

# Progressive Domain Expansion Network for Single Domain Generalization

Lei Li<sup>12</sup>, Ke Gao<sup>1\*</sup>, Juan Cao<sup>1\*</sup>, Ziyao Huang<sup>12</sup>, Yepeng Weng<sup>12</sup>,  
Xiaoyue Mi<sup>12</sup>, Zhengze Yu<sup>12</sup>, Xiaoya Li<sup>12</sup>, Boyang xia<sup>12</sup>

<sup>1</sup>Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

<sup>2</sup>University of Chinese Academy of Sciences, Beijing, China

{lilei17b, caojuan, huangziyao19f, wengyepeng19s, mixiaoyue19s}@ict.ac.cn

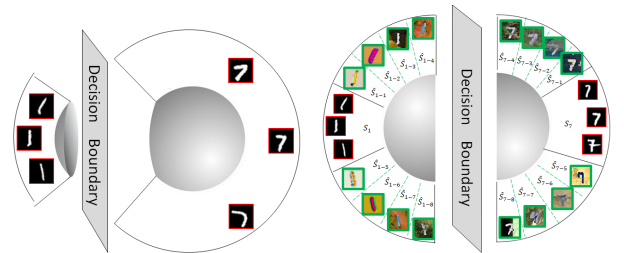
kegao512@gmail.com, {yuzhengze, lixiaoya18s, xiaboyang20s}@ict.ac.cn

## Abstract

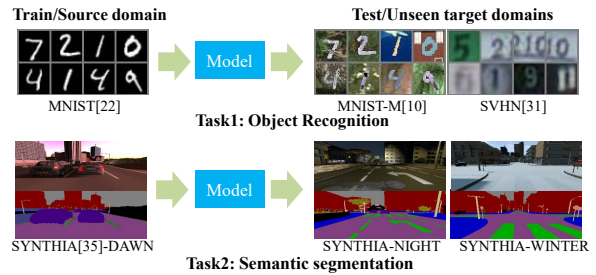
Single domain generalization is a challenging case of model generalization, where the models are trained on a single domain and tested on other unseen domains. A promising solution is to learn cross-domain invariant representations by expanding the coverage of the training domain. These methods have limited generalization performance gains in practical applications due to the lack of appropriate safety and effectiveness constraints. In this paper, we propose a novel learning framework called progressive domain expansion network (PDEN) for single domain generalization. The domain expansion subnetwork and representation learning subnetwork in PDEN mutually benefit from each other by joint learning. For the domain expansion subnetwork, multiple domains are progressively generated in order to simulate various photometric and geometric transforms in unseen domains. A series of strategies are introduced to guarantee the safety and effectiveness of the expanded domains. For the domain invariant representation learning subnetwork, contrastive learning is introduced to learn the domain invariant representation in which each class is well clustered so that a better decision boundary can be learned to improve its generalization. Extensive experiments on classification and segmentation have shown that PDEN can achieve up to 15.28% improvement compared with the state-of-the-art single-domain generalization methods. Codes will be released soon at <https://github.com/lileicv/PDEN>

## 1. Introduction

In this paper, we define domains as various distributions of objects appearance caused by different external conditions (such as weather, background, illumination etc.) or in-



(a) The traditional decision boundary learned with the original training domain. (b) The new decision boundary learned with our progressively expanded domains.



(c) Domain generalization scenario

Figure 1. The illustration of our PDEN for single domain generalization. The tiny images in (a) and (b) with red border denote the source domain and the one with green border denote the expanded domains with our PDEN.

trinsic attributes (such as color, texture, pose etc.), as shown in Fig. 1. The performance of a deep model usually drops when applied to unseen domains. For example, The accuracy of the CNN model (trained on MNIST) on MNIST test set is 99%, while that on SVHN test set is only 30%. Model generalization is important to machine learning.

Two solutions have been proposed to deal with the above issue, namely, domain adaptation [10, 30, 36, 9] and domain generalization [29, 11, 12, 16]. Domain adaptation aims to generalize to a known target domain whose labels are unknown. Distribution alignment (e.g., MMD) and style

\*Corresponding author

transfer(e.g., CycleGAN) are frequently used in these methods to learn domain-invariant features. However, it requires data from the target domain to train the model, which is difficult to achieve in many tasks due to lack of data.

Domain generalization, which not requires access to any data from the unseen target domain, can solve these problems. The idea of domain generalization is to learn a domain-agnostic model from one or multiple source domains. Particularly, in many fields we are usually faced with the challenge of giving a single source domain, which is defined as single domain generalization [33]. Recently, studies have made progress on this task [34, 45, 38, 40, 33, 46]. All of these methods, which are essentially data augmentation, improve the robustness of the model to the unseen domain by extending the distribution of the source domain. Specifically, additional samples are generated by manually selecting the augmentation type[45, 34] or by learning the augmentation through neural networks[33, 46].

Data augmentation has proved to be an important means for improving model generalization [44]. However, such methods require the selection of an augmentation type and magnitude based on the target domain, which is difficult to achieve in other tasks. They cannot guarantee the safety and effectiveness of synthetic data or even reduce accuracy. [42, 20].

In this paper, we propose the progressive domain expansion network (PDEN) to solve the single domain generalization problem. Task models and generators in PDEN mutually benefit from each other through joint learning. Safe and effective domains are generated by the generator under the precise guidance of the task model. The generated domains are progressively expanded to increase the coverage and improve the completeness. Contrastive learning is introduced to learn the cross-domain invariant representation with all the generated domains. It is noteworthy that we can flexibly replace the generator in PDEN to achieve different types of domain expansion.

Our main contributions are as follows:

- We propose a novel framework called progressive domain expansion network (PDEN) for single domain generalization. The PDEN contains domain expansion subnetwork and domain invariant representation learning subnetwork, which mutually benefit from each other by joint learning.
- For the domain expansion subnetwork, multiple domains are progressively generated to simulate various photometric and geometric transforms in unseen domains. A series of strategies are introduced to guarantee the safety and effectiveness of these domains.
- For the domain invariant representation learning subnetwork, contrastive learning is introduced to learn the domain invariant representation in which each class is well clustered so that a better decision boundary can

be learned to improve its generalization.

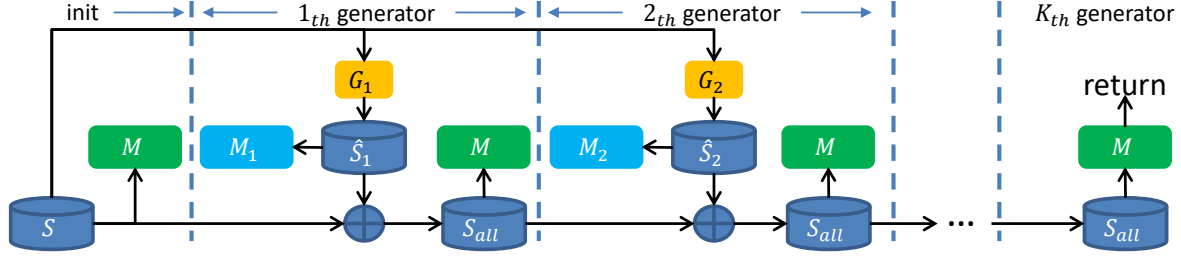
- Extensive experiments on classification and segmentation have shown the superior performance of our method. The proposed method can achieve up to 15.28% improvement compared with other single-domain generalization methods.

## 2. Related Work

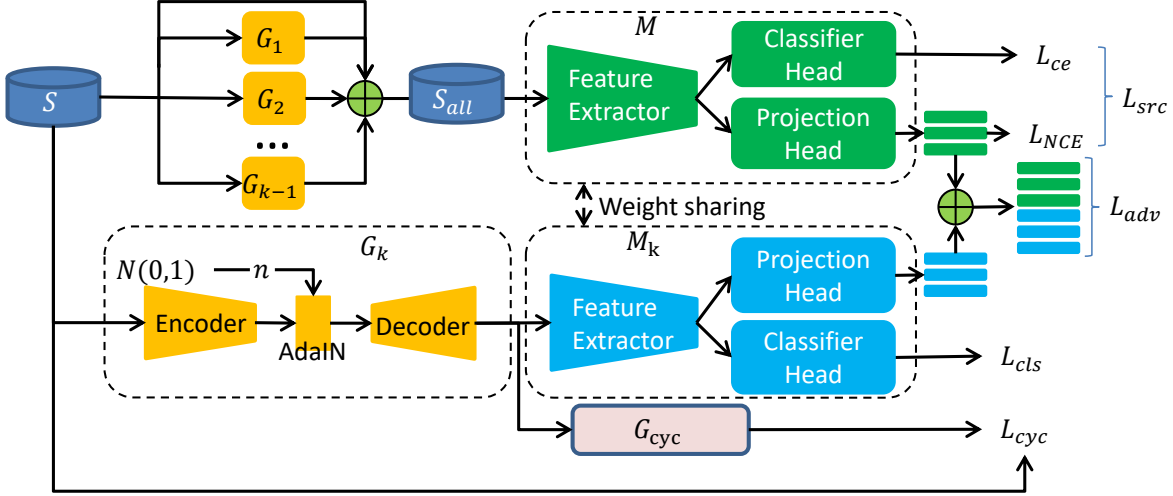
**Domain Adaptation.** In recent years, many domain adaptation methods [10, 30, 36, 9] have been proposed to solve the problem of domain drift between source and target domain, including feature-based adaptation[10], instance-based adaptation [6] and model parameter based adaptation [14]. The domain adaptation method in deep learning is mainly to align the distribution of source domain and target domain, including two kinds of methods: MMD based adaptation method[39, 25] and adversarial based method[10]. DDC[39] is first proposed to solve the domain adaptation problems in deep networks. DDC fixes the weights of the first 7 layers in AlexNet, and MMD is used on the 8th layer to reduce the distribution difference between the source domain and target domain. DAN[25] increased the number of adaptive layers (three in front of the classifier head) and introduced MK-MMD instead of MMD. AdaBN[24] proposed to measure the distribution of the source domain and target domain in BN layer. With the emergence of GAN, a lot of domain adaptation methods based on adversarial learning have been developed. DANN[10] is the first research work to reduce the distribution difference between the source domain and target domain by adversarial learning. DSN[1] assumes that each domain includes a domain-shared distribution and a domain-specific distribution. Based on this assumption, DSN learned the shared feature and the domain-specific feature respectively. DAAN[43] measures the marginal distribution and conditional distribution with a learnable weight.

**Domain Generalization.** Domain generalization is more challenging than domain adaptation. Domain generalization aims to learn the model with data from the source domain and the model can be generalized to unseen domains.

Domain generalization can be categorized as such several research interests: Domain alignment[29, 28, 23, 8] and domain ensemble[26]. Domain alignment methods assume that there is a distribution shared by different domains. These methods map the distribution from different domains to the shared one. CCSA[28] propose the contrastive semantic alignment loss to minimize the distance between data with the same label but from different domains and maximize the distance between the data with different class labels. In MMD-AAE[23], the feature distribution of source domain and target domain are aligned by MMD, and then



(a) Progressive domain expansion one by one



(b) Expand the  $k$ th unseen domain

Figure 2. Illustration of the proposed method PDEN. (a) We show how the domain is progressively extended. We trained the task model  $M$  and unseen domain generator  $G$  alternately.  $G$  is trained to synthesize the unseen domain  $\hat{S}$  under the guidance of  $M$ . Each synthetic domain will be added to the source domain. The task model  $M$  will be finetuned after the source domain is updated. (b) We show the network structure of PDEN. Note that  $M$  and  $M_k$  share the weights, and different  $G$  have the same structure but do not share the weights.

the feature representation is matched to the prior Laplace distribution by AAE. Model ensemble[26] methods train models for each source domain in the training set, and then ensemble their outputs according to the confidence of each model.

Single-domain generalization assumes that the training set only contains samples from just one source domain. Recent, many studies have made progress on this task [34, 45, 38, 40, 33, 46]. These methods are generally applied to synthesize more samples in image space or feature space to expand the range of data distribution in the training set. BigAug[45] observed that the differences in medical images (such as T2 MRI) are mainly different in 3 aspects: image quality, image appearance, and spatial configuration. They augment more variants for the 3 aspects by data augmentation. However, such methods require the selection of an augmentation type and magnitude based on the target domain, which is difficult to achieve in other tasks. GUD[40] and MADA[33] synthesize more data through adversarial learning to promote the model’s robustness. However, on

the one hand, the augmentation type is relatively simple; on the other hand, too much adversarial examples used for training will damage the performance of the classifier.

**Contrastive Learning.** Contrastive learning is a kind of unsupervised pre-training method for image recognition, which is popular these years. The key idea of contrastive learning is to train a model by bringing the positive pairs closer and pushing apart negative pairs. SimCLR[3] generates positive pairs by imposing strong augmentation on whole images. CPC[32] utilizes augmentation on image patches and uses the patch-level views for loss optimization.

### 3. Method

The PDEN proposed in this paper is used for single domain generalization. Suppose the source domain is  $\mathcal{S} = \{x_i, y_i\}_{i=1}^{N_S}$ , the target domain is  $\mathcal{T} = \{x_i, y_i\}_{i=1}^{N_T}$ , where  $x_i, y_i$  is the  $i$ -th image and class label,  $N_S, N_T$  represent the number of samples in source domain and target domain respectively. The aim is to train the model with only  $\mathcal{S}$  then

it can be generalized to the unseen  $\mathcal{T}$ .

### 3.1. The task model $M$

The overall model architecture of PDEN is shown in Fig. 2 (b), including the task net  $M$  and unseen domain generator  $G$ . In this section, we will introduce the task model in the PDEN.

There are 3 parts in  $M$ . (1) Feature extractor  $F : \mathcal{X} \rightarrow \mathcal{H}$ , where  $\mathcal{X}$  is the image space and  $\mathcal{H}$  is the feature space.  $F$  is a stack of convolution layers followed by the pooling layers and activation layers. The output of  $F$  is a 1-d vector obtained by global pooling. (2) Classifier head  $C : \mathcal{H} \rightarrow \mathcal{Y}$ , where  $\mathcal{Y}$  is the label space. Here we focus on the classification task, so the task head  $C$  is optimized by cross-entropy loss. In our experiment,  $C$  is a stack of fully connected layers followed by nonlinear activation layers, and the last activation layer in  $C$  is softmax. (3) Projection head  $P : \mathcal{H} \rightarrow \mathcal{Z}$ , where  $\mathcal{Z}$  is the hidden space in which the contrastive loss will be calculated.  $P$  contains only one full connection layer in our experiments. We normalize the output vector of  $P$  to lie on a unit hypersphere, which enables using an inner product to measure similarity in the  $\mathcal{Z}$  space.

### 3.2. The Unseen Domain Generator $G$

$G$  can convert the original image  $x$ (original domain  $\mathcal{S}$ ) to a new image  $\hat{x}$ (unseen domain  $\hat{\mathcal{S}}$ ) as follows:

$$\begin{aligned} \hat{x} &= G(x, n), n \sim N(0, 1) \\ \hat{\mathcal{S}} &= \{(G(x_i, n), y_i) | (x_i, y_i) \in \mathcal{S}\} \end{aligned} \quad (1)$$

where  $\hat{x}$  has the same semantic information as  $x$ , but the domains of  $\hat{x}$  and  $x$  is different.

$G$  can be a variety of structures depending on related downstream tasks, such as AutoEncoder [18], HRNet [37], spatial transform network(STN) [15] or a combination of these networks.

**Autoencoder as  $G$ :** In our experiment, we mainly use the Autoencoder with AdaIN [17] as the generator, as the  $G_k$  shown in Fig 2. The generator  $G$  contains the encoder  $G_E$ , the AdaIN and the decoder  $G_D$ . In AdaIN, there are two fully-connected layers  $L_{fc1}, L_{fc2}$ :

$$\begin{aligned} AdaIN(z, n) &= L_{fc1}(n) \frac{z - \mu(z)}{\sigma(z)} + L_{fc2}(n) \\ G(x, n) &= G_D(AdaIN(G_E(x), n)) \end{aligned} \quad (2)$$

where  $n \sim N(0, 1)$ . Fig.3(a) shows the unseen domains generated by Autoencoder.

**STN as  $G$ :** The Autoencoder can be replaced by the STN[15] as the generator. The STN is a geometry-aware module which can transform the spatial structure of the image. Fig.3(b) shows the unseen domains generated by STN.

PDEN is a framework in which generators can be replaced with different structures depending on the tasks. In our experiment, the autoencoder is applied.



(a) Domains generated by our domain expansion subnetwork with autoencoder.



(b) Domains generated by our domain expansion subnetwork with STN.

Figure 3. The domains generated by our domain expansion subnetwork.

### 3.3. Progressive Domain Expansion

In order to improve the completeness of the generated domains and expand its coverage, we progressively generate  $K$  unseen domains  $\{\hat{\mathcal{S}}_k = G_k(\mathcal{S})\}_{k=1}^K$  with the learnable generator  $G$ . The task model  $M$  is trained with these unseen domains to learn the cross-domain invariant representation. We train the task model and generator alternately, as shown in Fig. 2.

Take the  $k$ th domain expansion as an example. First, the generator  $G$  and task model  $M$  are jointly trained to synthesize safe and effective unseen domains  $\hat{\mathcal{S}}_k$  by minimize Eq.(9). Then, the task model  $M$  will be retrained with the updated data set  $\mathcal{S} \cup \{\hat{\mathcal{S}}_i\}_{i=1}^k$  by minimize Eq.(3). The performance of  $M$  will be improved, so  $M$  can guide the generator  $G_{k+1}$  to synthesize better unseen domains. The algorithm is shown in Alg.1.

### 3.4. Domain Alignment and Classification

In this section, we will introduce how to learn cross-domain invariant representation. Given a minibatch  $\mathcal{B} = \{x_i, y_i\}_{i=1}^{2N}$ ,  $x_i$  is the source image,  $x_i^+ = G(x_i, n)$  is the synthetic image originating from  $x_i$  ( $x_i$  and  $x_i^+$  have the same semantic information, but come from different domains),  $y_i$  is the class label.  $M$  is optimized by:

$$\begin{aligned} L_{ce}(\hat{y}_i, y_i) &= \min_{F,C} - \sum_m y_i^m \log(\hat{y}_i^m) \\ L_{NCE}(z_i, z_i^+) &= \min_{F,C} - \log \frac{\exp(z_i \cdot z_i^+)}{\sum_{j=1, j \neq i}^{2N} \exp(z_i \cdot z_j)} \\ L_{src} &= L_{ce}(\hat{y}_i, y_i) + L_{NCE}(z_i, z_i^+) \end{aligned} \quad (3)$$

where  $y_i^m$  is the  $m$ th dimension of  $y_i$ ;  $\hat{y}_i = C(F(x_i))$ ;  $z_i = P(F(x_i))$ .

$L_{ce}$  is the cross-entropy loss used for classification.  $L_{NCE}$  is the InfoNCE loss[32] used for contrastive learning. In the minibatch  $\mathcal{B}$ ,  $z_i$  and  $z_i^+$  have the same semantic information but come from different domains. By minimizing  $L_{NCE}$ , the distance between  $z_i$  and  $z_i^+$  will be smaller.

In other words, samples from different domains with the same semantic information will be closer in the  $\mathcal{Z}$  space.  $L_{NCE}$  will guide  $F$  to learn domain-invariant representation.

### 3.5. Unseen Domain $\hat{\mathcal{S}}$ Generation

In this section, we will show how to generate  $k$ th unseen domain  $\hat{\mathcal{S}}_k$  from  $\mathcal{S}$  via the generator  $G_k$  (For convenience, we use  $G, \hat{\mathcal{S}}$  instead of  $G_k, \hat{\mathcal{S}}_k$ ).  $\hat{\mathcal{S}}$  satisfy the constraints of safety and effectiveness. Safety means the generated samples contain the domain-invariant information. Effectiveness means the generated samples contain various unseen domain-specific information.

**Safety.**  $\hat{\mathcal{S}}$  is safe if all the  $x \in \hat{\mathcal{S}}$  can be predicted correctly by task model  $M$ . Formally, we optimize:

$$L_{cls} = \min_{G, F, C} L_{ce}(C(F(G(x, n))), y), n \sim N(0, 1) \quad (4)$$

Cycle consistency loss[47] is introduced to further ensure the safety of  $\hat{\mathcal{S}}$ .  $\hat{\mathcal{S}}$  is safe if it can be converted to  $\mathcal{S}$  by an generator  $G_{cyc}$ .  $G_{cyc}$  has the same structure as  $G$ , but no noise input. Formally, we optimize:

$$L_{cyc} = \min_{G, G_{cyc}} \|x - G_{cyc}(G(x, n))\|_2 \quad (5)$$

**Effectiveness.** Adversarial learning is introduced to generate effective unseen domains. The generator  $G$  and task model  $M$  are learned jointly. The task model  $M$  which extracts the domain-share representation is always trained to minimize the InfoNCE loss. The generator  $G$  is trained to maximize the InfoNCE loss. Through adversarial training,  $G$  will generate unseen domains from which  $M$  can't extract domain shared representation, and  $M$  will be better able to extract cross-domain invariant representations. The loss can be defined as:

$$\tilde{L}_{adv} = \min_G -L_{NCE}(P(F(x)), P(F(G(x, n)))) + \min_{F, P} L_{NCE}(P(F(x)), P(F(G(x, n)))) \quad (6)$$

However, the loss function Equ. 6 is difficult to converge. As the first item in  $\tilde{L}_{adv}$  gets smaller, the gradient gets larger. Therefore, we use the following equation to approximate  $\tilde{L}_{adv}$ .

$$L_{NCE2}(z_i, z_i^+) = \sum_i^{2N} \log \left( 1 - \frac{\exp(z_i \cdot z_i^+)}{\sum_{j=1, j \neq i}^{2N} \exp(z_i \cdot z_j)} \right)$$

$$L_{adv} = \min_G -L_{NCE2}(P(F(x)), P(F(G(x, n)))) + \min_{F, P} L_{NCE}(P(F(x)), P(F(G(x, n)))) \quad (7)$$

We also use a loss function to encourage  $G$  to generate more diverse samples.

$$L_{div} = \min_G -\|G(x, n_1) - G(x, n_2)\|_2 \quad (8)$$

---

### Algorithm 1 PDEN

---

**Input:** Source domain dataset  $\mathcal{S}$ ; Pre-train task model  $M$ ; Number of synthetic domains  $K$

**Output:** learned task model  $M$

```

1: Initialize:  $\mathcal{S}_{all} \leftarrow \mathcal{S}$ 
2: for  $k=1, \dots, K$  do
3:   initialize the weights of  $G_k$  randomly
4:   for  $t=1, \dots, T$  do ▷ Train  $G_k$  to get  $\hat{\mathcal{S}}_k$ 
5:     Sample  $(x_i, y_i)$  from  $\mathcal{S}$ 
6:      $(x_i^+, y_i) \leftarrow (G_k(x_i, n), y_i)$ 
7:     train  $G$  and  $M$  using Eq.(9)
8:   Synthetic  $k$ th unseen domain  $\hat{\mathcal{S}}_k$  using Eq.(1)
9:    $\mathcal{S}_{all} = \mathcal{S} \cup \hat{\mathcal{S}}_k$ 
10:  for  $t=1, \dots, T$  do ▷ Retrain  $M$ 
11:    Sample  $(x_i, y_i)$  from  $\mathcal{S}_{all}$ 
12:    train  $M$  using Eq.(3)
13: return  $M$ 

```

---

where  $n_1, n_2 \sim N(0, 1)$ , and  $n_1 \neq n_2$ . To sum up, the loss function of training generate  $G$  is as follow:

$$L_{unseen} = L_{cls} + w_{cyc} \cdot L_{cyc} + w_{adv} \cdot L_{adv} + w_{div} \cdot L_{div} \quad (9)$$

The weight of  $L_{cls}$  is always 1,  $w_{cyc}, w_{adv}, w_{div}$  are the weights of  $L_{cyc}, L_{adv}, L_{div}$ .

## 4. Experiment

### 4.1. Datasets and Evaluate

Follow [33, 40], we evaluated our approach on Digits, CIFAR10-C and SYNTHIA.

**Digits Dataset:** Digits dataset contains 5 datasets: MNIST[22], MNSIT-M[10], SVHN[31], USPS[7], SYNDIGIT[10]. Each dataset is considered as a domain. We use MNIST as the source domain and the other four data sets as the target domains. The first 10,000 images in MNIST are used to train the model.

**CIFAR10-C Dataset:** We use the CIFAR10[21] as the source domain and the CIFAR10-C[13] as the target domain. CIFAR10-C is a benchmark dataset to evaluate the robustness of classification models. CIFAR10-C dataset consists of test images with 19 corruption types, which are algorithmically generated. The corruptions come from 4 categories and each type of corruption has 5 levels of severity.

**SYNTHIA Dataset:** The SYNTHIA VIDEO SEQUENCES[35] dataset is used for traffic scene segmentation. The dataset consists of 3 locations: Highway, New York ish and Old European Town. Each location consists of the same traffic situation but under different weather/illumination/season conditions (we use Dawn, Fog, Spring, Night and Winter in our experiment). Following the protocol in[40], we train our model on one domain

| Method                         | Manual Data Augmentation | SVHN                             | MNIST-M                          | SYNDIGIT                         | USPS                    | Avg.                             |
|--------------------------------|--------------------------|----------------------------------|----------------------------------|----------------------------------|-------------------------|----------------------------------|
| ERM[19]                        | False                    | 27.83                            | 52.72                            | 39.65                            | 76.94                   | 49.29                            |
| CCSA, WVU, 2017[28]            | False                    | 25.89                            | 49.29                            | 37.31                            | 83.72                   | 49.05                            |
| d-SNE, UH, 2019[41]            | False                    | 26.22                            | 50.98                            | 37.83                            | <b>93.16</b>            | 52.05                            |
| JiGen, Huawei, London, 2019[2] | False                    | 33.80                            | 57.80                            | 43.79                            | 77.15                   | 53.14                            |
| GUD, Stanford, 2018[40]        | False                    | 35.51                            | 60.41                            | 45.32                            | 77.26                   | 54.62                            |
| MADA, UDel, 2020[33]           | False                    | 42.55                            | 67.94                            | 48.95                            | 78.53                   | 59.49                            |
| PDEN                           | False                    | <b>62.21</b> (19.66 $\uparrow$ ) | <b>82.20</b> (14.26 $\uparrow$ ) | <b>69.39</b> (20.44 $\uparrow$ ) | 85.26(6.73 $\uparrow$ ) | <b>74.77</b> (15.28 $\uparrow$ ) |
| AutoAugment, Google, 2018 [4]  | True                     | 45.23                            | 60.53                            | 64.52                            | 80.62                   | 62.72                            |
| RandAugment, Google, 2020 [5]  | True                     | 54.77                            | 74.05                            | 59.60                            | 77.33                   | 66.44                            |
| PDEN                           | False                    | <b>62.21</b> (7.44 $\uparrow$ )  | <b>82.20</b> (8.15 $\uparrow$ )  | <b>69.39</b> ( 9.79 $\uparrow$ ) | 85.26(7.93 $\uparrow$ ) | <b>74.77</b> (8.33 $\uparrow$ )  |

Table 1. Experiment results on Digits dataset. All the models are trained on MNIST. The top half of the table is the comparison with other single domain generalization methods. None of these methods use manual data augmentation. The following section of the table is the comparison with other methods which use manual data augmentation.

and evaluate on the other domains. For each domain, we randomly sample 900 images from the left front camera and all the images are resized to  $192 \times 320$  pixels.

**Evaluate:** For Digit and CIFAR10 datasets, we compute the mean accuracy on each unseen domain. For SYNTHIA dataset, we use the standard mean Intersection over Union(mIoU) to evaluate the performance on each unseen domain.

## 4.2. Evaluation of Single Domain Generalization

We compare our method with the following state-of-the-art methods. (1) Empirical Risk Minimization(ERM) [19] is the baseline method trained with only the cross-entropy loss. (2) CCSA [28] aligns samples from different domains of the same category to get the robust feature space for domain generalization. (3) d-SNE[41] minimizes the maximum distance between sample pairs of the same class and maximizes the minimum distance among sample pairs of different categories. (4) GUD [40] proposes an adversarial data augmentation method to synthesize more hard samples which can improve the robustness of the classifier. (5) MADA [33] minimizes the distance of semantic space and maximize the distance of pixel space to generate more effective samples. (6) JiGen [2] proposes a multi-task learning method that combines the target recognition task and the Jigsaw classification task to improve the cross-domain generalization of the model. (7) AutoAugment(AA) [4] proposes a method to automatically searches improved data augmentation policies for the specific data set. (8) Based on AA, RandAugment(RA) [5] has a better data augment policies, which greatly reduces the policies space .

**Comparison on Digits:** We train the model with the first 10,000 images in the MNIST train set, validate the model on the MNIST test set, and evaluate the model on the MNIST-M, SVHN, USPS, and Syndigits datasets. We calculate the mean accuracy on each data set as the evaluation index. We

first compared with the single-domain generalization methods, as shown in the top half of Table 1. To be fair, we did not use any manual data augmentation. We observed that our method performs much better than other methods on SVHN, MNIST-M and USPS. On USPS, the performance of our method is comparable to others, mainly because the USPS is more similar to MNIST. The d-SNE[41] performs well on USPS, but bad on other data sets. We also compare with the data augmentation methods as shown in the bottom half of Table 1. The hyperparameters are consistent with those in the original paper. We found that our method performs better than these methods. What’s more, our approach is orthogonal to these data augmentation techniques.

**Comparison on CIFAR10:** We train all the models on the CIFAR10 train set, validate the models on the CIFAR10 test set, and evaluate the models on the CIFAR10-C. The experimental results across five levels of corruption severity are shown in Tab2. Our approach performs better than other single-domain generalization methods such as GUD and MADA. The severer the corruption, the more our approach surpasses MADA. Compared to approaches using manual data augmentation, our approach performs as well as they do at lower corruption levels and outperform them at higher corruption levels. We also show the experimental results across different types of corruptions with the 5th level severity in Tab 3. Our approach has higher average accuracy than other approaches. In some corruption types, the RandAugment approach performs better than us. However, it is important to note that there is no manual data augmentation in our approach, and our approach can be used together with RandAugment.

**Comparison on SYNTHIA:** Follow the protocol in [33], we conducted three experiments, using Highway-Dawn, Highway-Fog and Highway-Spring as the source domain respectively, and taking all the weather in New York ish and Old European Town as the unseen target do-

| Method   | Level1       | Level2       | Level3       | Level4       | Level5       |
|----------|--------------|--------------|--------------|--------------|--------------|
| ERM[19]  | 87.8         | 81.5         | 75.5         | 68.2         | 56.1         |
| GUD[40]  | 88.3         | 83.5         | 77.6         | 70.6         | 58.3         |
| MADA[33] | 90.5         | 86.8         | 82.5         | 76.4         | 65.6         |
| AA[4]    | 91.42        | 87.88        | 84.10        | 78.46        | 71.13        |
| RA[5]    | <b>91.74</b> | 88.89        | 85.82        | 81.03        | 74.93        |
| PDEN     | 90.62        | <b>88.91</b> | <b>87.03</b> | <b>83.71</b> | <b>77.47</b> |

Table 2. Experiment results on CIFAR10-C dataset across different levels.

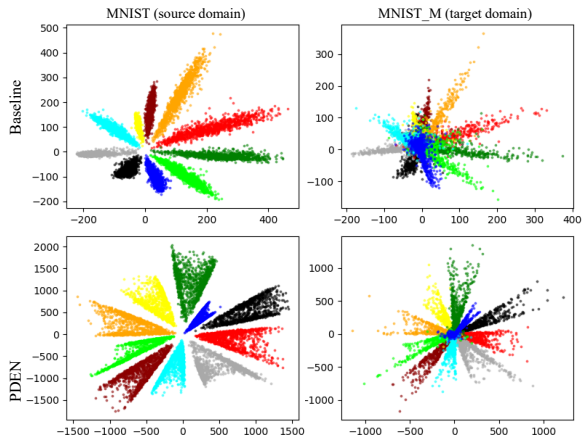


Figure 4. Visualization of different domains in the feature space. Rows 1 and 2 represent the feature space of baseline model and ours, respectively. Columns 1 and 2 represent the feature distribution of MNIST and MNIST\_M, respectively.

mains. The scene segmentation results(mIoU) are show in Tab 4. Our approach improves the average mIoU compared to other approaches. When the source domain is highway-Dawn or Highway-fog, the improvement is greater.

### 4.3. Additional Analysis

**Validation of  $K$ :** We study the effect of the hyper-parameters  $K$  on the Digits dataset. We use the MNIST as the source domain, and take MNIST-M, SVHN, USPS and SYNDIGIT as the unseen target domains. The experimental result is shown in Fig.5(a). We report the classification accuracy on the target domains when  $K = 1, 2, \dots, 20$ . The accuracy is increased rapidly when  $K$  is small, and gradually converges when  $K$  is large. In experiments on Digits, we set  $K = 20$ . In the Digits experiment in MADA[33], their approach performed best at  $K = 3$  and decreased as  $K$  grew. This indicates that the domain generated by our approach is safer than MADA.

**Validation of  $w_{adv}$ :** We study the effect of the hyper-parameters  $w_{adv}$  on Digits dataset. The experimental results are shown in Fig.5(b). We report the classification accuracy on target domains when  $w_{adv} = 0.02, 0.05, 0.08, 0.1, 0.13, 0.16, 0.2$ . We found that the accuracy increases with the increase of  $w_{adv}$  on the unseen

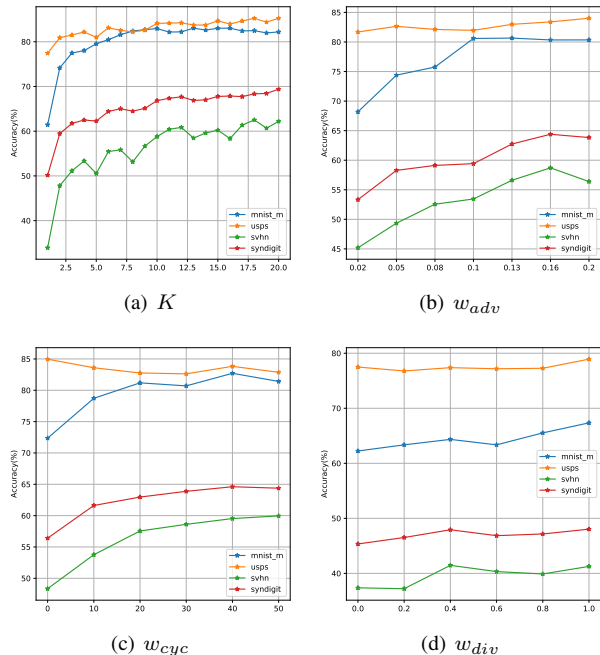


Figure 5. Hyper-parameters tuning of  $K, w_{adv}, w_{cyc}$  on the Digit dataset.

target domains.

**Validation of  $w_{cyc}$ :** We study the effect of the hyper-parameters  $w_{cyc}$  on Digits dataset. The experimental result is shown in Fig.5(c). We report the classification accuracy on MNIST-M, USPS, SVHN and SYNDIGIT when  $w_{cyc} = 0, 10, 20, 30, 40, 50$ . On MNIST-M, SVHN and SYNDIGIT, the accuracy increases with the increases of  $w_{cyc}$ . On USPS, the classification accuracy did not change significantly (fluctuated within a small range) with the increase of  $w_{cyc}$ , mainly because of the high similarity between USPS and MNIST.

**Validation of  $w_{div}$ :** We illustrate the effect of the hyper-parameters  $w_{div}$  in Fig.5(d). For all the unseen domain in Digit dataset, the classification accuracy increases with the increases of  $w_{div}$ .

**Visualization of the feature space:** Fig.4 illustrates the difference in 2-d feature spaces between PDEN and the baseline models. For PDEN, the sample distribution of target domain is consistent with that of source domain. For the baseline model, most of the target samples are mixed in the feature space so that it is difficult to classify them.

### 4.4. Evaluation of of Few-shot Domain Adaptation

We also compared our methods in the experimental setting of few-shot domain adaptation [27]. In few-shot domain adaptation, data from the source domain  $\mathcal{S}$  and a few samples from the target domain  $\mathcal{T}$  are used to train the model.

We use MNIST as the source domain and SVHN as the

|           | Weather      |              |              | Blur         |              |              | Noise        |              |              | Digital      |              |              | Avg.         |
|-----------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|           | Fog          | Snow         | Frost        | Zoom         | Defocus      | Glass        | Speckle      | Shot         | Impulse      | Jpeg         | Pixelate     | Spatter      |              |
| ERM[19]   | 65.92        | 74.36        | 61.57        | 59.97        | 53.71        | 49.44        | 41.31        | 35.41        | 25.65        | 69.90        | 41.07        | 75.36        | 56.15        |
| CCSA[28]  | 66.94        | 74.55        | 61.49        | 61.96        | 56.11        | 48.46        | 40.12        | 33.79        | 24.56        | 69.68        | 40.94        | 77.91        | 56.31        |
| d-SNE[41] | 65.99        | 75.46        | 62.25        | 58.47        | 53.71        | 50.48        | 45.30        | 39.93        | 27.95        | 70.20        | 38.46        | 73.40        | 56.96        |
| GUD[40]   | 68.29        | 76.75        | 69.94        | 62.95        | 56.41        | 53.45        | 38.45        | 36.87        | 22.26        | 74.22        | 53.34        | 80.27        | 58.26        |
| MADA[33]  | 69.36        | 80.59        | 76.66        | 68.04        | 61.18        | 61.59        | 60.88        | 60.58        | 45.18        | 77.14        | 52.25        | 80.62        | 65.59        |
| AA[4]     | 84.61        | 81.04        | 72.32        | 83.94        | 84.38        | 52.29        | 52.14        | 45.40        | 52.54        | 73.65        | 36.12        | 89.13        | 71.13        |
| RA[5]     | <b>85.99</b> | 80.13        | 74.97        | <b>88.60</b> | <b>89.33</b> | 57.70        | 60.50        | 56.03        | 55.64        | 74.92        | 37.36        | <b>90.42</b> | 74.93        |
| PDEN      | 69.64        | <b>81.81</b> | <b>84.50</b> | 83.73        | 82.15        | <b>60.13</b> | <b>79.31</b> | <b>81.28</b> | <b>66.79</b> | <b>85.24</b> | <b>70.82</b> | 79.38        | <b>77.47</b> |

Table 3. The experimental result on CIFAR10-C. The model is trained on the clean data of CIFAR10 and evaluate on CIFAR10-C. We compared the accuracy of 19 types of corruption(only 12 corruptions are shown in the table) at level 5(the severest) in different methods.

| Source Domain  | Method   | New York ish |              |              |              |              | Old European Town |              |              |              |              | Avg.         |
|----------------|----------|--------------|--------------|--------------|--------------|--------------|-------------------|--------------|--------------|--------------|--------------|--------------|
|                |          | Dawn         | Fog          | Night        | Spring       | Winter       | Dawn              | Fog          | Night        | Spring       | Winter       |              |
| Highway/Dawn   | ERM[19]  | 27.80        | 2.73         | 0.93         | 6.80         | 1.65         | 52.78             | 31.37        | 15.86        | 33.78        | 13.35        | 18.70        |
|                | GUD[40]  | 27.14        | 4.05         | 1.63         | 7.22         | 2.83         | 52.80             | 34.43        | 18.19        | 33.58        | 14.68        | 19.66        |
|                | MADA[33] | 29.10        | 4.43         | 4.75         | 14.13        | 4.97         | 54.28             | 36.04        | 23.19        | 37.53        | 14.87        | 22.33        |
|                | PDEN     | <b>30.63</b> | <b>21.74</b> | <b>16.76</b> | <b>26.10</b> | <b>19.91</b> | <b>54.93</b>      | <b>47.55</b> | <b>36.97</b> | <b>43.98</b> | <b>23.83</b> | <b>32.24</b> |
| Highway/Fog    | ERM[19]  | 17.24        | 34.80        | 12.36        | 26.38        | 11.81        | 33.73             | 55.03        | 26.19        | 41.74        | 12.32        | 27.16        |
|                | GUD[40]  | 18.75        | <b>35.58</b> | 12.77        | 26.02        | 13.05        | 37.27             | <b>56.69</b> | 28.06        | 43.57        | 13.59        | 28.53        |
|                | MADA[33] | 21.74        | 32.00        | 9.74         | 26.40        | 13.28        | 42.79             | 56.60        | 31.79        | 42.77        | 12.85        | 29.00        |
|                | PDEN     | <b>25.61</b> | 35.16        | <b>17.05</b> | <b>32.45</b> | <b>21.03</b> | <b>45.67</b>      | 54.91        | <b>37.38</b> | <b>48.29</b> | <b>20.80</b> | <b>33.83</b> |
| Highway/Spring | ERM[19]  | 26.75        | 26.41        | 18.22        | 32.89        | 24.60        | 51.72             | 51.85        | 35.65        | 54.00        | 28.13        | 35.02        |
|                | GUD[40]  | 28.84        | <b>29.67</b> | 20.85        | 35.32        | 27.87        | 52.21             | <b>52.87</b> | 35.99        | 55.30        | 29.58        | 36.85        |
|                | MADA[33] | <b>29.70</b> | 31.03        | 22.22        | <b>38.19</b> | 28.29        | 53.57             | 51.83        | 38.98        | <b>55.63</b> | 25.29        | 37.47        |
|                | PDEN     | 28.17        | 27.67        | <b>27.53</b> | 34.30        | <b>28.85</b> | <b>53.75</b>      | 51.53        | <b>46.87</b> | <b>55.63</b> | <b>30.61</b> | <b>38.49</b> |

Table 4. Segmentation experiment results on SYNTHIA. We report the mIoU. All the models are trained on Highway and tested in New York ish and Old European Town.

target domain. We first train the model on mnist with the proposed PDEN, and then finetune the model with few samples from SVHN. The model will be evaluated on SVHN, as shown in Fig. 6. We found that finetuning with few samples from the target domain can significantly improve the performance of the model on the target domain. Compared with MADA, the proposed PDEN performs better in this case.

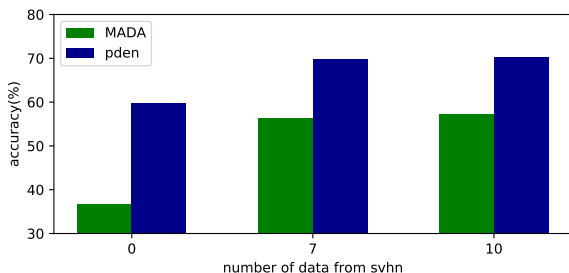


Figure 6. Few-shot domain adaptation experiment. We train the model with all the samples from MNIST and few samples from SVHN, and test the model on SVHN.

## 5. Conclusion

In this paper, we propose a single domain generalization learning framework to learn the domain-invariant feature, which can generalize the model to the unseen domains. We learn the generator to synthetic unseen domains, which share the same semantic information as the source domain. The domain-invariant representation can be learned by aligned the source and unseen domain distribution. We mine the hard unseen domains in which the domain-invariant representation cannot be extracted by the task model. The model will be more robust by adding these generated domains to the training set. The novel method PDEN proposed in this paper provides a promising direction to solve the single-domain generalization problem.

## Acknowledgements

This work was supported by the National Key Research and Development Program of China (2018YFC0825202), and the National Natural Science Foundation of China (U1703261,61871004), and Beijing Municipal Natural Science Foundation Cooperation Beijing Education Committee: No. KZ 201810005002.



## References

- [1] Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. Domain separation networks. In *Advances in neural information processing systems*, pages 343–351, 2016. [2](#)
- [2] Fabio M Carlucci, Antonio D’Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. Domain generalization by solving jigsaw puzzles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2229–2238, 2019. [6](#)
- [3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020. [3](#)
- [4] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 113–123, 2019. [6](#), [7](#), [8](#)
- [5] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 702–703, 2020. [6](#), [7](#), [8](#)
- [6] Wenyuan Dai, Qiang Yang, Gui-Rong Xue, and Yong Yu. Boosting for transfer learning. In *Proceedings of the 24th international conference on Machine learning*, pages 193–200, 2007. [2](#)
- [7] JS Denker, WR Gardner, HP Graf, D Henderson, RE Howard, W Hubbard, LD Jackel, HS Baird, and I Guyon. Advances in neural information processing systems 1. chapter neural network recognizer for hand-written zip code digits. 1989. [5](#)
- [8] Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. Domain generalization via model-agnostic learning of semantic features. *Advances in Neural Information Processing Systems*, 32:6450–6461, 2019. [2](#)
- [9] Geoffrey French, Michal Mackiewicz, and Mark Fisher. Self-ensembling for visual domain adaptation. In *International Conference on Learning Representations*, number 6, 2018. [1](#), [2](#)
- [10] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015. [1](#), [2](#), [5](#)
- [11] Muhammad Ghifary, W Bastiaan Kleijn, Mengjie Zhang, and David Balduzzi. Domain generalization for object recognition with multi-task autoencoders. In *Proceedings of the IEEE international conference on computer vision*, pages 2551–2559, 2015. [1](#)
- [12] Thomas Grubinger, Adriana Birlutiu, Holger Schöner, Thomas Natschläger, and Tom Heskes. Multi-domain transfer component analysis for domain generalization. *Neural processing letters*, 46(3):845–855, 2017. [1](#)
- [13] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. [5](#)
- [14] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. [2](#)
- [15] Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. In *Advances in neural information processing systems*, pages 2017–2025, 2015. [4](#)
- [16] Yunpei Jia, Jie Zhang, Shiguang Shan, and Xilin Chen. Single-side domain generalization for face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8484–8493, 2020. [1](#)
- [17] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4401–4410, 2019. [4](#)
- [18] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. [4](#)
- [19] Vladimir Koltchinskii. *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems: Ecole d’Été de Probabilités de Saint-Flour XXXVIII-2008*, volume 2033. Springer Science & Business Media, 2011. [6](#), [7](#), [8](#)
- [20] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013. [2](#)
- [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. [5](#)
- [22] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. [5](#)
- [23] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5400–5409, 2018. [2](#)
- [24] Yanghao Li, Naiyan Wang, Jianping Shi, Xiaodi Hou, and Jiaying Liu. Adaptive batch normalization for practical domain adaptation. *Pattern Recognition*, 80:109–117, 2018. [2](#)
- [25] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015. [2](#)
- [26] Massimiliano Mancini, Samuel Rota Bulo, Barbara Caputo, and Elisa Ricci. Best sources forward: domain generalization through source-specific nets. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 1353–1357. IEEE, 2018. [2](#), [3](#)
- [27] Saeid Motiian, Quinn Jones, Seyed Iranmanesh, and Gianfranco Doretto. Few-shot adversarial domain adaptation. In *Advances in Neural Information Processing Systems*, pages 6670–6680, 2017. [7](#)
- [28] Saeid Motiian, Marco Piccirilli, Donald A Adjeroh, and Gianfranco Doretto. Unified deep supervised domain adaptation and generalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5715–5725, 2017. [2](#), [6](#), [8](#)
- [29] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant fea-

- ture representation. In *International Conference on Machine Learning*, pages 10–18, 2013. 1, 2
- [30] Zak Murez, Soheil Kolouri, David Kriegman, Ravi Ramamoorthi, and Kyungnam Kim. Image to image translation for domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4500–4509, 2018. 1, 2
- [31] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bis-sacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 5
- [32] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 3, 4
- [33] Fengchun Qiao, Long Zhao, and Xi Peng. Learning to learn single domain generalization. 2020. 2, 3, 5, 6, 7, 8
- [34] Eduardo Romera, Luis M. Bergasa, Jose M. Alvarez, and Mohan Trivedi. Train here, deploy there: Robust segmentation in unseen domains. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, 2018. 2, 3
- [35] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M. Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 5
- [36] Rui Shu, Hung Bui, Hirokazu Narui, and Stefano Ermon. A DIRT-t approach to unsupervised domain adaptation. In *International Conference on Learning Representations*, 2018. 1, 2
- [37] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5693–5703, 2019. 4
- [38] Thanh Dat Truong, Chi Nhan Duong, Khoa Luu, Minh Triet Tran, and Minh Do. Beyond domain adaptation: Unseen domain encapsulation via universal non-volume preserving models. 2018. 2, 3
- [39] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014. 2
- [40] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John C Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 5334–5344. Curran Associates, Inc., 2018. 2, 3, 5, 6, 7, 8
- [41] Xiang Xu, Xiong Zhou, Ragav Venkatesan, Gurumurthy Swaminathan, and Orchid Majumder. d-sne: Domain adaptation using stochastic neighborhood embedding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2497–2506, 2019. 6, 8
- [42] Quanzeng You, Jiebo Luo, Hailin Jin, and Jianchao Yang. Robust image sentiment analysis using progressively trained and domain transferred deep networks. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pages 381–388, 2015. 2
- [43] Chaohui Yu, Jindong Wang, Yiqiang Chen, and Meiyu Huang. Transfer learning with dynamic adversarial adaptation network. In *2019 IEEE International Conference on Data Mining (ICDM)*, pages 778–786. IEEE, 2019. 2
- [44] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *Int. Conf. Learn. Represent.*, 2017. 2
- [45] Ling Zhang, Daguang Xu, Ziyue Xu, Xiaosong Wang, and Andriy Myronenko. Generalizing deep learning for medical image segmentation to unseen domains via deep stacked transformation. *IEEE Transactions on Medical Imaging*, PP(99):1–1, 2020. 2, 3
- [46] Long Zhao, Ting Liu, Xi Peng, and Dimitris Metaxas. Maximum-entropy adversarial data augmentation for improved generalization and robustness. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2, 3
- [47] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 5