

# Adaptation Across Extreme Variations using Unlabeled Bridges

Shuyang Dai<sup>1</sup>, Kihyuk Sohn<sup>2</sup>,

<sup>1</sup>Duke University

Yi-Hsuan Tsai<sup>2</sup>, Lawrence Carin<sup>1</sup>,

<sup>2</sup>NEC Labs America

Manmohan Chandraker<sup>2,3</sup>

<sup>3</sup>UC San Diego

---

## Abstract

We tackle an unsupervised domain adaptation problem for which the domain discrepancy between labeled source and unlabeled target domains is large, due to many factors of inter- and intra-domain variation. While deep domain adaptation methods have been realized by reducing the domain discrepancy, these are difficult to apply when domains are significantly different. We propose to decompose domain discrepancy into multiple but smaller, and thus easier to minimize, discrepancies by introducing unlabeled bridging domains that connect the source and target domains. We realize our proposed approach through an extension of the domain adversarial neural network with multiple discriminators, each of which accounts for reducing discrepancies between unlabeled (bridge, target) domains and a mix of all precedent domains including source. We validate the effectiveness of our method on several adaptation tasks including object recognition and semantic segmentation.

## 1 Introduction

With advances in supervised deep learning, many vision problems have realized significant performance improvements [6, 13, 19, 27, 33, 36, 37, 40]. While the success is driven by several factors, such as improved deep learning architectures [19, 23] or optimization techniques [11, 25, 26], it is strongly dependent on the existence of large-scale labeled training data [10]. Unfortunately, such a dataset may not be available for each application domain. This demands new ways of knowledge transfer from existing labeled data to individual target applications, potentially with access to large-scale *unlabeled* data from the application domain.

Unsupervised domain adaptation (UDA) [0, 5] has been proposed to improve the generalization ability of classifiers, using unlabeled data from the target domain. Deep domain adaptation that realizes UDA in a deep learning framework has been successful in several vision tasks [0, 12, 22, 24, 32, 35]. The core idea is to reduce the discrepancy metric between the two domains, measured by the domain discriminator [12] or MMD kernel [24] at certain representation of deep networks. Ideally, the discriminator learns the transformation mechanisms between the two domains. However, it could be difficult to model such dynamics when there are many factors of inter- and intra-domain variation applied to transform the source domain into the target domain.

In this paper, we aim to solve unsupervised domain adaptation challenges when domain discrepancy is large due to variation across the source and target domains. Figure 1 provides an

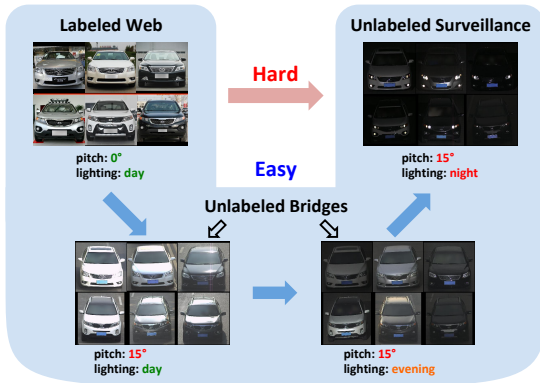


Figure 1: Unsupervised domain adaptation is challenging when the target domain is significantly different from the source domain due to many convoluted factors of variation. We introduce *bridging domains* composed of unlabeled images with some common factors to the source (e.g., lighting) and the target domain (e.g., viewpoint, image resolution).

As in Figure 1, the domain on the bottom left shares a consistent lighting condition (day) with the source, while the viewpoint is similar to that of the target domain. We note that there could be multiple bridging domains, such as the one on the bottom right of Figure 1, whose lighting intensity is between that of the first bridging domain and the target domain.

To utilize unlabeled bridging domains, we propose to extend the domain adversarial neural network [12] using multiple domain discriminators, each of which accounts for learning and reducing the discrepancy between unlabeled (bridging, target) domains and the mix of all precedent domains. We justify our learning framework by deriving a bound on the target error, that contains the source error and a list of discrepancies between unlabeled domain and the mix of precedent domains, including the source. This bound captures the intuition that judicious choices of bridge domains should not introduce large discrepancies. We hypothesize that the decomposition of a single, large discrepancy into multiple, small ones leads to a series of easier optimization problems, culminating in better alignment of source and target domains. We illustrate this intuition in Figure 2 on a variant of the two-moons dataset.

While works on unsupervised discovery of latent domains exist [14, 15, 17], it still remains a hard, unsolved problem. Firstly, we focus on the complementary and also unsolved problem of devising adversarial formulations that exploit given bridging domains. We observe that such domain information is often easily available in practice, for example, image meta-data such as timestamps, geo-tags and calibration parameters suffice to inform about illumination, weather or perspective. Moreover, we exploit different methods on measuring domain discrepancy [12, 18] or out-of-distribution (OOD) sample detection [20] to discover latent domains in an unsupervised manner, i.e., without domain information.

## 2 Related Work

**Unsupervised Domain Adaptation.** The proper reduction of discrepancy across domains [31] is a longstanding challenge. Specifically, an appropriate metric is required in order to measure the difference in between domains [2]. Recent works use kernel-based methods such as maximum mean discrepancy (MMD) [25] and optimal transport (OT) [9] to measure the domain difference in the feature space. Others adopt the idea of adversarial training [12, 33] which is inspired by the generative adversarial network (GAN) [16]. This training

illustrative example of adapting from labeled images of cars from the internet to recognize cars for surveillance applications at night. Two dominant factors, the perspective and illumination, make this a difficult adaptation task. As a step towards solving these problems, we introduce *unlabeled domain bridges* whose factors of variation are partially shared with the source domain, while the others are in common with the target domain. As in Figure 1, the domain on the bottom left shares a consistent lighting condition (day) with the source, while the viewpoint is similar to that of the target domain. We note that there could be multiple bridging domains,

procedure allows the feature representations to be indistinguishable between the source and target domain, aligning the two. One example of using adversarial training on UDA problems is the domain adversarial neural network (DANN) [12]. It trains a discriminator that distinguishes domains, while also learning a feature extractor to fool the discriminator by providing domain-invariant feature.

**Multiple Domains.** In [17], models are proposed for multiple-source UDA problems based on a domain adversarial learning. While the intuition is to utilize extra source domains that are available, the adaptation process is in practice favored toward the source domain that is closely related to the target domain [29]. Our method shares the similar high-level idea, in which relevant domains should guide the adaptation. In contrast, *unlabeled* bridging domains that share factors of variation with both source and target domains are utilized to guide the two domains, aligned with the bridging domain. Similar to our proposed approach, the benefit of having intermediate domains to guide transfer learning is shown in [14], but in the context of semi-supervised label propagation, requiring labeled data from the target domains.

### 3 Method

Our proposed domain adaptation framework is built atop DANN, utilizing *unlabeled bridging domains* to enhance the adaptation performance when the source and target domains are significantly different due to factors of variations.

**Notation.** Denote  $\mathcal{D}_S$  and  $\mathcal{D}_T$  as the source and target domains, respectively, from which data  $x$  are drawn. Output label  $y \sim \mathcal{Y}$  has  $N$  categories. The model contains: 1) a feature extractor  $f: \mathcal{D} \rightarrow \mathbb{R}^K$ , with parameter  $\theta_f$ , that maps  $x$  into a feature vector  $f(x)$ ; 2) the domain discriminator  $d: \mathbb{R}^K \rightarrow (0, 1)$ , with parameter  $\theta_d$ , that tells whether  $f(x)$  is from  $\mathcal{D}_S$  or  $\mathcal{D}_T$ ; and 3) the classifier  $C: \mathbb{R}^K \rightarrow \mathcal{Y}$ , with parameter  $\theta_C$ , that gives a predicted label  $\hat{y} = C(f(x))$ .

#### 3.1 Domain Adversarial Neural Network

The domain adversarial neural network transfers a classifier learned from the labeled source domain to the unlabeled target domain by learning domain-invariant features. It is realized by first learning the domain-related information and leveraging it with features extracted from the input. DANN uses a domain discriminator  $d$  to control the amount of domain-related information in the extracted feature. The discriminator is updated by maximizing the following:

$$\mathcal{L}_d = \mathbb{E}_{x \sim \mathcal{D}_S} \log d(f(x)) + \mathbb{E}_{x \sim \mathcal{D}_T} \log(1 - d(f(x))). \quad (1)$$

In comparison, the feature extractor  $f$  wants to confuse the discriminator  $d$  to remove any domain-specific information. Moreover, to make sure the extracted feature is task-related,  $f$  is trained to generate features that can be correctly classified by the classifier  $C$  trained by minimizing the following:

$$\mathcal{L}_C = \mathbb{E}_{(x,y) \sim \mathcal{D}_S \times \mathcal{Y}} [-y \log C(f(x))], \quad (2)$$

and a learning objective for feature extractor is as follows:

$$\min_{\theta_f, \theta_C} \mathcal{L}_C + \lambda \mathcal{L}_d. \quad (3)$$

While [12] introduces a gradient reversal layer to jointly train all parameters, we do alternating update of GANs [16] between  $\theta_d$  and  $\{\theta_f, \theta_C\}$  in our implementation.

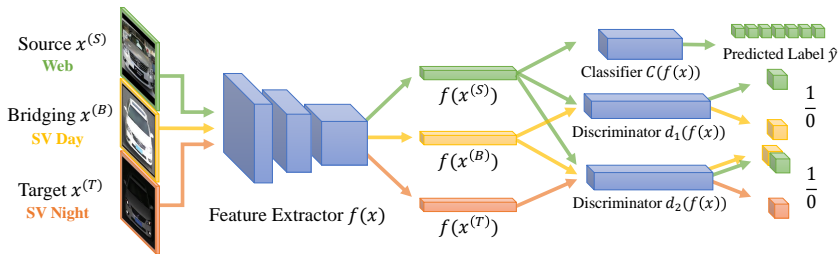


Figure 3: The learning framework with labeled source, unlabeled target, and unlabeled bridging domains for our extension of DANN using multiple discriminators. The model is composed of shared feature extractor  $f$ , classifier  $C$ , which is trained using labeled source examples, and two domain discriminators  $d_1$  and  $d_2$ .

### 3.2 Challenge in Domain Adversarial Learning

While deep domain adaptation algorithms are realized in different forms [6, 12, 68, 69, 44, 45], their theoretical motivation largely derives from the seminal work of [4]. In short, a theorem from that work states that the target domain task error  $\epsilon_T$  is bounded by the source error  $\hat{\epsilon}_S$  and the domain discrepancy:

$$\epsilon_T(h) \leq \hat{\epsilon}_S(h) + d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_T), \quad (4)$$

where  $h \in \mathcal{H}$  is a hypothesis and  $d_{\mathcal{H}\Delta\mathcal{H}}$  is written as:

$$\sup_{h, h' \in \mathcal{H}} |P_{\mathcal{D}_S}(h(x) \neq h'(x)) - P_{\mathcal{D}_T}(h(x) \neq h'(x))|.$$

Adversarial loss can be used to minimize the domain discrepancy to obtain a tighter bound. While it provides flexibility on the types of discrepancy, it is challenging to learn the right transformation from the source domain to the target domain when the two are far apart.

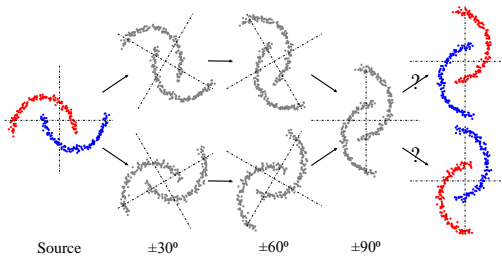


Figure 2: Translated two moons. The inter-twining moons (left) are considered as the source domain. Two moons are rotated and translated to the right to generate the target domain (right in gray).

Figure 2). In such a case, knowing what happens in the middle of the entire transformation process from source to target domains is critical, as these data points in the middle, even if they are unlabeled, can guide learning algorithms to easily disentangle transformation factors (e.g., clockwise rotation and translation to the right) from task-relevant factors.

### 3.3 Adaptation with Bridging Domain

We introduce additional sets of unlabeled examples, which we call bridging domains, that reside in the transformation pathway from labeled source to unlabeled target domains.

As motivation, consider a variant of the two-moons dataset, whose data points are translated to the right by the amount proportional to the rotation angle, as in Figure 2. The source domain is centered at the origin, while the target domain is moved to the right after being rotated by  $90^\circ$ , and given without labels. Adapting from source to target directly is difficult due to a significant change. Moreover, there are many ways to generate the same unlabeled target data points (e.g., rotate counterclockwise instead of clockwise, as in the bottom of

**DANN with a Single Bridging Domain.** Besides  $\mathcal{D}_S$  and  $\mathcal{D}_T$ , we denote  $\mathcal{D}_B$  as a bridging domain. Our framework is composed of feature extractor  $f(x)$  from an input  $x \in \mathcal{D}_S \cup \mathcal{D}_B \cup \mathcal{D}_T$  and classifier  $C(f(x))$  trained using classification loss in (2). Unlike DANN, which directly aligns  $\mathcal{D}_S$  and  $\mathcal{D}_T$ , we decompose the adaptation into two steps. First,  $\mathcal{D}_S$  and  $\mathcal{D}_B$  are aligned. This is an easier task than direct adaptation as in DANN, since there are less discriminating factors between  $\mathcal{D}_S$  and  $\mathcal{D}_B$ . Second, we adapt  $\mathcal{D}_T$  to the union of  $\mathcal{D}_S$  and  $\mathcal{D}_B$ . Similarly, the task is easier since it needs to discover remaining factors between  $\mathcal{D}_T$  and  $\mathcal{D}_S$  or  $\mathcal{D}_B$ , as some factors are already found from the previous step. To accommodate the two adaptation steps, we use two binary domain discriminators,  $d_1$  for learning discrepancy between  $\mathcal{D}_S$  and  $\mathcal{D}_B$ , and  $d_2$  between  $\mathcal{D}_S \cup \mathcal{D}_B$  and  $\mathcal{D}_T$ . Finally, this is realized with the following objectives:

$$\mathcal{L}_{d_1} = \mathbb{E}_{\mathcal{D}_S} \log d_1(f) + \mathbb{E}_{\mathcal{D}_B} \log(1 - d_1(f)), \quad (5)$$

$$\mathcal{L}_{d_2} = \mathbb{E}_{\mathcal{D}_S \cup \mathcal{D}_B} \log d_2(f) + \mathbb{E}_{\mathcal{D}_T} \log(1 - d_2(f)). \quad (6)$$

Both  $\mathcal{L}_{d_1}$  and  $\mathcal{L}_{d_2}$  are minimized to update their respective model parameters  $\theta_{d_1}$  and  $\theta_{d_2}$ . We update the classifier using (2) and the feature extractor to confuse discriminators as follows:

$$\min_{\theta_f, \theta_C} \mathcal{L}_C + \lambda_1 \mathcal{L}_{d_1} + \lambda_2 \mathcal{L}_{d_2}, \quad (7)$$

with two hyperparameters  $\lambda_1$  and  $\lambda_2$  to adjust the strengths of adversarial loss. We alternate updates between  $d_1, d_2$  and  $f, C$ . The proposed framework is visualized in Figure 3.

**Theoretical Insights.** To provide insights on how our learning objectives are constructed, we derive a bound on the target error while considering the unlabeled bridging domain:

$$\varepsilon_T(h) \leq \varepsilon_T(h_T^*) + \frac{1}{2} \varepsilon_B(h_B^*) + 2\gamma\alpha + \eta + d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_\alpha, \mathcal{D}_T) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_S, \mathcal{D}_B), \quad (8)$$

where  $h^* = \arg \min_{h \in \mathcal{H}} \varepsilon(h)$ ,  $\mathcal{D}_\alpha = \mathcal{D}_S \cup \mathcal{D}_B$ , and

$$\gamma\alpha = \min_{h \in \mathcal{H}} \{ \varepsilon_T(h) + \sum_{j=1}^N \alpha_j \varepsilon_j(h) \} \text{ with } \alpha = (0.5, 0.5).$$

Note that  $\varepsilon_T(h_T^*) + \frac{1}{2} \varepsilon_B(h_B^*) + 2\gamma\alpha \approx \hat{\varepsilon}_S(h)$ , making (8) similar to (4). The derivation is provided in the Supplementary Material.

The implications of (8) are two-fold: Firstly, to keep the bound tight, we need to assure that both domain discrepancies are small. This motivates the design of our proposed adversarial learning framework discussed earlier. More importantly, we argue that the individual components of *decomposed discrepancies are much easier to optimize* than the one in (4) when the bridging domain is chosen properly.

**Unsupervised Bridging Domain Discovery.** While there are many real-world problems where the bridging domains naturally arise (e.g., the illumination condition of the surveillance images, which can be obtained from the mean pixel values), it is not always available. In such cases, one may resort to the unsupervised discovery of latent domains [14, 17, 21].

To find out whether an unlabeled image of the target domain belongs to the bridging domain, one may measure the closeness of individual target examples to the source domain. For example, we propose to pretrain a standard DANN and exploit the discriminator score  $d_{\text{pre}}(f_{\text{pre}}(x))$ ,  $x \in \mathcal{D}_T$  to quantify the closeness. Since the discriminator converges at equilibrium of source and target distributions [16], this requires an early stopping in practice [8].

Alternatively, we can use off-the-shelf algorithms to compute the distance between individual target examples to the source domain. Given a feature extractor  $f_{\text{pre}}$  trained on the

source examples, one may compute the MMD between  $f_{\text{pre}}(x), x \in \mathcal{D}_T$  and  $\{f_{\text{pre}}(x)\}_{x \in \mathcal{D}_S}$ . In addition, out-of-distribution (OOD) sample detection methods [24] are good candidates as they provide the score quantifying how likely an example belongs to the source domain.

**DANN with Multiple Bridging Domains.** Our framework can be extended to the case for which multiple unlabeled bridging domains exist, which is desirable to span larger discrepancies between source and target domains. To formalize, we denote  $\mathcal{D}_0 = \mathcal{D}_S, \mathcal{D}_{M+1} = \mathcal{D}_T$  as source and target domains, and  $\mathcal{D}_m, m = 1, \dots, M$  as unlabeled bridging domains with  $\mathcal{D}_m$  closer to source than  $\mathcal{D}_{m+1}$ . We introduce  $M+1$  domain discriminators  $d_1, \dots, d_{M+1}$ , each of which is trained by maximizing the following objective:

$$\mathcal{L}_{d_m} = \mathbb{E}_{\bigcup_{i=0}^{m-1} \mathcal{D}_i} \log d_m(f) + \mathbb{E}_{\mathcal{D}_m} \log(1 - d_m(f)), \quad (9)$$

and the learning objective for  $f$  and  $C$  is given as follows:

$$\min_{\theta_f, \theta_C} \mathcal{L}_C + \sum_{m=1}^{M+1} \lambda_m \mathcal{L}_{d_m}. \quad (10)$$

## 4 Experiments

We evaluate our methods mainly on three adaptation tasks: digit classification, object recognition, and semantic scene segmentation. For the recognition task, we use the Comprehensive Cars (CompCars) [46] dataset to recognize car models in the surveillance domain at night using labeled images from the web domain. For the scene segmentation task, synthetic images of the GTA5 dataset [34] are given as the source domain and the task is to perform adaptation on Foggy Cityscapes [35]. In the Supplementary Material, we provide more results on the two-moons toy dataset as described previously and the digit classification task.

### 4.1 Toy Experiment with Two Moons

Created for binary classification problem, the inter-twinning moons 2D dataset suits our model if we consider different rotated versions of the standard two entangled moons as different domains. In this experiment, we consider a hard adaptation from the original data to the ones that are rotated  $90^\circ$  (clockwise or counter-clockwise), while intermediate rotation such as  $30^\circ$  and  $60^\circ$  can be considered as bridging domains. Moreover, as discussed in Section 3.2, the domains do not share the same centers and are proportionally translated according to the rotated angle. We follow the same network architecture as in [14], with one hidden layer of 15 neurons followed by sigmoid non-linearity. The performance is summarized in Table 1.

Model	$0^\circ$	$30^\circ$	$60^\circ$	$90^\circ$
0→90	80.88±1.71	-	-	56.98±4.47
0→30→90	87.23±3.64	95.66±4.18	-	60.98±7.41
0→60→90	79.19±1.21	-	89.66±3.61	80.67±9.47
0→30→60→90	78.75±1.56	82.33±8.71	87.33±3.83	<b>86.97±2.17</b>

Table 1: Average classification accuracy on test set of each domain. Results for the baseline and different bridging domain combinations are included.

One observation is that when  $90^\circ$  is involved as a target domain, the source domain accuracy is sacrificed a lot, which may be because of the limited network capacity. While the adaptation achieves only 56.90%, which is almost a random guess, with source-to-target model (0→90), the proposed method clearly demonstrates its effectiveness, achieving 86.97% on the target domain.

## 4.2 Digit Classification

Different digit datasets are considered as separated domains. MNIST [28] provides a large amount of hand written digit images in gray scale. SVHN [60] contains colored digit images of house numbers from street view. MNIST-M [17] is enriched from MNIST using randomly selected colored image patches in BSD500 [10] as background. We consider adaptation from labeled MNIST to unlabeled SVHN, while using MNIST-M as an unlabeled bridging domain. Given the differences between MNIST and SVHN, MNIST-M seems appropriate bridging domain (similar appearance of foreground digits to MNIST but color statistics to SVHN).

We compare our model with the baseline model, i.e., a standard DANN from source to target without bridging domain. A DANN model that adapts to the mixture of bridge and target domains as a single target is included for comparison. We present results in Table 2. When the bridging domain is involved, the average accuracy on SVHN (target) significantly improves upon the baseline model. Moreover, our proposed model achieves higher performance than the model with mixture of unlabeled domains, demonstrating benefits from the bridging domain.

Model	MNIST-M	SVHN
MNIST→SVHN	-	71.02
MNIST→MNIST-M+SVHN	96.27	78.07
MNIST→MNIST-M→SVHN	97.07	<b>81.28</b>

Table 2: Accuracy on MNIST-M and SVHN test sets averaged over 10 runs. We report the performance of the standard DANN, the DANN model using mixture of unlabeled domains as a single target (MNIST→MNIST-M+SVHN), and our proposed model.

## 4.3 Recognizing Cars in SV Domain at Night

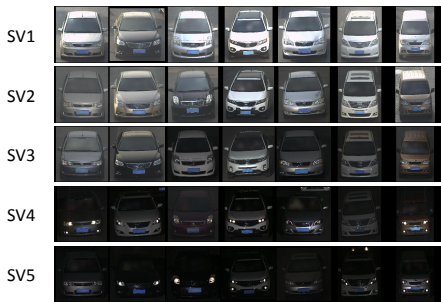


Figure 4: Sample images of CompCars surveillance (SV) domain from light (SV1) to dark (SV5) illumination conditions.



Figure 5: t-SNE plots of CompCars web and SV domains from light (SV1) to dark (SV5) illumination conditions using baseline features.

**Dataset and Experimental Setting.** Two sets of images are provided in the CompCars dataset: 1) the web-nature images are collected from car forums, public websites and search engines, and 2) the surveillance-nature images are collected from surveillance cameras. The dataset is composed of 52,083 web images across 431 car models and 44,481 SV images across 181 car models, with these categories of the SV set being inclusive of 431 categories from the web set. We consider a set of adaptation problems from labeled web to unlabeled SV images. This is challenging as SV images have different perspective and illumination variations from web images.

We use an illumination condition as a metric for adaptation difficulty and partition the SV set into SV1–5 based on the illumination condition of each image.<sup>1</sup> SV1 contains the brightest images, whereas SV5 contains the darkest ones.

We visualize samples from SV1–5 in Figure 4, and confirm the domain discrepancy between

<sup>1</sup>We compute the mean pixel-intensity and sort/threshold images to construct SV1–5 with roughly the same sizes. In practice, the illumination condition may be obtained from metadata, such as recorded time.

web and SV1–5 domains through t-SNE plots in Figure 5. More details of model architecture and training are in the Supplementary Material.

We present two experimental protocols. First, we evaluate on an adaptation task from web to SV night (SV4–5) using SV day (SV1–3) as one domain bridge. We demonstrate the difficulty of adaptation when two domains are far from each other, and show the importance of bridging domain and the effectiveness our adaptation method. Second, we adapt to extreme SV domain (SV5) using different combinations of one or multiple bridging domains (SV1–4) and characterize the properties of an effective bridging domain.

**Evaluation with a Single Bridging Domain.** We demonstrate the difficulty of adaptation when domains are far apart and show that the performance of adversarial DA can be enhanced using bridging domains. In particular, night images (SV4–5) are considered as unlabeled target domain and day images (SV1–3) as unlabeled bridging domain. We compare the following models in Table 3: baseline model trained on labeled web images, DANN from source to target (Web→SV4–5), from source to mixture of bridge and target (Web→SV1–5), and the proposed model from source to bridge to target (Web→SV1–3→SV4–5).

While the DANN adapted to the target domain (SV4–5) improves upon the baseline model, the performance is still far from adequate when compared to the performance of day images. By introducing unlabeled bridging domain, we observe significant improvement in accuracy on the target domain, achieving 77.84% using standard DANN adapted to the mixture of bridging and target domains and 78.78% using our proposed method.

We further conduct an experiment using SV4 as a bridge domain and SV5 as a target domain and compare with the naively trained model (Web→SV4–5). As in Table 3, the proposed model (Web→SV4→5) outperforms the DANN by a large margin (49.83% to 61.37% on SV4–5 test set). This is because it is difficult to determine the adaptation curriculum as both domains are distant from the source domain (see Figure 5), which is different from the previous experiment where there are sufficient amount of day images that are fairly close to the source domain, for discriminator to figure out the curriculum.

To better understand the advantage of our proposed training scheme, we monitor the validation accuracy on SV4–5 of the standard (Web→SV1–5) and the proposed (Web→SV1–3→4–5) models and plot curves in Figure 6. Interestingly, we observe a large fluctuation in the performance of night images using the standard DANN. In contrast, our method allows stable performance earlier in the training, which implies that knowing the curriculum [21] (*i.e.*, adaptation difficulty) is important. Our method with multiple discriminators effectively utilizes such information.

Model	SV1–3	SV4–5
Web (source only)	72.67	19.87
Web→SV4–5	68.90±1.28	49.83±0.70
Web→SV4→5	<b>74.03±0.71</b>	<b>61.37±0.30</b>
Web→SV1–5	<b>83.29±0.14</b>	77.84±0.34
Web→SV1–3→4–5	<b>82.83±0.40</b>	<b>78.78±0.33</b>

Table 3: Accuracy and standard error over 5 runs on SV test sets for models without (Web→SV4–5) and with (Web→SV4→5, Web→SV1–3→4–5) bridging domain. Baselines include a model using mixture of bridge and target domains as a single target domain (Web→SV1–5).

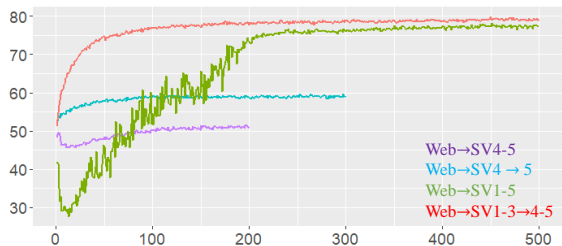


Figure 6: Validation accuracy on SV4–5 (night)



Model	SV5
Web→SV5	37.83±0.51
Web→SV4→5	58.40±0.60
Web→SV1→5	69.69±0.99
Web→SV3→4→5	74.01±0.52
Web→SV2→3→4→5	75.15±0.18
Web→SV1→2→3→4→5	75.47±0.20

Table 4: Accuracy and standard error over 5 runs on SV5 test set.

Split	OOD (0.69)	MMD (0.79)	$d$ -Score (0.85)
2-way	76.43±0.28	78.32±0.27	78.62±0.39
3-way	76.14±0.34	78.47±0.31	78.29±0.28
4-way	76.67±0.32	78.68±0.39	77.93±0.43
5-way	76.54±0.29	78.91±0.33	79.13±0.32

Table 5: Accuracy and standard error over 5 runs on SV4–5 test set with different unsupervised bridging domain discovery configurations. Split  $m$ -way means that we evenly split the unlabeled data into  $m$  domains.

### Which is a Good Bridging Domain?

We perform an ablation study to characterize the properties of a good bridging domain. Specifically, we would like to answer which is a more useful bridging domain: the one closer to the source domain or the one closer to the target domain. To this end, we compare two models, namely, Web→SV1→5 and Web→SV4→5. Note that SV4 is more similar to the target domain (SV5) in terms of visual attributes than SV1.

The results are summarized in Table 4. We observe much higher accuracy on the target domain (SV5) for the model using SV1 as a bridging domain (69.69%) than the one using SV4 (58.40%). We believe that the optimization of the adversarial loss for the second model ( $d_{\mathcal{H}\Delta\mathcal{H}}(\mathbf{W}, \text{SV4}) + d_{\mathcal{H}\Delta\mathcal{H}}(\mathbf{W} + \text{SV4}, \text{SV5})$ ) is more difficult than that of the first model ( $d_{\mathcal{H}\Delta\mathcal{H}}(\mathbf{W}, \text{SV1}) + d_{\mathcal{H}\Delta\mathcal{H}}(\mathbf{W} + \text{SV1}, \text{SV5})$ ) as SV4 is farther from the web domain than SV1. This implies that a good bridging domain decomposes the domain discrepancy between the source and target domains, so that the decomposed discrepancy losses is easily optimized.

**Evaluation with Multiple Bridging Domains.** Our theoretical motivation suggests that, if we have many bridging domains whose generalization error between any two neighboring domains is small, we can also reduce the generalization error between source and target.

We test our hypothesis through experiments that adapt to SV5 with different bridging domain configurations. Specifically, bridging domains are included one by one from SV4 to SV1, and finally reach adaptation with four bridging domains. As in shown Table 4, DANN fails at adaptation without domain bridge (Web→SV5). While including SV4 as the target domain raises adaptation difficulty, using it as a bridging domain (Web→SV4→5) greatly improves the performance on the SV5 test set. Including SV3 as an additional bridging domain (Web→SV3→4→5) shows additional improvement, confirming our hypothesis. While adding SV2 and SV1 as bridging domains leads to an extra improvement, the margin is not as large as including SV3 and SV4. The reason is that SV3 is already close to the web domain as SV1 or SV2 (see Figure 4 and 5), and there is little benefit of introducing an extra bridge.

**Evaluation with Unsupervised Bridge Discovery.** We perform several unsupervised approaches for bridging domain discovery, including discriminator scores of the DANN, MMD, and OOD sample detection [20]. These approaches provide scores indicating the closeness to the source domain distribution for each target example, based on which we can split the target domain into multiple bridging domains (details in the Supplementary Material).

To evaluate the performance of our methods on unsupervised bridging domain discovery, we first compute the AUC between the predicted score and the ground-truth day/night labels and report results in Table 5. While we observe the highest AUC using discriminator of the DANN (0.85), it requires early stopping based on the AUC. In comparison, MMD (0.79) and OOD (0.69) show slightly lower AUC, but are preferred as additional side information is not required. Overall, we observe that using  $d$ -Score or MMD, the performance on SV4–5 test set is improved compared to not using any domain bridges (Table 3 Web→SV1–5).

## 4.4 Foggy Scene Segmentation

**Dataset and Experimental Setting.** We use the GTA5 dataset [24], a synthetic dataset of street-view images, containing 24,966 images of size  $1914 \times 1052$ , as labeled source domain. Unlike previous works [21, 42], we adapt to Foggy Cityscapes [5], a derivative from the real scene images of Cityscapes [8] with a fog simulation, and use Cityscapes as well as Foggy Cityscapes with lighter foggy levels (0.01) as unlabeled bridging domains.

The task is to categorize each pixel into one of 19 semantic categories on the test set images of Foggy Cityscapes with 0.02 foggy level. We consider several models for comparison, such as the traditional DANN ( $GTA5 \rightarrow F_{0.02}$ ), with one bridging domain ( $GTA5 \rightarrow City \rightarrow F_{0.02}$ ), or with two of them ( $GTA5 \rightarrow City \rightarrow F_{0.01} \rightarrow F_{0.02}$ ). One consideration is that we partition 2975 training images of Cityscapes equally for each domain to prevent the case where the algorithm finds an exact correspondence between images from different unlabeled domains.

**Evaluation on Semantic Segmentation.** We utilize the adaptation method by [42] as our base model, which reduces the domain discrepancy at structured output spaces. The same discriminator architecture is used for multiple adversarial losses in our framework.

We first conduct experiments with one bridging domain to adapt to Foggy Cityscapes 0.02 ( $F_{0.02}$ ). We construct two partitions of unlabeled data for Cityscapes and  $F_{0.02}$ . Mean intersection-over-union (IoU) averaged over 5 runs, using different partition for each run, is reported in rows 2–4 of Table 6. While significant improvement in mIoU is observed by directly adapting to the target domain ( $GTA5 \rightarrow F_{0.02}$ ), our framework using a bridging domain further enhances the performance on the final target domain from 33.08 to 34.82. We also conduct a baseline model by merging Cityscapes with the target domain ( $GTA5 \rightarrow City + F_{0.02}$ ), but the performance is not good, indicating that naively merging two domains with different properties may be suboptimal for adversarial adaptation.

We then experiment with two bridging domains by introducing Foggy Cityscapes 0.01 ( $F_{0.01}$ ) as a bridge between Cityscapes and  $F_{0.02}$ . The setting is similar, but we use  $\frac{1}{3}$  of entire images for each unlabeled domain. Table 6 rows 5–7 validate our hypothesis that additional bridging domains are beneficial, improving mIoU from 34.13 to 35.31. While using the same number of overall unlabeled images during the training, we also observe benefit of using two bridging domains (7<sup>th</sup> row) than one (4<sup>th</sup> row).

## 5 Conclusions

This paper aims to simplify adaptation problems with extreme domain variations, using unlabeled bridging domains. A novel framework based on DANN is developed by introducing additional discriminators to account for decomposed many, but smaller discrepancies of the source-to-target domain discrepancy. Several adaptation tasks in computer vision are considered, demonstrating the effectiveness of our framework with bridging domains.

Model	# images	mIoU on $F_{0.02}$
GTA5 (source only)	–	27.5
GTA5 $\rightarrow$ $F_{0.02}$ [42]		33.08 $\pm$ 0.32
GTA5 $\rightarrow$ City + $F_{0.02}$	1487	32.62 $\pm$ 0.58
GTA5 $\rightarrow$ City $\rightarrow$ $F_{0.02}$		<b>34.82</b> $\pm$ 0.39
GTA5 $\rightarrow$ $F_{0.02}$ [42]		33.16 $\pm$ 0.35
GTA5 $\rightarrow$ City $\rightarrow$ $F_{0.02}$	991	34.13 $\pm$ 0.78
GTA5 $\rightarrow$ City $\rightarrow$ $F_{0.01} \rightarrow$ $F_{0.02}$		<b>35.31</b> $\pm$ 0.06

Table 6: mIoU and standard error over 5 runs on Foggy Cityscapes 0.02 ( $F_{0.02}$ ) test set. We partition unlabeled training data into 2 for models in rows 2–4, resulting in 1487 images per domain, while for models in rows 5–7, we partition them into 3, resulting in 991 images per domain.

## References

- [1] Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation. *TPAMI*, 2011.
- [2] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. In *NIPS*, 2007.
- [3] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 2010.
- [4] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML*, 2009.
- [5] Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. Domain separation networks. In *NIPS*, 2016.
- [6] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *TPAMI*, 2018.
- [7] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *CVPR*, 2018.
- [8] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016.
- [9] Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for domain adaptation. *TPAMI*, 2015.
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [11] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 2011.
- [12] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 2016.
- [13] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- [14] Boqing Gong, Kristen Grauman, and Fei Sha. Reshaping visual datasets for domain adaptation. In *NIPS*, 2013.
- [15] Boqing Gong, Kristen Grauman, and Fei Sha. Learning kernels for unsupervised domain adaptation with applications to visual object recognition. In *ICCV*, 2014.
- [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014.

- [17] Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *ICCV*, 2011.
- [18] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar): 723–773, 2012.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [20] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *ICLR*, 2017.
- [21] Judy Hoffman, Dequan Wang, Fisher Yu, and Trevor Darrell. Fcns in the wild: Pixel-level adversarial and constraint-based adaptation. *arXiv preprint arXiv:1612.02649*, 2016.
- [22] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *ICML*, 2018.
- [23] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, June 2018.
- [24] Naoto Inoue, Ryosuke Furuta, Toshihiko Yamasaki, and Kiyoharu Aizawa. Cross-domain weakly-supervised object detection through progressive domain adaptation. In *CVPR*, June 2018.
- [25] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML*, 2015.
- [26] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [28] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- [29] Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation with multiple sources. In *NIPS*, 2009.
- [30] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop*, 2011.
- [31] Sinno Jialin Pan, Qiang Yang, et al. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 2010.
- [32] Sujoy Paul, Yi-Hsuan Tsai, Samuel Schuster, Amit K. Roy-Chowdhury, and Manmohan Chandraker. Domain adaptive semantic segmentation using weak labels. In *European Conference on Computer Vision (ECCV)*, 2020.

- [33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *NIPS*, 2015.
- [34] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *ECCV*, 2016.
- [35] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Semantic foggy scene understanding with synthetic data. *IJCV*, 2018.
- [36] E. Shelhamer, J. Long, and T. Darrell. Fully convolutional networks for semantic segmentation. *PAMI*, 2017.
- [37] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [38] Kihyuk Sohn, Sifei Liu, Guanyu Zhong, Xiang Yu, Ming-Hsuan Yang, and Manmohan Chandraker. Unsupervised domain adaptation for face recognition in unlabeled videos. In *ICCV*, 2017.
- [39] Jong-Chyi Su, Yi-Hsuan Tsai, Kihyuk Sohn, Buyu Liu, Subhansu Maji, and Manmohan Chandraker. Active adversarial domain adaptation. In *WACV*, 2020.
- [40] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- [41] Ben Tan, Yangqiu Song, Erheng Zhong, and Qiang Yang. Transitive transfer learning. In *SIGKDD*, 2015.
- [42] Yi-Hsuan Tsai, Wei-Chih Hung, Samuel Schulter, Kihyuk Sohn, Ming-Hsuan Yang, and Manmohan Chandraker. Learning to adapt structured output space for semantic segmentation. In *CVPR*, June 2018.
- [43] Yi-Hsuan Tsai, Kihyuk Sohn, Samuel Schulter, and Manmohan Chandraker. Domain adaptation for structured output via discriminative patch representations. In *ICCV*, 2019.
- [44] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- [45] Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Simultaneous deep transfer across domains and tasks. In *ICCV*, 2015.
- [46] Linjie Yang, Ping Luo, Chen Change Loy, and Xiaoou Tang. A large-scale car dataset for fine-grained categorization and verification. In *CVPR*, 2015.
- [47] Han Zhao, Shanghang Zhang, Guanhang Wu, José MF Moura, Joao P Costeira, and Geoffrey J Gordon. Adversarial multiple source domain adaptation. In *NIPS*, 2018.