain invariant feature representation for the main task. In this way, we have labels for the main task in source domain, but we do not require labels for such task in the target domain. In other words, via self-supervised learning, we perform unsupervised domain adaptation.

Accordingly, and using object recognition and semantic segmentation of urban scenes as challenging main-task use cases, the main contributions of this work are three-fold:

- We proposed a generic method for domain adaptation with self-supervised visual representation learning.
- Focusing on the image rotation prediction pretext learning task, we proposed several variations and studied their domain adaptation performance.

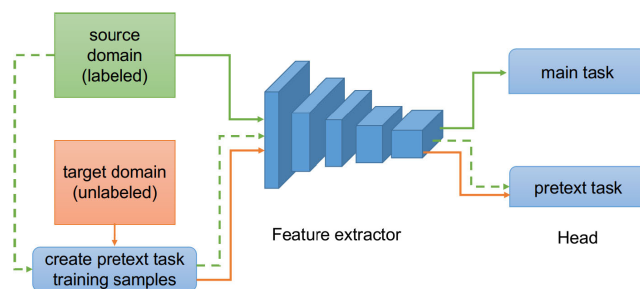


FIGURE 1. Proposed self-supervised domain adaptation framework. We learn a domain invariant feature representation by incorporating a pretext learning task which can automatically create labels from target domain images. The pretext and main task (e.g. object recognition or semantic segmentation) are learned jointly via multi-task learning. Solid lines indicate the forwarded data flow and the dash lines indicate optional data flow.

- We proposed additional strategies to further boost the self-supervised domain adaptation, including prediction layer alignment and batch normalization calibration.

This paper is organized as follows. In Section II, we review related self-supervised representation learning and domain adaptation methods. In Section III, we explain the proposed method. In Section IV, we conduct experiments on domain adaptation for object recognition as well as semantic segmentation, via our method. Finally, Section V summarizes the work and future directions.

II. RELATED WORK

A. SELF-SUPERVISED VISUAL REPRESENTATION LEARNING

An extensive review of deep learning-based self-supervised general visual feature learning methods from images or videos is provided in [24]. The recent work of self-supervised representation learning mainly focus on the design of pretext tasks. The work of [26] gives a comprehensive study of some state-of-the-art methods. A pretext task of predicting the relative location of image patches was first proposed in [12], where the patch ID is the supervision/label. This initial patch-based method has been followed by several variants [25], [40], [42]. Other works incorporate image colorization [65] or image inpainting [44] as pretext tasks. Yet other works focus on automatic ways of creating image samples with corresponding labels; for instance, in [4] the labels are classes derived from unsupervised image clustering, and in [17] the labels are image rotation angles since from an original image four possible rotations were created. As compared in [26], the rotation prediction based method [17] has shown promising results for learning high-level image representations. The rotation based method is further improved in [15] by decoupling rotation related and unrelated features. Therefore, in this work, we employ this pretext task as well as the location of image patches in line with [12]. In [23], relative depth prediction is used as a self-supervised proxy task, which has shown improvements to the downstream tasks, including

semantic segmentation and car detection. However, it relies on the video data in order to obtain the relative depth.

B. UNSUPERVISED DOMAIN ADAPTATION

There have been numerous domain adaptation methods proposed for object recognition since [48]. After the pioneer work of [21], [67], semantic segmentation has also aroused increasing interests. Among existing domain adaptation methods, some try to align domains at input level, including GAN-based methods [5], [20], [31] and image stylization ones [13], [61], [68]. Some focus at feature level adaptation [9], [21], [22], [45], [49], and others on adapting the output space [39], [50], [55], [59]. According to recent surveys [60], [66], most methods are built on the principle of domain adversarial training [16], with differences on how to incorporate it to the training of the segmentation network. Among the adaptation strategies we use as complement to self-supervision, the prediction layer alignment is similar to adversarial training for output space alignment.

In [69], iterative self-labeling and fine-tuning with spatial urban-scene location priors are used to perform the domain adaptation. In [67], a curriculum learning style is applied, where super-pixels are computed in source and target domains and their distributions must match as auxiliary task during semantic segmentation training. The use of such auxiliary task is similar in spirit to our multi-task learning approach with pretext tasks as nexus between source and target domains. However, neither our auxiliary tasks nor our complementary adaptation strategies are restricted to semantic segmentation, and they are way simpler than computing super-pixels. Comparing to these work, our method is not specifically designed for semantic segmentation but generic for various computer vision tasks.

In [2], the self-supervised learning method jigsaw puzzle is used for object recognition domain generalization and adaptation. As we will see in the experimental section, our method outperforms the jigsaw puzzle based method on both object recognition and semantic segmentation tasks. For semantic segmentation, we compare our results to [13], [20], [22], [49], [55], [61], [67]–[69]. The final semantic segmentation accuracy we obtain in target domain is superior to most of these methods, only behind [69] which is specific for semantic segmentation, and still not being far apart. Moreover, although it is out of the scope of this paper, our method can be complementary to some of the ones aforementioned, such as those based on adapting the input images via GANs.

III. METHOD

In this section, we first introduce our generic framework of self-supervised domain adaptation. Then, we present the considered pretext tasks. Finally, we introduce domain adaptation steps which complement self-supervision.

A. SELF-SUPERVISED DOMAIN ADAPTATION

1) OVERVIEW OF THE FRAMEWORK

Taking semantic segmentation as an example of main task, but without lose of generality, our method is shown in

FIGURE 1; where E denotes an encoder network (feature extractor) and S a decoder network (specific of the main task), so that $E + S$ is a CNN for semantic segmentation. This CNN is trained end-to-end with source domain labeled samples, $\{X_s, Y_s\}$. We denote by P the network added to support the creation of a model for solving the pretext task. This model consists in the CNN $E + P$, where E is shared with the CNN of the main task. The pretext task training samples, $\{X_t, Y_t\}$, are automatically created from the target domain images so that the training of $E + P$ is also supervised.

The complete domain adaptation method is drawn in Algorithm 1, where we can see how the self-supervised domain adaptation is a joint training of models to perform the pretext and main tasks. During the forward propagation, both source and target domain samples pass through the shared encoder. After, the losses of the main task \mathcal{L}_{seg} and pretext task \mathcal{L}_p are computed, they are back-propagated and accumulated at the encoder. Because the encoder is trained with both source and target domain samples, it learns domain invariant feature representations. In the testing phase, we feed the target domain images to the encoder and pass the features to the decoder of the main task to obtain the predictions.

Algorithm 1 Self-Supervised Domain Adaptation

Data: Labeled source domain images: $\{X_s, Y_s\}$, and unlabeled target domain images: $\{X_t\}$

Result: Model trained for main task in target domain
Create samples for pretext task: $\{X_t, Y_t\}$;

$i = 0$;

while $i < \max_iters$ **do**

 Load target mini-batch $\{x_t^i, y_t^i\}$;

 Forward pass and compute \mathcal{L}_p ;

 Back-propagate \mathcal{L}_p gradients by P and E ;

 Update weights of P ;

 Load source mini-batch $\{x_s^i, y_s^i\}$;

 Forward pass and compute \mathcal{L}_{seg} ;

 Back-propagate \mathcal{L}_{seg} gradients by S and E ;

 Accumulate gradients from \mathcal{L}_p and \mathcal{L}_{seg} for E ;

 Update weights of E and S ;

end

It is also possible to create pretext task samples with the source domain data, *i.e.*, dash lines in FIGURE 1. In this case, the pretext model can be trained with both source and target domain pretext task samples. We investigate this in Section IV.

2) PRETEXT TASKS

In this section, we first introduce the image rotation prediction pretext task. Inspired by the image-patch based methods [12], [25], [40], [42], we also take into account the spatial layout of the image and propose a new pretext task.

a: IMAGE ROTATION PREDICTION AS PRETEXT TASK

We select image rotation prediction as pretext task due to its simplicity and superior performance on visual representation learning to other proposals [17]. Given a set of N_t training images from target domain $D_t = \{x_i^t\}_{i=0}^{N_t}$, similar to [17], we define the set of geometric transformations as 2D image rotations by 0, 90, 180 and 270 degrees. We denote the rotation function by $g(x_i^t, r)$, $r \in [0, 3]$ rotates image x_i^t by $r * 90$ degrees. The geometric transformation prediction model P takes feature map from E as input and outputs a probability distribution over all possible geometric transformations. The self-supervised training objective that the geometric transformation model must learn to solve is:

$$\min_{\theta_e, \theta_p} \frac{1}{N_t} \sum_{i=1}^{N_t} \mathcal{L}_p(x_i^t, \theta_e, \theta_p), \quad (1)$$

where θ_e and θ_p are the parameters of the encoder E and pretext network P respectively, \mathcal{L}_p is the loss function defined as:

$$\mathcal{L}_p = -\frac{1}{4} \sum_{r=0}^3 \log(P(E(g(x_i^t, r), \theta_e), \theta_p)). \quad (2)$$

By learning to predict the image orientations, the convolutional neural networks also implicitly learn to localize salient objects in the images, recognize their orientations and object types [17]. Such implicitly learned knowledge contains semantic information of the target domain images which is expected to improve the cross-domain feature representation power of the encoder network. In other words, the pretext task with target domain images helps the encoder to learn domain invariant feature representation, thus, helps to achieve domain adaptation.

The work of [17] uses full images from ImageNet [11]. However, the images from a specific domain are usually biased to particular structures or patterns, especially at a full image level. If we train a rotation prediction model with full images, the training process could find a trivial solution and, thus, not being able to learn a domain invariant feature representation. To avoid this problem, we first randomly crop an image patch from the full image and then rotate this patch. In this way, we create more difficult and diverse samples for the pretext task.

b: SPATIAL-AWARE ROTATION PREDICTION AS PRETEXT TASK

Beyond image rotation, we further propose to take into account the image spatial layout to create a more complex pretext task. As depicted in FIGURE 2, instead of randomly cropping a patch from the full image, we first split the full image into four regions. From each region, we apply cropping and rotation operations as in the previous pretext task. We call this strategy spatial-aware rotation prediction. The dimension of a label is then extended from 4 (rotation angles) to 16 (spatial locations times rotation angles). This scheme encodes

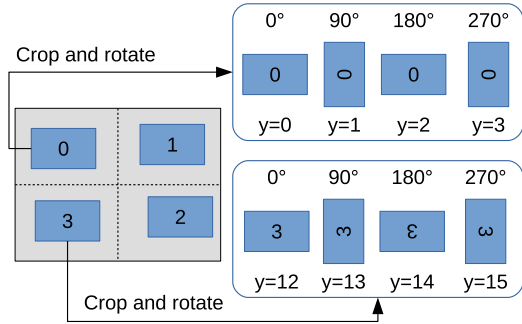


FIGURE 2. Region-based cropping and rotation.

the geometry transform as well as spatial layout information, which results in a more complex pretext task.

3) OBJECTIVE FUNCTION FOR DOMAIN ADAPTATION

Given a set of N_s labeled training images from the source domain $D_s = \{x_i^s, y_i^s\}_{i=0}^{N_s}$, the segmentation network takes as input the feature maps from $E(x_i^s)$ and outputs the segmentation predictions: $O_i^s = S(E(x_i^s), \theta_e), \theta_s) \in \mathbb{R}^{H \times W \times C}$, where C is the number of semantic categories, H and W are the height and width of the output respectively, and θ_e and θ_s convey the parameters of E and S , respectively. The semantic segmentation training objective that we need to solve for E and S is:

$$\min_{\theta_e, \theta_s} \frac{1}{N_s} \sum_{i=1}^{N_s} \mathcal{L}_{seg}(x_i^s, \theta_e, \theta_s), \quad (3)$$

where the segmentation loss is the cross-entropy loss, defined as:

$$\mathcal{L}_{seg} = -\sum_{h,w} \sum_{c \in C} y_i^s(h, w, c) \log(O_i^s(h, w, c)). \quad (4)$$

With Eq. (1) and Eq. (3), the objective function that self-supervised domain adaptation must solve is:

$$\min_{\theta_e, \theta_p, \theta_s} \frac{1}{N_s} \sum_{i=1}^{N_s} \mathcal{L}_{seg}(x_i^s, \theta_e, \theta_s) + \frac{\lambda_p}{N_t} \sum_{j=1}^{N_t} \mathcal{L}_p(x_j^t, \theta_e, \theta_p), \quad (5)$$

where λ_p is the weight to balance the two losses. In this work, we simply set $\lambda_p = 1$ for our experiments. The training process follows Algorithm 1.

B. COMPLEMENTARY ADAPTATION STEPS

In this section, we introduce two different strategies to complement self-supervised domain adaptation, including adversarial training for prediction layer alignment and batch normalization.

1) PREDICTION LAYER ALIGNMENT

The proposed pretext task learning is able to perform domain adaptation at feature level, however, the predicted semantic labels may still not be well aligned. There have been some

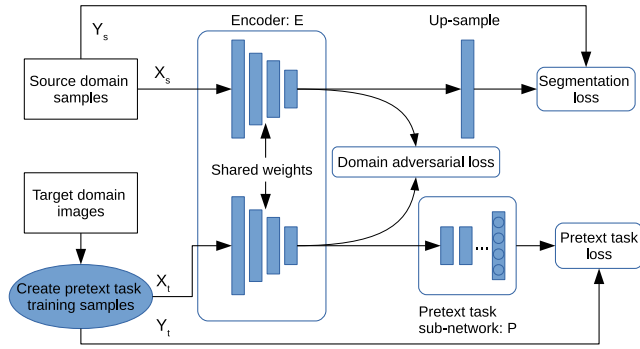


FIGURE 3. Self-supervised domain adaptation with prediction layer alignment.

previous work tackling this problem [39], [55]. In this work, we also consider to align the prediction layer to improve the domain adaptation performance. The main idea is illustrated in FIGURE 3. For semantic segmentation, we simplified the decoder by a single up-sampling layer. In this way, the last layer of the encoder is corresponding to the prediction layer. By placing a domain discriminator after the prediction layer, the commonly used domain adversarial training can be employed. We denote by D the discriminator and θ_d for its parameters. Given an input image x_i , the discriminator takes as input the feature maps from the encoder $E(x_i)$ and performs the binary classification to distinguish whether the feature map is from the source image or the target one, $\mathbf{Z}_i = D(E(x_i))$, $\mathbf{Z}_i \in \mathbb{R}^{H \times W \times 2}$. The training of D is a standard supervised training, which minimizes the following 2-D cross-entropy loss:

$$\mathcal{L}_d(x_i, \theta_d) = - \sum_{h,w} [(1-z) \log \mathbf{Z}_i(h, w, 0) + z \log \mathbf{Z}_i(h, w, 1)], \quad (6)$$

where h, w are indexing the output layer, $z = 0$ indicates that the sample is drawn from the target domain, and $z = 1$ if it is drawn from the source domain.

In order to learn a domain invariant feature representation, we want the encoder to fool the domain discriminator D , which is equivalent to minimize the following adversarial loss function:

$$\mathcal{L}_{adv}(x_i, \theta_e) = - \sum_{h,w} [(1-z) \log \mathbf{Z}_i(h, w, 1) + z \log \mathbf{Z}_i(h, w, 0)]. \quad (7)$$

\mathcal{L}_{adv} encourages to fool D by optimizing θ_e while \mathcal{L}_d encourages to improve the classification accuracy of D by optimizing θ_d . The optimization of Eq. (6) and Eq. (7) is essentially a domain adversarial training. Combining the self-supervised domain adaptation objective function Eq. (5), the overall optimization problem that we solve is as following:

$$\min_{\theta_e, \theta_p, \theta_s, \theta_d} \frac{1}{N_s} \sum_{i=1}^{N_s} \mathcal{L}_{seg}(x_i^s, \theta_e, \theta_s)$$

$$+ \frac{\lambda_p}{N_t} \sum_{i=1}^{N_t} \mathcal{L}_p(x_i^t, \theta_e, \theta_p) + \frac{\lambda_{adv}}{N_t} \sum_{i=1}^{N_t} \mathcal{L}_{adv}(x_i, \theta_e) + \frac{\lambda_d}{N_t + N_s} \sum_{i=1}^{N_t + N_s} \mathcal{L}_d(x_i, \theta_d), \quad (8)$$

where λ_{adv} and λ_d are the weights to balance the corresponding losses. These hyper-parameters are tuned on the validation set and then fixed for all experiments. In this work, we set $\lambda_{adv} = 0.01$ and $\lambda_d = 1.0$ for our experiments. We show how this prediction layer alignment improves the self-supervised domain adaptation in Section IV-C4.

2) BATCH NORMALIZATION CALIBRATION

The batch normalization (BN) is originally designed to reduce the internal covariate shift and speedup the training of deep neural networks. Given a mini-batch $\mathcal{B} = \{z_1, \dots, z_m\}$ as input, BN layer first calculates the mean and variance by $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m z_i$, $\sigma^2 = \frac{1}{m} \sum_{i=1}^m (z_i - \mu_{\mathcal{B}})^2$. Each example is then normalized by $\hat{z}_i = \frac{z_i - \mu_{\mathcal{B}}}{\sqrt{\sigma^2 + \epsilon}}$, where ϵ is a constant added to the mini-batch variance for numerical stability. The normalized values are then scaled and shifted by $\lambda \hat{z}_i + \beta$ to produce the output, where λ and β are learnable parameters.

For a trained source domain model, $\mu_{\mathcal{B}}$ and σ^2 are statistics from source domain images, which may cause domain shift when applied with target domain images. Although during domain adaptation training, both of source and target data are passed through the BN layers, the statistics from both domain can be still ambiguous for the target domain model. What we proposed in this work is to re-calibrate these statistics to reduce the domain shift. Given a pretrained network, we keep all the learnable parameters fixed and feed forward the target domain training images. During this forward propagation, we re-calculate the mean and variance values of each BN layer.

Our BN calibration is similar to the AdaBN method [30]. However, AdaBN adopts an online algorithm to estimate the mean and variance, while we simply use the common moving average mean and variance available in existing deep learning frameworks. AdaBN is applied at the inference stage, i.e. to the testing images, while we use BN calibration as a post training process with target domain training images.

IV. EXPERIMENTS AND RESULTS

In this section, we conduct experiments to validate the proposed domain adaptation method for both object recognition and semantic segmentation.

A. IMPLEMENTATION DETAILS

We implement the proposed method using the PyTorch framework on a single GTX 1080 Ti GPU with 11 GB memory. For object recognition, we use the code base of

TABLE 1. The abbreviations of the the evaluated methods.

Methods	Description
Ours(Jigsaw)	Self-supervised domain adaptation with pretext task of solving jigsaw puzzle. [2]
Ours(Rot)	Self-supervised domain adaptation with image rotation prediction pretext task.
Ours(MixRot)	Same as Rot but mixing with source domain samples in the pretext task learning.
Ours(SPRot)	Self-supervised domain adaptation with spatial-aware rotation prediction pretext task.
Ours(Adv)	Adversarial domain adaptation with prediction layer alignment.
Ours(Rot+Adv)	Rot with Adv as complementary strategy.
Ours(Rot+Adv+BN)	Rot with Adv and batch normalization calibration as complementary strategies.

TABLE 2. Accuracy (%) on Office dataset (ResNet-50).

Method	$A \rightarrow W$	$D \rightarrow W$	$W \rightarrow D$	$A \rightarrow D$	$D \rightarrow A$	$W \rightarrow A$	Avg.
ResNet-50 [19]	68.4±0.2	96.7±0.1	99.3±0.1	68.9±0.2	62.5±0.3	60.7±0.3	76.1
JDDA [6]	82.6±0.4	95.2±0.2	99.7±0.0	79.8±0.1	57.4±0.0	66.7±0.2	80.2
DAN [33]	80.5±0.4	97.1±0.2	99.6±0.1	78.6±0.2	63.6±0.3	62.8±0.2	80.4
RTN [35]	84.5±0.2	96.8±0.1	99.4±0.1	77.5±0.3	66.2±0.2	64.8±0.3	81.6
DANN [16]	82.0±0.4	96.9±0.2	99.1±0.1	79.7±0.4	68.2±0.4	67.4±0.5	82.2
ADDA [56]	86.2±0.5	96.2±0.3	98.4±0.3	77.8±0.3	69.5±0.4	68.9±0.5	82.9
JAN [36]	85.4±0.3	97.4±0.2	99.8±0.2	84.7±0.3	68.6±0.3	70.0±0.4	84.3
MADA [45]	90.0±0.1	97.4±0.1	99.6±0.1	87.8±0.1	70.3±0.3	66.4±0.3	85.2
GTA [51]	89.5±0.5	97.9±0.3	99.8±0.4	87.7±0.5	72.8±0.3	71.4±0.4	86.5
CDAN [34]	94.1±0.1	98.6±0.1	100.0±0.0	92.9±0.2	71.0±0.3	69.3±0.3	87.7
Ours(Jigsaw)	86.9±0.8	98.6±0.5	100.0±0.0	82.9±1.0	62.9±1.2	61.2±0.7	82.1
Ours(Rot)	90.1±0.8	98.1±0.3	100.0±0.0	88.6±0.7	65.1±0.8	65.0±0.6	84.5
Ours(MixRot)	88.6±0.7	98.0±0.3	99.9±0.1	86.1±0.3	65.7±1.0	65.0±0.7	83.9
Ours(SPRot)	87.3±0.6	98.5±0.7	100.0±0.0	85.1±0.7	65.0±0.7	62.8±0.8	83.1
Ours(Adv)	80.4±1.1	98.4±0.4	100.0±0.0	80.1±1.9	65.1±0.8	64.9±0.5	81.8
Ours(Rot+Adv)	88.4±0.7	97.9±0.3	100.0±0.0	85.5±0.7	67.6±0.6	65.4±0.6	84.1
Ours(Rot+Adv+BN)	88.6±0.7	98.0±0.3	100.0±0.0	85.7±0.3	68.0±0.6	65.5±0.3	84.3

JiGen [2].² We use the default hyper-parameters and ResNet-18 and ResNet-50 architectures. The deep networks used in our semantic segmentation experiments are ResNet-101 based DeepLab-v2 [7] and dilated residual networks (DRN) [64]. Specifically, we take the commonly used DRN-26 architecture in order to compare to other state-of-the-art methods. Both networks are initialized with ImageNet [11] pretrained weights.

B. DOMAIN ADAPTATION FOR OBJECT RECOGNITION

We first evaluate the proposed domain adaptation method with state-of-the-art methods on Office [48] dataset. Office is the most widely used dataset for visual domain adaptation, with 4652 images and 31 categories collected from three distinct domains: Amazon (**A**), Webcam (**W**) and DSLR (**D**). We evaluate all methods on six domain adaptation tasks, $A \rightarrow W$, $D \rightarrow W$, $A \rightarrow D$, $D \rightarrow A$ and $W \rightarrow A$. We follow the standard protocols for unsupervised domain adaptation in [36], and use all labeled source examples and all unlabeled target examples. We compare the average classification accuracy based on three random experiments. The abbreviations of the evaluated strategies are listed in Table 1. The results on Office dataset based on ResNet-50 are reported in Table 2. **Rot** achieves the best average accuracy among all evaluated strategies and reaching the accuracies of state-of-the-art methods. **MixRot** and **SPRot** obtain similar accuracies which are very close to **Rot**.

MixRot even outperforms **Rot** on $D \rightarrow W$ task. Because **Adv** has very low accuracy comparing to **Rot**, **Rot+Adv** does not improve **Rot**. **Rot+Adv+BN** shows consistent improvements to **Rot+Adv**, but still can not outperform the simplest method **Rot**.

To better analysis the domain adaptation performance of **Rot**, we visualize by t-SNE [57] the learned deep features in FIGURE 4 on task $A \rightarrow W$. FIGURE 4 (a) and (b) show that categories are better discriminated by **Rot** than the non-adapted model ResNet-50. FIGURE 4 (c) and (d) show that the source and target domains are aligned much better by **Rot** than ResNet-50.

For object recognition, we also evaluate on the multiple source domain adaptation dataset PACS [29] dataset, which has 7 object categories and 4 domains (Photo, Art Paintings, Cartoon and Sketches). FIGURE 5 shows sample images from PACS dataset. We follow the same experimental settings as [2] and trained our model considering three domains as source datasets and the remaining one as target. Following [2], we also compare to the domain discovery method DDiscovery [38] and Dial [3]. We set three different random seeds and run each experiment three times. The final result is the average over the three repetitions. To make a fair comparison, we run jigsaw puzzle method with the same random seeds and denoted by Ours(jigsaw). The results are shown in Table 3. We obtained similar conclusions to the Office dataset experiment. Our image rotation based self-supervised domain adaptation **Rot** outperforms all baselines. **MixRot** outperforms **Rot** on adaptation to art

²<https://github.com/fmcarlucci/JigenDG>

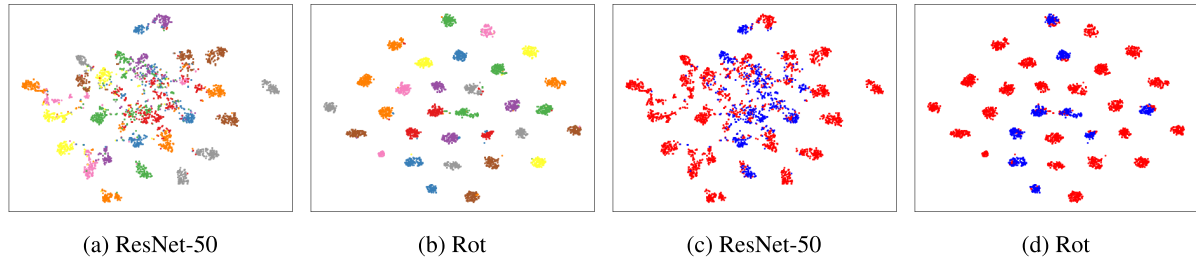


FIGURE 4. The t-SNE [57] visualization of deep features in $A \rightarrow W$ task. (a)(b) are generated from category information and each color in (a)(b) represents a category. (c)(d) are generated from domain information. Red and blue points represent samples of source and target domains, respectively.

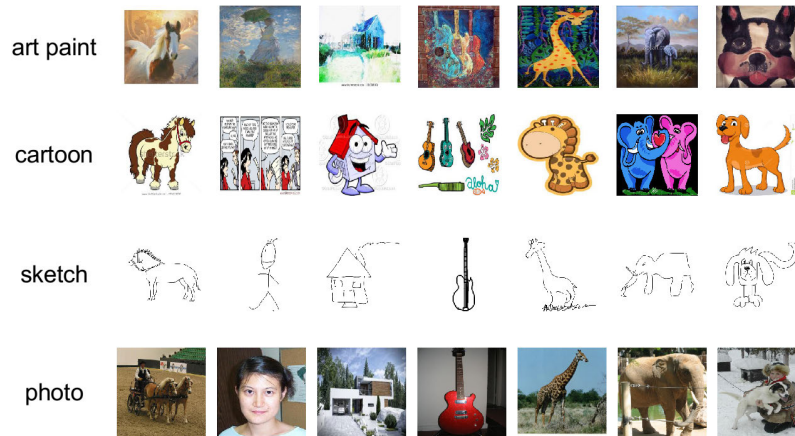


FIGURE 5. Sample images from PACS dataset. Each row represents a domain and each column represents a category.

TABLE 3. Multi-source Domain Adaptation results on PACS (ResNet-18). Three domains are used as source datasets and the remaining one as target.

Method	art paint.	cartoon	sketches	photo	Avg.
SRC [38]	74.70	72.40	60.10	92.90	75.03
Dial [3]	87.30	85.50	66.80	97.00	84.15
DDiscovery [38]	87.70	86.90	69.60	97.00	85.30
SRC [2]	77.85	74.86	67.74	95.73	79.05
JiGen [2]	84.88	81.07	79.05	97.96	85.74
Ours(SRC)	79.33	76.75	64.40	96.39	79.22
Ours(Jigsaw)	84.93	83.85	69.04	93.92	82.94
Ours(Rot)	88.67	86.39	74.93	98.02	87.00
Ours(MixRot)	89.35	84.14	74.49	98.24	86.56
Ours(SPRot)	86.57	87.95	67.06	98.26	84.96
Ours(Adv)	80.68	77.45	67.80	96.81	80.69
Ours(Rot+Adv)	86.47	86.18	70.53	97.90	85.27
Ours(Rot+Adv+BN)	86.87	86.52	71.59	98.01	85.75

painting and photo. **SPRot** outperforms **Rot** on adaptation to cartoon and photo. But their overall performance, *i.e.*, average accuracies are still lower than **Rot**, showing that **Rot** is the most robust method. Again, due to the relatively too low performance of **Adv**, **Rot+Adv** can not further improve **Rot**. As in the previous experiment, **Rot + Adv + BN** consistently improves **Rot + Adv**.

For domain adaptation on PACS, we also show t-SNE [57] visualization of deep features in FIGURE 6. From (a) and (b), we can see that **Rot** has much better

discriminativity on categories than non-adapted method **SRC**. From (c) and (d), **Rot** shows clearly much better domain alignment than non-adapted method **SRC**. The t-SNE visualization reveals the effectiveness of **Rot** on domain adaptation.

C. DOMAIN ADAPTATION FOR SEMANTIC SEGMENTATION

For semantic segmentation, we adapt semantic segmentation models from the source domain of synthetic images to the target domain of real-world images. For the synthetic datasets, we use SYNTHIA [46] and GTA5 [54], and for the target domain, we use the Cityscapes dataset [10]. The GTA5 [54] dataset is rendered from the Grand Theft Auto V video game. It consists of 24996 images with resolution of 1914×1052 and has 19 classes compatible with Cityscapes dataset. We use the full set of GTA5 as our source domain training set. For SYNTHIA dataset, we use the SYNTHIA-RAND-CITYSCAPES set [46] as the source domain training set, which contains 9400 images. We evaluate with the 16 common classes for SYNTHIA to Cityscapes domain adaptation. The training set of Cityscapes has 2975 images which are used as unlabeled target domain training samples. The validation set of Cityscapes has 500 samples which are used as our testing set.

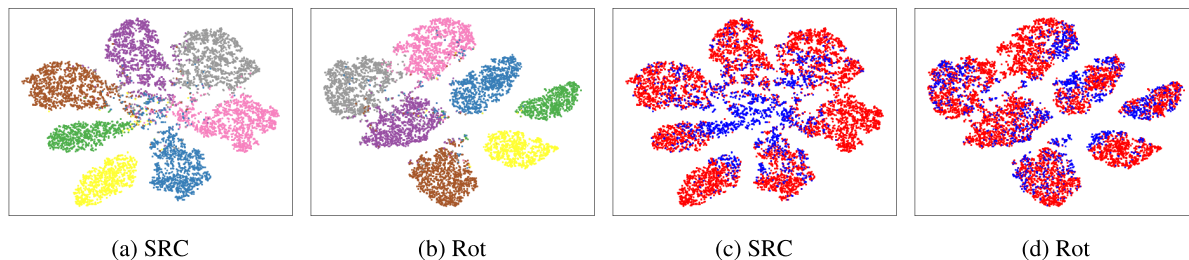


FIGURE 6. The t-SNE [57] visualization of deep features in multi-source DA task (art painting is used as target domain). (a)(b) are generated from category information and each color in (a)(b) represents a category. (c)(d) are generated from domain information. Red and blue points represent samples of source and target domains, respectively.

TABLE 5. Comparison with jigsaw puzzle method (ResNet-101).

Method	SYN2CS			GTA2CS		
	SRC	Adapt	Gain	SRC	Adapt	Gain
Jigsaw puzzle	30.5	34.3	3.8	35.0	36.2	1.2
Ours(Rot)	30.5	36.1	5.6	35.0	41.2	6.2

We conduct ablation studies to understand the impact of each component of our self-supervised domain adaptation. If not otherwise specified, all the experiments in this section use ResNet-101 as backbone network and the domain adaptation is from GTA5 to Cityscapes.

1) PRETEXT TASK LEARNING STRATEGIES

The first two rows in Table 4 show their domain adaptation results. As can be seen from Table 4, mixing source domain training data in the pretext learning (**MixRot**) shows even inferior results, which may be because the source samples are dominated in the mixed samples (24996 vs 2975), which makes the model more source domain oriented and reduces domain invariant representation power.

Next, we would like to know whether the proposed spatial-aware rotation prediction pretext task is better than the simple rotation prediction strategy, *i.e.*, the **Rot** method. Table 4 displays the results of the spatial-aware rotation prediction pretext task as **SPRot**. It turns out that the more difficult pretext task learning leads to worse domain adaptation performance. In our practice, the pretext task learning of **SPRot** has more difficulties to converge than **Rot**, and this may result in the failure of learning good feature representations. Therefore, how to design a proper pretext task for domain adaptation still needs more exploration.

We also compare our method **Rot** to the jigsaw puzzle based self-supervision [2]. The results are shown in Table 5, where SYN2CS denotes SYNTHIA to Cityscapes domain adaptation and GTA2CS for GTA5 to Cityscapes. **Rot** outperforms the jigsaw puzzle for both SYN2CS and GTS2CS. Especially for GTA2CS, jigsaw puzzle has shown very limited gain (1.2 percentage point) while **Rot** still achieved 6.2 percentage point.

2) INPUT IMAGE SIZE FOR PRETEXT TASK LEARNING

As the images from Cityscapes dataset have large resolution (*e.g.*, 1024×2048). We are interested in what cropping size is

best for the self-supervised learning. In Table 4, we compare three different cropping sizes. The smallest cropping size (128×128) shows worst performance due to too small field of view to learn good representations. Comparing the remaining two cropping sizes, we see that the larger one (400×400) does not further improve the performance. In fact, when we use the full image as input, the pretext learning easily gets stuck in a trivial solution, *i.e.* 100% prediction accuracy. As a result, the final model fails to perform domain adaptation. Thus, we believe that a proper cropping size is important to control the difficulty of learning pretext tasks.

3) FEATURE EXTRACTION LAYER

By default, the pretext task takes as input the features extracted from the last layer of the encoder. However, whether the last layer is the best for domain adaptation is unclear. In this section, we train self-supervised domain adaption models with different feature extraction layers. We mainly compare the feature extraction from the middle and the end of the encoder. Table 4 shows the corresponding results, where **Middle** represents the feature extraction from middle layer and **Final** uses features from the end layer of the encoder. As, in this case, the decoder of the segmentation network is simply an up-sampling layer without any learnable parameter, the **Final** layer is actually the prediction layer of the segmentation network. As we can see from the results, the model **Middle** shows slightly better results and we think the pretext task learning is not very sensitive to the choice of feature extraction layers.

4) EVALUATION OF COMPLEMENTARY STRATEGIES

Table 6 shows the results with different complementary strategies. The source domain model is denoted by **SRC** and the model trained with target domain samples is denoted by **TAR**, which represent the lower and upper bound of the accuracy respectively. **Rot** is our baseline method. **+Adv** is with prediction layer alignment (Section III-B1), which improves **Rot** by 1.1 percentage points. **+BN** is with BN calibration (Section III-B2), which does not show improvement over **Rot**. But when combined **Adv** and **BN**, we obtain the best results, improving **Rot** by 2.1 percentage points. Table 7 shows more results with other architectures and

TABLE 4. Domain adaptation performance under different pretext task settings (ResNet-101).

		GTA5 → Cityscapes																		
	mIoU	road	sidewalk	building	wall	fence	pole	light	sign	veg	terrain	sky	person	rider	car	truck	bus	train	mbike	bike
Rot	41.2	87.6	25.7	77.5	19.8	16.8	29.0	32.1	20.5	79.9	32.9	75.3	58.2	26.0	79.0	23.3	31.6	2.1	26.9	37.7
MixRot	40.3	79.4	13.7	77.6	20.1	19.4	27.9	36.3	30.6	83.3	29.2	74.4	60.2	29.2	64.9	27.8	18.1	0.3	28.4	44.2
SPRot	37.7	80.6	19.2	76.1	17.8	16.0	29.4	32.9	20.2	77.5	19.6	74.2	59.0	28.1	67.8	31.2	12.4	0.4	26.8	25.9
128x128	38.6	77.8	13.5	78.9	18.6	19.1	25.8	34.3	28.8	77.4	19.2	72.1	60.2	27.3	67.0	31.5	8.3	0.7	32.1	40.7
256x256	41.2	87.6	25.7	77.5	19.8	16.8	29.0	32.1	20.5	79.9	32.9	75.3	58.2	26.0	79.0	23.3	31.6	2.1	26.9	37.7
400x400	40.3	79.4	13.7	77.6	20.1	19.4	27.9	36.3	30.6	83.3	29.2	74.4	60.2	29.2	64.9	27.8	18.1	0.3	28.4	44.2
Middle	41.2	87.6	25.7	77.5	19.8	16.8	29.0	32.1	20.5	79.9	32.9	75.3	58.2	26.0	79.0	23.3	31.6	2.1	26.9	37.7
Final	40.4	84.1	25.0	79.5	15.5	15.5	29.5	30.5	27.8	82.1	21.7	80.3	54.3	26.0	70.1	29.5	29.2	0.2	26.3	40.9

TABLE 6. Evaluation of complementary strategies (ResNet-101).

		GTA5 → Cityscapes																		
	mIoU	road	sidewalk	building	wall	fence	pole	light	sign	veg	terrain	sky	person	rider	car	truck	bus	train	mbike	bike
Rot	41.2	87.6	25.7	77.5	19.8	16.8	29.0	32.1	20.5	79.9	32.9	75.3	58.2	26.0	79.0	23.3	31.6	2.1	26.9	37.7
+Adv	42.3	84.9	31.9	80.4	19.0	21.7	28.2	34.7	27.7	82.8	26.5	72.7	58.2	25.3	82.1	18.7	42.1	1.2	26.0	39.7
+BN	41.0	86.7	32.6	78.7	20.4	20.6	27.0	28.6	15.8	82.4	38.0	74.6	57.6	24.0	80.1	23.6	29.3	0.8	23.2	34.8
+Adv+BN	43.3	87.3	35.0	80.0	20.2	21.8	28.7	32.3	25.8	83.3	29.3	73.7	58.7	25.7	83.2	27.5	43.6	1.8	27.5	38.1
SRC	35.0	77.5	12.3	71.3	8.1	18.8	26.6	32.4	19.6	73.7	11.3	67.9	55.2	24.3	73.3	16.9	9.9	0.9	26.4	39.5
TAR	65.3	96.5	74.3	88.0	48.8	41.2	42.3	47.0	60.7	88.5	52.8	90.5	68.6	48.9	91.1	68.5	69.5	46.3	51.3	65.0

TABLE 7. Comparison with the state-of-the-art methods. The result of AdaptSegNet [55] here is from the single resolution version as our output adversarial method is built on top of this version.

Method	Network	SYN2CS			GTA2CS			Mechanism
		SRC	Adapt	Gain	SRC	Adapt	Gain	
CyCADA [20]	DRN-26	-	-	-	21.7	39.5	17.8	Input adversary
Stylization [13]	DRN-26	22.0	35.0	13.0	22.9	38.3	15.4	Input stylization
DCAN [61]	ResNet-101	28.0	36.5	8.5	29.8	38.5	8.7	Input stylization
FCAN [68]	ResNet-101	-	-	-	29.2	46.6	17.4	Input stylization + feature adversary
CrDoCo [9]	DRN-26	22.9	33.4	10.5	22.9	45.1	22.2	Input adversary + feature adversary
DISE [5]	ResNet-101	-	41.5	-	-	45.4	-	Input adversary + output adversary
BDL [31]	ResNet-101	-	51.4	-	33.6	48.5	14.9	Input adversary + self-labelling
ADR [49]	ResNet-50	-	-	-	25.3	33.3	8.0	Feature adversary
GAM [22]	DRN-26	-	-	-	-	40.2	-	Feature adversary
AdaptSegNet [55]	ResNet-101	-	-	-	36.6	41.4	4.8	Output adversary
ADVENT [59]	ResNet-101	-	41.2	-	-	45.5	-	Output adversary
CLAN [37]	ResNet-101	38.6	47.8	9.2	36.6	43.2	6.6	Output adversary
CBST [69]	ResNet-38	29.2	42.5	13.3	35.4	47.0	11.6	Self-labelling
CURC [67]	DRN-26	21.9	28.2	6.3	-	-	-	Curriculum
Ours(Rot)	ResNet-101	30.5	36.1	5.6	35.0	41.2	6.2	Self-supervision
	DRN-26	25.1	28.9	3.8	29.4	34.6	5.2	
Ours(Adv)	ResNet-101	30.5	35.4	4.9	35.0	40.4	5.4	Output adversary
	DRN-26	25.1	29.9	4.8	29.4	33.1	3.7	
Ours(Rot+Adv)	ResNet-101	30.5	39.3	8.8	35.0	42.3	7.3	Self-supervision + output adversary
	DRN-26	25.1	31.7	6.6	29.4	36.1	6.7	
Ours(Rot+Adv+BN)	ResNet-101	30.5	38.8	8.3	35.0	43.3	8.3	Self-supervision + output adversary + BN
	DRN-26	25.1	30.4	5.3	29.4	36.2	6.8	

datasets, where +**Adv** has consistent improvements to **Rot** but +**Adv+BN** gets saturated for the SYN2CS problem.

To understand why BN does not have consistent improvements, we further conduct experiments using only BN calibration for domain adaptation. FIGURE 7 shows the results on multiple datasets using multiple networks. BN calibration alone achieves surprisingly good results, and the best domain adaptation gain even reaches 6.8 percentage points. However, when combined with **Rot** or **Rot+Adv**, it only improves 1 or 2 percentage points. This might be because **Rot** and **Rot+Adv** have already learned domain invariant

representation that effectively reduces the covariate shift and BN calibration could not contribute more to the adapted model. The reason that **Adv** gives consistent rise to the base method is because **Adv** further aligns the predicted label distributions which is more complementary adaptation to the **Rot** than the provided by **BN**.

5) QUALITATIVE ANALYSIS

Following [37], we also visualize the learned feature representations by t-SNE [57] in FIGURE 8. For the non-adapted features, the classes are not discriminated well. They are

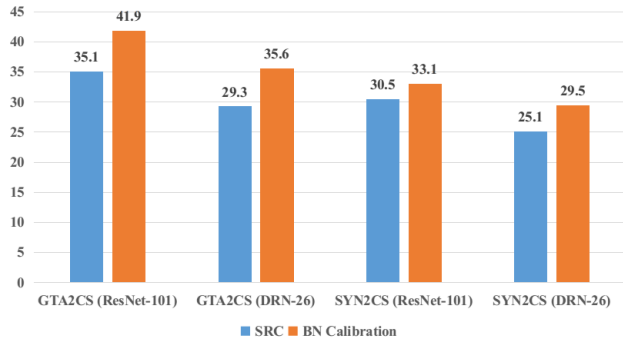


FIGURE 7. Domain adaptation performance of BN calibration. The vertical axis denotes mIoU accuracy.

discriminated better by the **Rot** adaptation and **Rot+Adv+BN** discriminates the classes best. In FIGURE 9, we illustrate some qualitative results of our models. Without domain adaptation, the source domain model **SRC** produces noisy segmentation. **Rot** shows significant improvements over **SRC** in terms of segmentation quality. The results of **Rot+Adv+BN** is less noisy and more accurate in details than **Rot**.

D. COMPARISON TO THE STATE-OF-THE-ART

Lastly, we compare our method to some recently published state-of-the-art works which use similar architectures to ours. The results are shown in Table 7. The compared methods cover large varieties of domain adaptation mechanisms, including input/feature/output level alignment methods, curriculum and self-labeling based methods. Some of these methods are also surveyed in [66]. We refer the readers to [66] for more details. The results in Table 7 show that our adapted models (**Adapt**) achieve comparable accuracies to the state-of-the-art. It is worth noting some of these state-of-the-art methods obtain worse results than we obtain when training with the source data alone (SRC columns), so their relative gain is higher. On the other hand, with this work we aim at encouraging the use of pretext tasks for domain adaption of

semantic segmentation models, which, as mentioned before, can be a complementary idea to others. We also find that a deeper network (ResNet-101) can achieve better domain adaptation gain than the shallow one (DRN-26).

E. DISCUSSION

The experimental results reveal several insightful observations. (1) The current deep learning methods learn good feature representations for single domain but can not remove cross-domain discrepancy. (2) Using self-supervised representation learning can help to reduce domain shift, and the simplest image rotation prediction pretext task **Rot** can even achieve comparable performance to the state-of-the-art domain adaptation methods. (3) **Rot** turns out to be more robust than other alternatives, e.g., **MixRot**, **SPRot** and **Jigsaw**. (4) **Rot** is complementary to existing domain adaptation methods, e.g., adversarial based and batch normalization based ones.

The reasons that self-supervised learning helps domain adaptation are as following: (1) the self-supervised learning involves source and target domain samples in a common supervised learning process which can help to learn cross-domain feature representations. This can be verified from the feature visualization on domain alignment, e.g., FIGURE 4 and FIGURE 6. (2) Because the self-supervised learning and the main task are in a joint multi-task learning process, the model of the main task also learns from the cross-domain feature representations. As a result, the final model achieves domain adaptation on target domain.

Based on our experiments, we also have following findings about how to design a good self-supervised pretext task for domain adaptation: (1) As a common practice on many deep learning tasks, a deeper architecture can achieve better self-supervised domain adaptation performance than a shallow one. (2) Better performance on representation learning, better performance on domain adaptation, e.g., **Rot** vs **Jigsaw**. (3) More complex pretext task does not lead to better

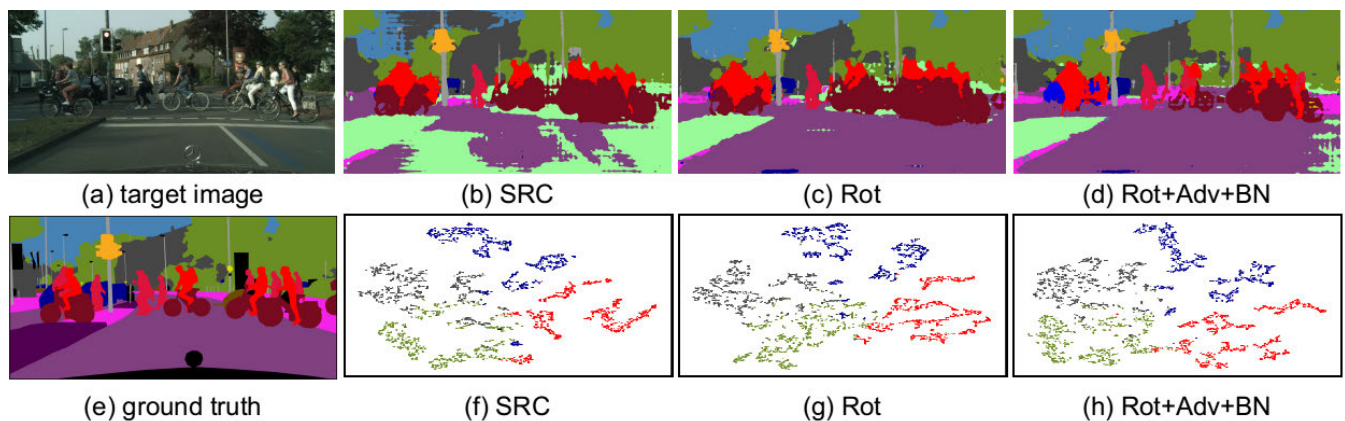


FIGURE 8. The t-SNE [57] visualization of the learned features. (a): A target domain image. (b): A non-adapted results. (c) Adapted result of Rot. (d) Adapted results of Rot+Adv+BN. (e) ground truth segmentation. We map the high-dimensional features of (b), (c) and (d) to a 2-D space with t-SNE [57] shown in (f), (g) and (h). For clear illustration, we visualize features of 4 classes, including building (grey), car (blue), vegetation (green) and rider (red).

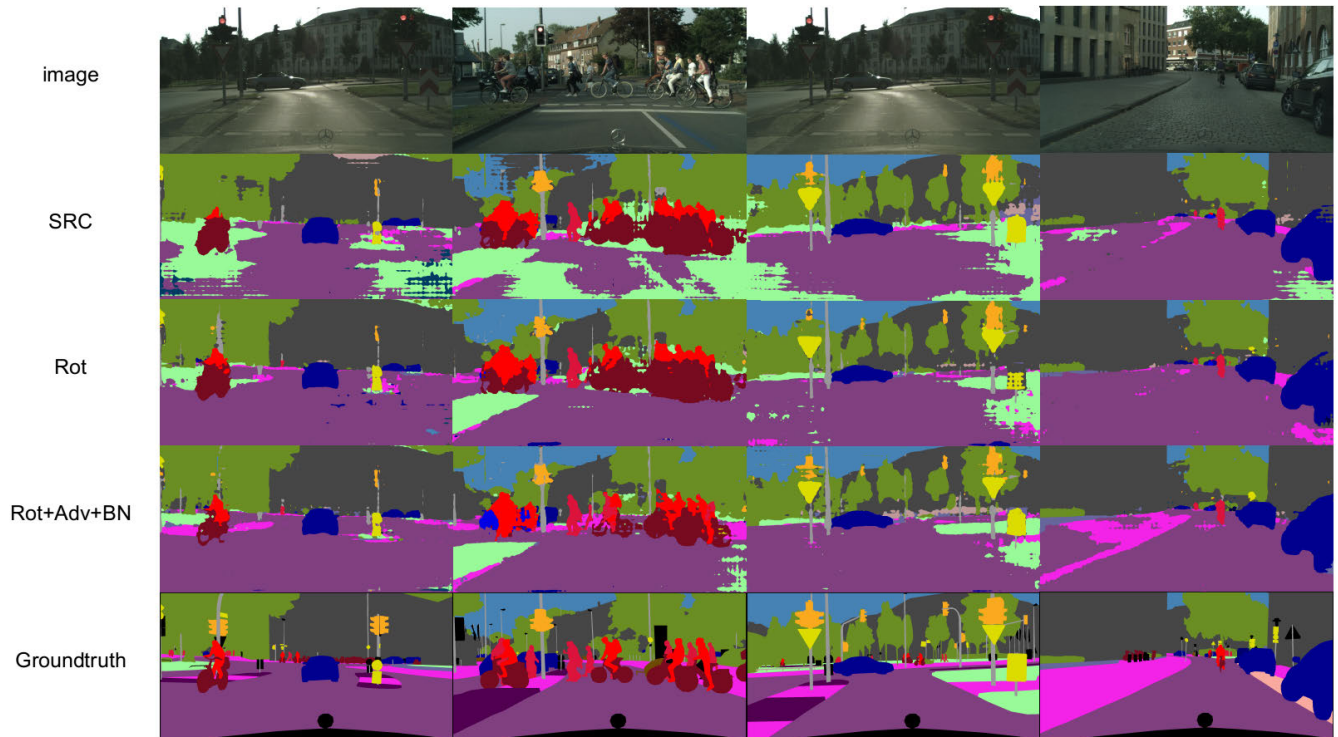


FIGURE 9. Qualitative results of GTA2CS domain adaptation. The first row and last row are the input images and corresponding groundtruths respectively. The second row shows results from source domain model. The third and fourth rows are results from the proposed Rot and Rot+Adv+BN models. ResNet-101 is used as backbone.

domain adaptation performance, *e.g.*, **SPRot** is outperformed by **Rot**. (4) **Adv** and **Adv+BN** can further improve **Rot** if the performance of **Adv** is not worse than **Rot**. In this work, we only investigated several simple self-supervised learning strategies, we believe that for better self-supervised domain adaptation there are still large space to explore.

V. CONCLUSION

In this work, we have explored self-supervised learning for domain adaptation. We have shown that a simple image rotation prediction (pretext task) self-supervision can achieve state-of-the-art domain adaptation performance. We have studied several pretext tasks as well as complementary domain adaptation strategies. Taking object recognition and semantic segmentation of urban scenes as relevant use cases, we have performed an ablative analysis of the different components included in our overall domain adaptation procedure. As future work, we would like to investigate more pretext tasks and to apply our method to other relevant vision tasks.

ACKNOWLEDGMENT

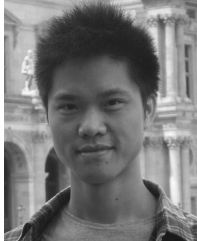
As CVC member, Antonio thanks the Generalitat de Catalunya CERCA Program and its ACCIO agency.

REFERENCES

- [1] Y. Abramson and Y. Freund, "Active learning for visual object recognition," Univ. California San Diego, San Diego, CA, USA, UCSD Tech. Rep., 2005. [Online]. Available: <http://cseweb.ucsd.edu/~yfreund/papers/index.html>
- [2] F. M. Carlucci, A. D'Innocente, S. Bucci, B. Caputo, and T. Tommasi, "Domain generalization by solving jigsaw puzzles," in *Proc. CVPR*, Jun. 2019, pp. 2229–2238.
- [3] F. M. Carlucci, L. Porzi, B. Caputo, E. Ricci, and S. R. Bulò, "Just dial: Domain alignment layers for unsupervised domain adaptation," in *Proc. Int. Conf. Image Anal. Process.*, 2017, pp. 357–369.
- [4] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Proc. ECCV*, 2018, pp. 132–149.
- [5] W.-L. Chang, H.-P. Wang, W.-H. Peng, and W.-C. Chiu, "All about structure: Adapting structural information across domains for boosting semantic segmentation," in *Proc. CVPR*, Jun. 2019, pp. 1900–1909.
- [6] C. Chen, Z. Chen, B. Jiang, and X. Jin, "Joint domain alignment and discriminative feature learning for unsupervised deep domain adaptation," in *Proc. AAAI*, 2019, pp. 1–9.
- [7] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, Apr. 2017.
- [8] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool, "Domain adaptive faster R-CNN for object detection in the wild," in *Proc. CVPR*, Jun. 2018, pp. 3339–3348.
- [9] Y.-C. Chen, Y.-Y. Lin, M.-H. Yang, and J.-B. Huang, "CrDoCo: Pixel-level domain transfer with cross-domain consistency," in *Proc. CVPR*, Jun. 2019, pp. 1791–1800.
- [10] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proc. CVPR*, Jun. 2016, pp. 3213–3223.
- [11] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. CVPR*, Miami, FL, USA, 2009, pp. 1–8.
- [12] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *Proc. ICCV*, Dec. 2015, pp. 1422–1430.

- [13] A. Dundar, M.-Y. Liu, T.-C. Wang, J. Zedlewski, and J. Kautz, "Domain stylization: A strong, simple baseline for synthetic to real image domain adaptation," 2018, *arXiv:1807.09384*. [Online]. Available: <https://arxiv.org/abs/1807.09384>
- [14] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL visual object classes (VOC) challenge," *Int. J. Comput. Vis.*, vol. 88, no. 2, pp. 303–338, 2010.
- [15] Z. Feng, C. Xu, and D. Tao, "Self-supervised representation learning by rotation feature decoupling," in *Proc. CVPR*, Jun. 2019, pp. 10364–10374.
- [16] Y. Gani, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," in *Proc. ICML*, 2015, pp. 1–35.
- [17] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *Proc. ICLR*, 2018, pp. 1–16.
- [18] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. CVPR*, Columbus, OH, USA, Jun. 2014, pp. 580–587.
- [19] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, Jun. 2016, pp. 770–778.
- [20] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell, "CYCADA: Cycle-consistent adversarial domain adaptation," in *Proc. ICML*, 2018, pp. 1–15.
- [21] J. Hoffman, D. Wang, F. Yu, and T. Darrell, "FCNs in the wild: Pixel-level adversarial and constraint-based adaptation," in *Proc. CVPR*, Dec. 2016, pp. 1–9.
- [22] H. Huang, Q. Huang, and P. Krähenbühl, "Domain transfer through deep activation matching," in *Proc. ECCV*, 2018, pp. 590–605.
- [23] H. Jiang, G. Larsson, M. M. G. Shakhnarovich, and E. Learned-Miller, "Self-supervised relative depth learning for urban scene understanding," in *Proc. ECCV*, 2018, pp. 19–35.
- [24] L. Jing and Y. Tian, "Self-supervised visual feature learning with deep neural networks: A survey," 2019, *arXiv:1902.06162*. [Online]. Available: <https://arxiv.org/abs/1902.06162>
- [25] D. Kim, D. Cho, D. Yoo, and I. S. Kweon, "Learning image representations by completing damaged jigsaw puzzles," in *Proc. WACV*, Mar. 2018, pp. 793–802.
- [26] A. Kolesnikov, X. Zhai, and L. Beyer, "Revisiting self-supervised visual representation learning," 2019, *arXiv:1901.09005*. [Online]. Available: <https://arxiv.org/abs/1901.09005>
- [27] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. NIPS*, Lake Tahoe, NV, USA, 2012, pp. 1097–1105.
- [28] H. Kuehne, H. Jhuang, E. Garrote-Contreras, T. Poggio, and T. Serre, "HMDB: A large video database for human motion recognition," in *Proc. ICCV*, 2011, pp. 2556–2563.
- [29] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales, "Deeper, broader and artier domain generalization," in *Proc. ICCV*, Oct. 2017, pp. 5542–5550.
- [30] Y. Li, N. Wang, J. Shi, J. Liu, and X. Hou, "Revisiting batch normalization for practical domain adaptation," in *Proc. PR*, 2018, pp. 1–12.
- [31] Y. Li, L. Yuan, and N. Vasconcelos, "Bidirectional learning for domain adaptation of semantic segmentation," in *Proc. CVPR*, Jun. 2019, pp. 6936–6945.
- [32] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, and P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *Proc. ECCV*, 2014, pp. 740–755.
- [33] M. Long, Y. Cao, J. Wang, and M. I. Jordan, "Learning transferable features with deep adaptation networks," in *Proc. ICML*, 2015, pp. 1–9.
- [34] M. Long, Z. Cao, J. Wang, and M. I. Jordan, "Conditional adversarial domain adaptation," in *Proc. NIPS*, 2018, pp. 1640–1650.
- [35] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Unsupervised domain adaptation with residual transfer networks," in *Proc. NIPS*, 2016, pp. 136–144.
- [36] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *Proc. ICML*, 2017, pp. 2208–2217.
- [37] Y. Luo, L. Zheng, T. Guan, J. Yu, and Y. Yang, "Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 2507–2516.
- [38] M. Mancini, L. Porzi, S. R. Bulò, B. Caputo, and E. Ricci, "Boosting domain adaptation by discovering latent domains," in *Proc. CVPR*, Jun. 2018, pp. 3771–3780.
- [39] J. Manders, T. van Laarhoven, and E. Marchiori, "Adversarial alignment of class prediction uncertainties for domain adaptation," in *Proc. ICPRAM*, 2019, pp. 1–11.
- [40] T. N. Mundhenk, D. Ho, and B. Y. Chen, "Improvements to context based self-supervised learning," in *Proc. CVPR*, Jun. 2018, pp. 9339–9348.
- [41] G. Neuhold, T. Ollmann, S. R. Bulò, and P. Kotschieder, "The mapillary vistas dataset for semantic understanding of street scenes," in *Proc. ICCV*, Oct. 2017, pp. 4990–4999.
- [42] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *Proc. ECCV*, 2016, pp. 69–84.
- [43] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 10, pp. 1345–1359, Oct. 2010.
- [44] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proc. CVPR*, 2016, pp. 2536–2544.
- [45] Z. Pei, Z. Cao, M. Long, and J. Wang, "Multi-adversarial domain adaptation," in *Proc. AAAI*, 2018, pp. 1–8.
- [46] G. Ros, L. Sellart, J. Materzynska, D. Vázquez, and A. López, "The SYNTHIA dataset: A large collection of synthetic images for semantic segmentation of urban scenes," in *Proc. CVPR*, Jun. 2016, pp. 3234–3243.
- [47] S. Roy, A. Unmesh, and V. P. Nambodiri, "Deep active learning for object detection," in *Proc. BMVC*, 2018, pp. 1–91.
- [48] K. Saenko, B. Hulis, M. Fritz, and T. Darrell, "Adapting visual category models to new domains," in *Proc. ECCV*, 2010, pp. 213–226.
- [49] K. Saito, Y. Ushiku, T. Harada, and K. Saenko, "Adversarial dropout regularization," in *Proc. ICLR*, 2018, pp. 1–14.
- [50] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada, "Maximum classifier discrepancy for unsupervised domain adaptation," in *Proc. CVPR*, Jun. 2018, pp. 3723–3732.
- [51] S. Sankaranarayanan, Y. Balaji, C. D. Castillo, and R. Chellappa, "Generate to adapt: Aligning domains using generative adversarial networks," in *Proc. CVPR*, Jun. 2018, pp. 8503–8512.
- [52] B. Settles, "Active learning," *Synth. Lectures Artif. Intell. Mach. Learn.*, vol. 6, no. 1, pp. 1–114, 2012.
- [53] K. Soomro, A. Zamir, and M. Shah, "UCF101: A dataset of 101 human actions classes from videos in the wild," 2012, *arXiv:1212.0402*. [Online]. Available: <https://arxiv.org/abs/1212.0402>
- [54] S. R. Stephan, R. V. Vineet, and V. Koltun, "Playing for data: Ground truth from computer games," in *Proc. ECCV*, 2016, pp. 102–118.
- [55] Y.-H. Tsai, W.-C. Hung, S. Schuster, K. Sohn, M.-H. Yang, and M. Chandraker, "Learning to adapt structured output space for semantic segmentation," in *Proc. CVPR*, Jun. 2018, pp. 7472–7481.
- [56] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proc. CVPR*, Jul. 2017, pp. 7167–7176.
- [57] L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *J. Mach. Learn. Res.*, vol. 9, pp. 2579–2605, Nov. 2008.
- [58] D. Vázquez, A. M. López, J. Marín, D. Ponsa, and D. Gerónimo, "Virtual and real world adaptation for pedestrian detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 4, pp. 797–809, Aug. 2014.
- [59] T.-H. Vu, H. Jain, M. Bucher, M. Cord, and P. Pérez, "ADVENT: Adversarial entropy minimization for domain adaptation in semantic segmentation," in *Proc. CVPR*, Jun. 2019, pp. 2517–2526.
- [60] M. Wang and W. Deng, "Deep visual domain adaptation: A survey," *Neurocomputing*, vol. 312, pp. 135–153, Oct. 2018.
- [61] Z. Wu, X. Han, Y.-L. Lin, M. G. Uzunbas, T. Goldstein, S. N. Lim, and L. S. Davis, "DCAN: Dual channel-wise alignment networks for unsupervised scene adaptation," in *Proc. ECCV*, 2018, pp. 518–534.
- [62] J. Xu, S. Ramos, D. Vázquez, and A. López, "Domain adaptation of deformable part-based models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 12, pp. 2367–2380, Jun. 2014.
- [63] J. Xu, D. Vázquez, K. Mikolajczyk, and A. López, "Hierarchical online domain adaptation of deformable part-based models," in *Proc. ICRA*, May 2016, pp. 5536–5541.
- [64] F. Yu, V. Koltun, and T. Funkhouser, "Dilated residual networks," in *Proc. CVPR*, Jun. 2017, pp. 472–480.
- [65] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *Proc. ECCV*, 2016, pp. 649–666.
- [66] Y. Zhang, P. David, H. Foroosh, and B. Gong, "A curriculum domain adaptation approach to the semantic segmentation of urban scenes," 2018, *arXiv:1812.09953*. [Online]. Available: <https://arxiv.org/abs/1812.09953>
- [67] Y. Zhang, P. David, and B. Gong, "Curriculum domain adaptation for semantic segmentation of urban scenes," in *Proc. ICCV*, Jul. 2017, pp. 2020–2030.

- [68] Y. Zhang, Z. Qiu, T. Yao, D. Liu, and T. Mei, "Fully convolutional adaptation networks for semantic segmentation," in *Proc. CVPR*, Jun. 2018, pp. 6810–6818.
- [69] Y. Zou, Z. Yu, B. V. K. Vijaya Kumar, and J. Wang, "Unsupervised domain adaptation for semantic segmentation via class-balanced self-training," in *Proc. ECCV*, 2018, pp. 289–305.



vision and machine learning.

JIAOLONG XU received the Ph.D. degree from the Computer Vision Center (CVC), Universitat Autònoma de Barcelona (UAB), in 2015, and the B.Sc. and M.Sc. degrees from the National University of Defense Technology, China, in 2008 and 2010, respectively. He is currently an Assistant Research Fellow with the Unmanned Systems Research Center, National Innovation Institute of Defense Technology (NIIDT), Beijing, China. His research interests include computer



LIANG XIAO received the Ph.D. degree in control science and engineering from the National University of Defense Technology (NUDT), Changsha, China, in 2017. He is currently an Assistant Research Fellow with the Unmanned Systems Research Center, National Innovation Institute of Defense Technology (NIIDT), Beijing, China. His research interests include computer vision, machine learning, and robotics.



ANTONIO M. LÓPEZ is currently the PI of the Autonomous Driving Lab, Computer Vision Center (CVC), UAB, where he is an Associated Professor with the Computer Science Department. He has a long research trajectory at the intersection of computer vision, computer graphics, machine learning, and autonomous driving. He has been deeply involved in the creation of the SYNTHIA dataset and the CARLA open-source simulator. He is actively working with industry partners to bring state-of-the-art techniques to the field of autonomous driving. He is granted by the Catalan ICREA Academia program.

...