

Cross-Domain Ensemble Distillation for Domain Generalization

Kyungmoon Lee^{1,2} Sungyeon Kim¹ Suha Kwak¹

¹POSTECH, Pohang, Korea ²NALBI Inc., Seoul, Korea
kyungmoon@nalbi.ai, {sungyeon.kim, suha.kwak}@postech.ac.kr
<http://cvlab.postech.ac.kr/research/XDED/>

Abstract. Domain generalization is the task of learning models that generalize to unseen target domains. We propose a simple yet effective method for domain generalization, named cross-domain ensemble distillation (XDED), that learns domain-invariant features while encouraging the model to converge to flat minima, which recently turned out to be a sufficient condition for domain generalization. To this end, our method generates an ensemble of the output logits from training data with the same label but from different domains and then penalizes each output for the mismatch with the ensemble. Also, we present a de-stylization technique that standardizes features to encourage the model to produce style-consistent predictions even in an arbitrary target domain. Our method greatly improves generalization capability in public benchmarks for cross-domain image classification, cross-dataset person re-ID, and cross-dataset semantic segmentation. Moreover, we show that models learned by our method are robust against adversarial attacks and image corruptions.

Keywords: domain generalization, knowledge distillation, flat minima

1 Introduction

Deep neural networks (DNNs) have brought remarkable advances in a number of research areas such as image classification [43], image synthesis [24], and reinforcement learning [54]. The huge success of DNNs depends heavily on the assumption that training and test data are sampled under the independent and identically distributed (i.i.d.) condition. However, this assumption often does not hold in real-world scenarios; a large error occurs due to the discrepancy between training and test data, also known as the domain shift problem. As a solution to this problem, domain generalization, the task of learning models that generalize to unseen target domains, is in the spotlight. A key to the success of domain generalization is to learn invariant features across domains. To this end, most previous methods align feature distributions of multiple domains by adversarial training [48,49], minimizing the dissimilarity between the distributions of source domains [55], or contrastive learning [40]. Then, a classifier is trained to predict the labels for the aligned source features in hopes that it will also generalize well

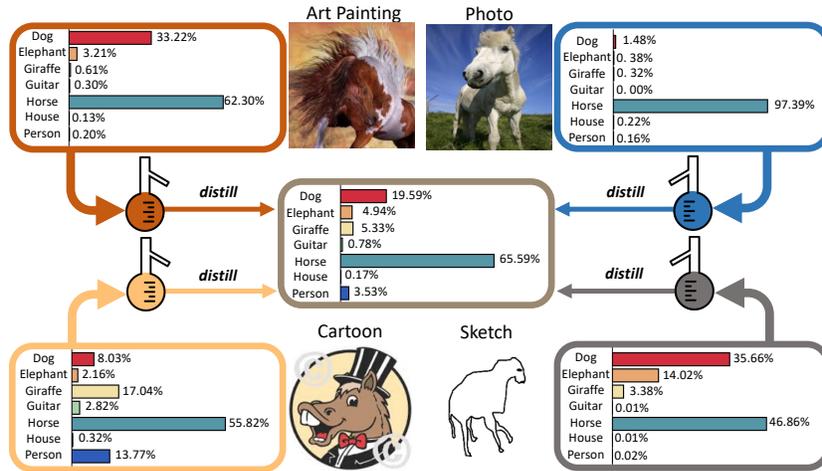


Fig. 1: Illustration of cross-domain ensemble distillation (XDED). Although the four images share the same class label, their predictions manifest different inter-class relations due to the visual gap between domains. XDED constructs an ensemble by averaging all predictions and matches it with each prediction.

for any target domain. However, this approach often drops performance when the target domain differs substantially from the source domains as the model is prone to overfit to the source domains.

Meanwhile, the relationship between the geometry of loss landscapes and generalization ability has attracted increasing attention [18,20,36,38]. In particular, converging to flat minima in loss landscapes is known as a key to achieve robustness against the loss landscape shift between training and test datasets. Inspired by the observation that higher posterior entropy helps a model converge to flat minima [10,60,89], entropy regularization techniques like self-knowledge distillation [87] and entropy maximization [9] have been proposed to increase entropy rather than forcing a model to completely fit training data (*i.e.*, one-hot labels) to induce low entropy. Since the degree of loss landscape shift is generally expected to be bigger in the case of domain generalization, it is more important to converge to flat minima in domain generalization. However, the benefit of flat minima in terms of domain generalization has not been actively studied yet.

In this paper, we propose a novel method, named *cross-domain ensemble distillation* (XDED), that learns domain-invariant features while encouraging convergence to flat minima for domain generalization. Specifically, XDED generates an ensemble of the output logits for the data with the same label but from different domains, and then penalizes each output for the mismatch with the ensemble (Fig. 1). By doing so, it enables a model to learn domain-invariant features by enforcing prediction consistency between the data with the same label but from different domains. Also, XDED increases the posterior entropy of each output distribution, which helps the model converge to flat minima as the

entropy regularization does. To the best of our knowledge, XDED is the first to achieve these two objectives simultaneously for domain generalization, and this contribution leads to significant performance improvement.

Since XDED is still limited to exploiting the information of only source domains, there is further room to reduce the domain gap with the target domain. Hence, we also introduce a de-stylization technique well-suited to domain generalization, called UniStyle. UniStyle suppresses domain-specific style bias simply by standardizing intermediate feature maps of input image during both training and test time. Thanks to UniStyle, our model produces style-consistent predictions not only for the source domains but also for the target domain, which greatly reduces the domain gap and boosts the effect of XDED.

Based on the recent theoretical result on the relationship between the domain generalization and the flatness of local minima [8], we first empirically show that the proposed framework can improve generalization capability by achieving two goals: promoting flat minima and reducing the domain gap. Next, we further demonstrate the superiority of our method through extensive experimental results. On the standard public benchmarks for cross-domain image classification, XDED significantly enhances generalization ability in both multi-source and single-source settings. We also validate the effectiveness of our method in various domain generalization scenarios by showing the non-trivial improvement on the DomainBed [26], cross-dataset person re-ID [90,91], and cross-dataset semantic segmentation experiments. Moreover, we demonstrate that models learned by our method also help achieve robustness against adversarial attacks and unseen image corruptions.

2 Related Work

Domain generalization. The goal of domain generalization is to learn domain-invariant features that well generalize to unseen target domains. For the purpose, existing methods match feature distributions of different domains by adversarial feature alignment [48,49] or reducing the difference between feature distributions of diverse source domains [55]. Recently, meta-learning frameworks [2,16,47] have been introduced to simulate the domain shift by dividing the meta-train and meta-test domains from source domains. Also, data augmentation methods have been proposed to generate more diverse data beyond those of given source domains [37,41,66,77,93]. Most similar to our framework, ensemble methods for domain generalization have been proposed [79,65,94]. They all train multiple modules such as exemplar SVMs [79], domain-specific BN [34] layers [65] or classifiers [94], and exploit the ensemble of learned modules for prediction in testing. However, we remark that our XDED utilizes the ensemble of model predictions as the soft label and transfers it to the model itself. Therefore, it does not demand any additional module during both training and testing.

Knowledge distillation (KD). KD was originally studied to transfer the knowledge of a deep model to a shallow model for model compression [31]. It has been also used for other purposes such as metric learning [59,42] and

network regularization [78,87,84]. In particular for network regularization, self-knowledge distillation (self-KD) has been studied; it distills knowledge from the model itself and enforces prediction consistency between a sample and its perturbed one or other samples. KD has been used for domain adaptation [52,19], and such method trains several teacher models from the source domains and distills the ensemble of their predictions to the student model. It unfortunately requires large memory due to multiple teachers, and are difficult to be extended to domain generalization as they demand target images in training. In contrast, our method improves generalization capability of a model on unseen domains without the need for target images and additional teacher models.

Flat minima in loss landscapes. Recent analyses have revealed that finding flat minima is crucial for model generalization [38,18,20]. In this context, multiple methods have been proposed to promote flat minima in loss landscapes since flat minima have an advantage over sharp minima in robustness against the loss landscape shift between training and test data. Among literature on ways of promoting flat minima (*e.g.*, weight averaging [35,8] and training strategies [20,10]), we focus on the high entropy-seeking approaches, on which XDED is based. Maximum Entropy [60,9] maximizes the entropy of an output distribution from a classifier. Similarly, KD-based methods also aim at inducing high entropy of the output distribution by penalizing the mismatch with the output distribution from that of another classifier such as differently initialized peer networks [89] or subnetworks within a network itself [87]. Although SWAD [8] has introduced the importance of flat minima in the area of domain generalization, we remark that SWAD belongs to weight averaging but does not focus on learning domain-invariant features, whereas XDED belongs to entropy regularization as well as is designed for learning domain-invariant features.

Bias towards styles. Recent studies [6,22] revealed that DNNs overly depend on a strong bias towards styles, and it is also confirmed in the domain generalization literature [12,95,37] that a visual domain is highly correlated to feature statistics. Hence, previous work defines image styles as the bias and attempts to remove the bias by style augmentation in the space of feature statistics [95,37], using another model that is intentionally biased to styles [56], or minimizing a whitening loss [12]. Distinct from these techniques, we show that a simple yet effective de-stylization technique leads to a smaller divergence measure between target and source domains without bells and whistles.

3 Our Method

3.1 Cross-Domain Ensemble Distillation

Review of knowledge distillation (KD). The goal of KD [31] is to transfer knowledge of a teacher model t to a student model s , usually a wide and deep model to a smaller one, for the purposes of model compression or model regularization. Given input data point x and its label $y \in \{1, \dots, C\}$, we denote the output logit of model as $z(x) = [z_1(x), \dots, z_C(x)]$. The posterior predictive

distribution of x is then formulated as:

$$P(y|x; \theta, \tau) = \frac{\exp(z_y(x)/\tau)}{\sum_{i=1}^C \exp(z_i(x)/\tau)}, \quad (1)$$

where the model is parameterized by θ and τ is a temperature scaling parameter. KD enforces to match the predictive distributions of s and t . Specifically, it is achieved by minimizing the Kullback-Leibler (KL) divergence between their predictive distributions as follows:

$$\mathcal{L}_{\text{KD}}(X; \theta_s, \tau) = \sum_{x_i \in X} \sum_{c=1}^C D_{\text{KL}}(P(c|x_i; \theta_t, \tau) || P(c|x_i; \theta_s, \tau)), \quad (2)$$

where X is a batch of input data, θ_t and θ_s are the parameters of a teacher and a student, respectively.

Cross-domain ensemble distillation. We propose a new KD method for domain generalization, called cross-domain ensemble distillation (XDED). XDED aims to construct the domain-invariant knowledge from the data of multiple domains. Specifically, XDED generates an ensemble of logits from the data with the same label but from different domains. Next, XDED penalizes each logit for the mismatch with the ensemble which is not biased towards a specific domain, which encourages learning domain-invariant features. Unlike the conventional KD, XDED does not require an additional network that increases training complexity (*e.g.*, extra parameters and training time) but distills the ensemble constructed by multiple samples to the model itself in the form of self-KD.

Formally, let X_y denote the set of samples that have the same class label y in a mini-batch. Then, we obtain an ensemble of logits from X_y by simply taking an average as:

$$\bar{z}(X_y) = \sum_{x_i \in X_y} \frac{z(x_i)}{|X_y|}. \quad (3)$$

Then, the predictive distribution for the ensemble created from data X_y is as:

$$\bar{P}(c|X_y; \theta, \tau) = \frac{\exp(\bar{z}_c(X_y)/\tau)}{\sum_{i=1}^C \exp(\bar{z}_i(X_y)/\tau)}, \quad (4)$$

The loss function of XDED is defined as follows:

$$\mathcal{L}_{\text{XDED}}(X_y; \theta, \tau) = \sum_{x_i \in X_y} \sum_{c=1}^C D_{\text{KL}}(\bar{P}(c|X_y; \hat{\theta}, \tau) || P(c|x_i; \theta, \tau)), \quad (5)$$

where $\hat{\theta}$ is a fixed copy of the parameter θ . Following [53,84], we stop the gradient to be propagated through $\hat{\theta}$ to prevent the model from falling into some trivial solutions. To sum up, we set our objective function as

$$\min_{\theta} L_{\theta} = \mathcal{L}_{\text{CE}}(X, Y; \theta) + \lambda \sum_{y \in \{Y\}} \mathcal{L}_{\text{XDED}}(X_y; \theta, \tau), \quad (6)$$

where X is a batch of input images, Y is a batch of corresponding class labels, \mathcal{L}_{CE} denotes the vanilla cross-entropy loss, and λ is a hyperparameter to balance \mathcal{L}_{CE} and $\mathcal{L}_{\text{XDED}}$. λ and τ are 5.0 and 4.0 throughout this paper.

3.2 UniStyle: removing and unifying style bias

To further regularize the model to produce style-consistent predictions, we propose a de-stylization technique that is well-suited to domain generalization. As source domain styles are not expected to appear at test time, we propose UniStyle to prevent the model from being biased towards the domain-specific styles, which reduces the domain gap with the target domain.

More specifically, following existing methods based on style transfer [17,32,70], we first represent a neural style as statistics of intermediate feature maps from the feature extractor. Formally, let $F \in \mathbb{R}^{C \times H \times W}$ denote an intermediate feature map of an image. Then, a neural style of the image is represented as the combination of channel-wise mean $\mu(F) \in \mathbb{R}^C$ and standard deviation $\sigma(F) \in \mathbb{R}^C$ of F as:

$$\mu_c(F) = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W F_{c,h,w}, \quad (7)$$

and

$$\sigma_c(F) = \sqrt{\frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W (F_{c,h,w} - \mu_c(F))^2}, \quad (8)$$

where $\mu(F) = [\mu_1(F), \dots, \mu_C(F)]$ and $\sigma(F) = [\sigma_1(F), \dots, \sigma_C(F)]$. Next, we simply standardize each feature to have constant channel-wise statistics, μ_W and σ_W as:

$$\text{UniStyle}(F) = \sigma_W \frac{F - \mu(F)}{\sigma(F)} + \mu_W, \quad (9)$$

where $\mu_W = \mathbf{0}$ and $\sigma_W = \mathbf{1}$ (*i.e.*, zero-mean standardization). Technically, UniStyle is a special case of InstanceNorm (IN) [70]. Nevertheless, we remark that UniStyle aims to remove domain-specific information without any learnable parameters to reduce the domain gap while IN learns channel-wise scaling and bias parameters for style transfer. Also, note that we empirically observed that UniStyle is effective when being applied at multiple early layers, which is aligned with recent studies [17,32] suggesting that the style information is usually captured at the early layers.¹

3.3 Analysis of Our Method

In this section, we analyze the effectiveness of XDED, especially through the link to the theoretical result and the supporting empirical evidences. We first begin with a theorem related to domain adaptation [3,4], which shows that the

¹ See the supplementary material for further analyses.

Table 1: Comparison of the entropy values. When each model is converged, the entropy value is calculated by averaging over all training samples.

Methods	OfficeHome (Clipart)		PACS (Cartoon)	
	Entropy	Accuracy	Entropy	Accuracy
ResNet-18	0.25	49.4	0.01	75.9
MixStyle [95]	0.35	53.4	0.03	78.8
XDED	0.92	55.2	0.38	81.7

expected risk on the target domain is bounded by that on the source domain and the divergence between these domains. To find a model parameter $\theta \in \Theta$ for domain generalization, Cha *et al.* [8] considered a robust empirical loss:

$$\hat{\varepsilon}_S^\gamma(\theta) := \max_{\|\Delta\| \leq \gamma} \hat{\varepsilon}_S(\theta + \Delta) \quad (10)$$

where $\hat{\varepsilon}_S(\theta)$ is an empirical risk over source domains S and γ is a radius which defines neighbor parameters of θ . Then, Cha *et al.* [8] proved that finding flat minima reduces the domain gap through the theorem below:

Theorem 1. Consider a set of N covers $\{\Theta_k\}_{k=1}^N$ such that the hypothesis space $\Theta \subset \cup_k \Theta_k$ where $\text{diam}(\Theta) := \sup_{\theta, \theta' \in \Theta} \|\theta - \theta'\|_2$, $N := \lceil (\text{diam}(\Theta)/\gamma)^d \rceil$ and d is dimension of Θ . Let v_k be a VC dimension of each Θ_k . Then, for any $\theta \in \Theta$, the following bound holds with probability at least $1 - \delta$,

$$\varepsilon_T(\theta) < \hat{\varepsilon}_S^\gamma(\theta) + \frac{1}{2I} \sum_{i=1}^I \text{Div}(S_i, T) + \max_{k \in [1, N]} \sqrt{\frac{v_k \ln(m/v_k) + \ln(N/\delta)}{m}}, \quad (11)$$

where $m = nI$ is the number of training samples and $\text{Div}(S_i, T)$ is the divergence between the source domain S_i and the target domain T .

We remark that, in Eq. (11), the test loss $\varepsilon_T(\theta)$ is bounded by three terms: (1) the robust empirical loss $\hat{\varepsilon}_S^\gamma(\theta)$, (2) the divergence $\text{Div}(S_i, T)$, and (3) a confidence bound depending on the radius γ and the number of training samples m . In the rest of this section, according to the above theorem, we provide a theoretical interpretation that our method enhances the generalization ability by lowering both $\hat{\varepsilon}_S^\gamma(\theta)$ and $\text{Div}(S_i, T)$ with the empirical evidences.

Promoting flat minima. We remark that XDED is motivated by recent entropy regularization methods [9,87,89] in pursuit of flat minima. It has been empirically demonstrated that these methods promote flat minima by inducing higher posterior entropy. It can be interpreted as relaxing the training procedure to learn richer information encoded in soft labels, which helps the model converge to flat minima more than forcing the model to completely fit one-hot labels. In this context, we also demonstrate that XDED clearly induces higher entropy as shown in Table 1. Considering that XDED is motivated by the observation that different domains manifest different inter-class relations due to the domain gap (Fig. 1), this is natural since our ensembles would integrate meaningful inter-class relations from multiple domains and the model learned with them would be led towards high entropy.

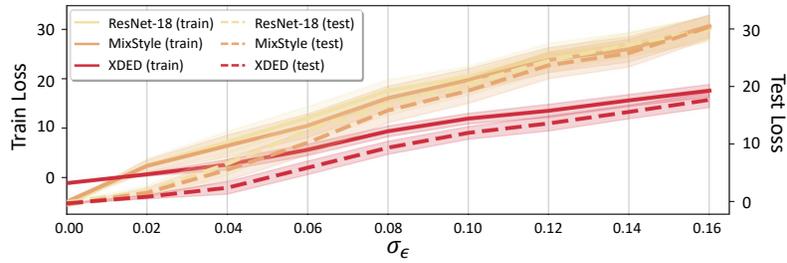


Fig. 2: Train/Test losses versus the weight perturbation while varying the standard deviation of the added Gaussian noise. Note that the results are produced with the target domain (Art of PACS) and the rest source domains, and the loss values are log-scaled.

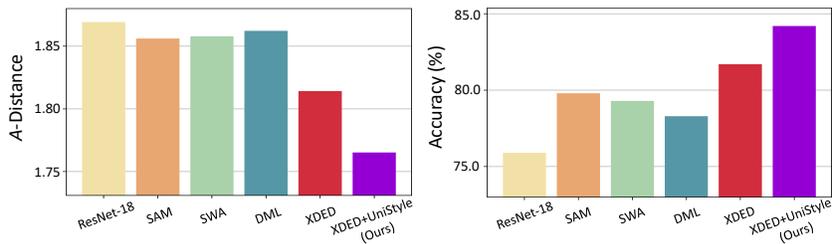


Fig. 3: Comparison to existing methods promoting flat minima. Each model is evaluated on Cartoon of PACS after being trained on the rest source domains. **Left:** The divergence (\mathcal{A} -distance) between the source domains and the target domain, **Right:** Generalization performance on the target domain.

Next, to investigate whether the model learned with XDED converges to flat minima indeed, we quantify the flatness of the local minima where the model converged by measuring the increase of loss values between θ and its neighborhoods, assuming that the model converged in flat minima would have smaller increases. Following [8,9,87,89], we measure the losses of the learned models before and after adding Gaussian noises to model parameters while varying the standard deviation of the noise σ_ϵ (*i.e.*, $\mathcal{L}_{\text{CE}}(X, Y; \theta + \epsilon)$ where $\epsilon \sim N(0, \sigma_\epsilon)$) with 100 runs. As a result, XDED demonstrates its robustness against the weight perturbation with smaller loss increases as shown in Fig. 2.

Domain-invariant feature learning. Here, we highlight that XDED also learns domain-invariant features via regularizing the consistency between the predictions from the data with the same label but from different domains and their ensemble. Thus, we compare XDED with existing methods promoting flat minima, which are dedicated to the flatness of local minima only. Specifically, to examine the effectiveness in reducing the divergence $\text{Div}(S_i, T)$, we measure \mathcal{A} -distance [3,39]. Due to the computational intractability, we calculated an ap-

proximated one [50,56]² As shown in Fig. 3 (Left), we observe that the existing methods promoting flat minima fail to reduce the distance while XDED clearly lowers the distance and UniStyle further enhances the result. Naturally, that result is connected to the quantitative superiority of our framework over existing flat minima-promoting methods (Fig. 3 (Right)).

4 Experiments

4.1 Generalization in image classification

Multi-source domain generalization. Specifically, for a fair comparison, we follow the leave-one-domain-out protocol [45] where we train a model on three domains and evaluate it on the remaining domain. For the benchmark datasets, we employ the PACS [45] and OfficeHome [72] that are widely-used benchmarks for domain generalization in image classification. PACS contains 9,991 images of 7 classes over 4 domains: Art Painting, Cartoon, Photo, and Sketch. OfficeHome includes 15,500 images of 65 classes over 4 domains: Artistic, Clipart, Product, and Real. We use ResNet-18 [27] as the backbone, and our UniStyle is applied to output feature maps of the first and second residual blocks for PACS and the first one only for OfficeHome.

Results. As summarized in Table. 2, we observe that our method not only significantly enhances the vanilla but also outperforms the latest competing methods. In particular, our method outperforms the second-best method on Cartoon of PACS and Clipart of OfficeHome by about 4.0% and 2.0%, respectively. these results justify the superiority of our method, which is simple yet effective.

Single-source domain generalization Thanks to the simple design of our proposed method, which does not explicitly require domain labels, our method can be transparently incorporated with single-source domain generalization where we only have access to a single source domain during training. Therefore, to further evaluate the impact of our method on single-source domain generalization, our model is trained on each single domain of PACS and evaluated on the remaining target domains.

Results. As shown in Table. 3, our model, on average, significantly outperforms other baselines by 8.7% in average accuracy. Besides, in all cases except for the case of $C \rightarrow S$, our model shows its superiority in performance. We believe this interesting result stems from the fact that our method is still able to help the model converge to flat minima and exploit the fine-grained relations between intra-domain samples even if only a single source domain is given.

DomainBed. We also conduct extensive experiments on the DomainBed [26] which is a testbed for domain generalization to compare state-of-the-art methods across several benchmark datasets. The rationale behind the DomainBed is that the domain generalization performances are too much dependent on the

² It is defined as $\hat{d}_{\mathcal{A}} = 2(1 - 2\epsilon_{\text{svm}})$ where ϵ_{svm} is the generalization error of a SVM-based two-class classifier trained to distinguish between target and source domains.

Table 2: Leave-one-domain-out generalization results on PACS and OfficeHome.

Methods	PACS					OfficeHome				
	Art	Cartoon	Photo	Sketch	Avg.	Artistic	Clipart	Product	Real	Avg.
ResNet-18	77.0	75.9	96.0	69.2	79.5	58.9	49.4	74.3	76.2	64.7
MMD-AE [48]	75.2	72.7	96.0	64.2	77.0	56.5	47.3	72.1	74.8	62.7
JiGen [7]	79.4	75.3	96.0	71.6	80.5	53.0	47.4	71.4	72.7	61.2
CrossGrad [66]	79.8	76.8	96.0	70.2	80.7	58.4	49.4	73.9	75.8	64.4
MASF [16]	80.2	77.1	94.9	71.6	81.0	-	-	-	-	-
Epi-FCR [47]	82.1	77.0	93.9	73.0	81.5	-	-	-	-	-
EISNet [74]	81.8	76.4	95.9	74.3	82.1	-	-	-	-	-
L2A-OT [93]	83.3	78.2	<u>96.2</u>	73.6	82.8	<u>60.6</u>	50.1	<u>74.8</u>	77.0	65.6
SagNet [56]	83.5	77.6	95.4	76.3	83.2	60.2	45.3	70.4	73.3	62.3
SelfReg [40]	82.3	78.4	<u>96.2</u>	77.4	83.6	-	-	-	-	-
MixStyle [95]	84.1	78.8	96.1	75.9	83.7	58.7	53.4	74.2	75.9	65.5
L2D [75]	81.4	79.5	95.5	80.5	84.2	-	-	-	-	-
FACT [77]	<u>85.3</u>	78.3	95.1	79.1	84.5	60.3	54.8	74.4	<u>76.5</u>	<u>66.5</u>
DSON [65]	84.6	77.6	95.8	<u>82.2</u>	85.1	59.3	45.7	71.8	74.6	62.9
RSC [33]	83.4	<u>80.3</u>	95.9	80.8	85.1	58.4	47.9	71.6	74.5	63.1
StyleNeophile [37]	84.4	79.2	94.9	83.2	<u>85.4</u>	59.5	<u>55.0</u>	73.5	75.5	65.8
Ours	85.6	84.2	96.5	79.1	86.4	60.8	57.1	75.3	<u>76.5</u>	67.4

Table 3: Single-source domain generalization accuracy (%) on PACS with a ResNet-18. (A: Art Painting, C: Cartoon, S:Sketch, P:Photo).

Methods	A \rightarrow C	A \rightarrow S	A \rightarrow P	C \rightarrow A	C \rightarrow S	C \rightarrow P	S \rightarrow A	S \rightarrow C	S \rightarrow P	P \rightarrow A	P \rightarrow C	P \rightarrow S	Avg.
ResNet-18	62.3	49.0	95.2	65.7	60.7	83.6	28.0	54.5	35.6	64.1	23.6	29.1	54.3
JiGen [7]	57.0	50.0	96.1	65.3	65.9	85.5	26.6	41.1	42.8	62.4	27.2	35.5	54.6
MixStyle [95]	65.5	49.8	<u>96.7</u>	69.9	64.5	85.3	27.1	50.9	32.6	67.7	<u>38.9</u>	39.1	57.4
RSC [33]	62.5	53.1	96.2	68.9	70.3	85.8	37.9	56.3	<u>47.4</u>	66.3	26.4	32.0	58.6
SelfReg [40]	65.2	55.9	96.6	72.0	<u>70.0</u>	<u>87.5</u>	37.1	54.0	46.0	67.7	28.9	33.7	59.5
SagNet [56]	<u>67.1</u>	<u>56.8</u>	95.7	<u>72.1</u>	69.2	85.7	<u>41.1</u>	<u>62.9</u>	46.2	<u>69.8</u>	35.1	<u>40.7</u>	61.9
Ours	74.6	58.1	96.8	74.4	69.6	87.6	43.3	65.6	50.3	71.4	54.3	51.5	66.5

hyperparameter tuning. For a fair comparison, we follow its standard protocols for training and evaluation.

Results. As shown in Table. 4, our method generally shows competitive performances and ranks second out of 15 methods on average accuracy. In particular, on CMNIST, our method substantially outperforms other competing methods. Since CMNIST is designed to simulate the domain shift by correlating the digit colors with the class labels, we conjecture that our improvement on CMNIST is attributed to the de-stylization effect of UniStyle, which would help the model decorrelate between the colors and labels.

4.2 Generalization in person re-ID

In this section, we further evaluate our method on person re-identification (re-ID), which is to match pedestrians across non-overlapping camera views.

Experimental setup. Here, we address domain generalization for person re-ID, where the test data is collected from cameras of the unseen dataset rather than

Table 4: Domain generalization accuracy (%) on DomainBed. The column ‘‘Terra’’ stands for TerraIncognita dataset. Note that we adopt leave-one-domain-out cross-validation as a model selection criteria.

Model selection: leave-one-domain-out cross-validation							
Methods	CMNIST	RMNIST	VLCS	PACS	OfficeHome	Terra	Avg.
ERM [71]	36.7	97.7	77.2	83.0	65.7	41.4	66.9
IRM [1]	40.3	97.0	76.3	81.5	64.3	41.2	66.7
GroupDRO [64]	36.8	97.6	77.9	83.5	65.2	44.9	66.7
Mixup [86]	33.4	97.8	77.7	83.2	67.0	48.7	67.9
MLDG [46]	36.7	97.6	77.2	82.9	66.1	46.2	67.7
CORAL [68]	39.7	97.8	78.7	82.6	68.5	46.3	68.9
MMD [48]	36.8	97.8	77.3	83.2	60.2	46.5	66.9
DANN [21]	40.7	97.6	76.9	81.0	64.9	44.4	67.5
CDANN [49]	39.1	97.5	77.5	78.8	64.3	39.9	66.1
MTL [5]	35.0	97.8	76.6	83.7	65.7	44.9	67.2
SagNet [56]	36.5	94.0	77.5	82.3	67.6	47.2	67.5
ARM [88]	36.8	98.1	76.6	81.7	64.4	42.6	66.7
VREx [44]	36.9	93.6	76.7	81.3	64.9	37.3	65.1
RSC [33]	36.5	97.6	77.5	82.6	65.8	40.0	66.6
Ours	46.5	97.7	74.8	83.8	65.0	42.5	68.4

Table 5: Generalization results on the cross-dataset person re-ID.

Methods	Market \rightarrow Duke		Duke \rightarrow Market	
	mAP	R@1	mAP	R@1
ResNet-50	19.3	35.4	20.4	45.2
RandomErase [92]	14.3	27.8	16.1	38.5
DropBlock [23]	18.2	33.2	19.7	45.3
MixStyle [95]	23.4	43.3	24.7	53.0
StyleNeophile [37]	26.3	46.5	27.2	55.0
Ours	27.4	49.3	30.1	59.0

from those of the training dataset. Specifically, the model trained to match people in the source dataset is evaluated by how well it matches pedestrian data of the unseen test set, which are disjoint from the source dataset. For datasets, we adopt two widely-used benchmarks: Market1501 (Market) [90] and DukeMTMC-reID (Duke) [62,91]. We use 32,668 images of 1,501 identities collected from 6 cameras and 36,411 images of 1,812 identities from 8 cameras for Market1501 and Duke, respectively. As for performance measures, we adopt mean average precision (mAP) and Recall@K (R@K). Following the prior work [95], we adopt ResNet-50 [27] as a backbone architecture. In these experiments, we apply UniStyle to the 1st, 2nd, and 3rd residual blocks of a model.

Comparison to other regularization methods. As shown in Table. 5, our method substantially outperforms other methods in mAP and Recall@1. Although RandomErase and Dropblock are effective for learning discriminative features, they fail to improve performance when encountering unseen domain data. Furthermore, by exploiting inter-class relations provided by different cam-

Table 6: mIoU (%) results on the cross-dataset semantic segmentation. GTA5 is for training, and Cityscapes, SYNTHIA, BDD, and Mapillary are test sets.

Methods (GTA5)	Cityscapes	BDD	Mapillary	SYNTHIA
DeepLabV3+ [11]	28.9	25.1	28.1	26.2
SW [58]	29.9	27.4	29.7	27.6
DRPC [82]	<u>37.4</u>	32.1	34.1	<u>28.0</u>
RobustNet [12]	36.5	35.2	40.3	28.3
Ours	39.2	<u>32.4</u>	<u>37.1</u>	<u>28.0</u>

Table 7: Ablation study of the proposed components on cross-domain tasks of image classification (Accuracy) and person re-ID (mAP).

Methods	Art	Clipart	Market → Duke
Vanilla	77.0	49.4	19.3
w/ UniStyle	81.2	50.4	<u>26.2</u>
w/ XDED	<u>83.3</u>	<u>55.2</u>	24.2
Ours	85.6	57.1	27.4

eras, our method shows its superiority over MixStyle and StyleNeophile which are designed for domain generalization but utilizes one-hot labels only.

4.3 Generalization in semantic segmentation

Experimental setup. Lastly, to investigate whether our method can be extended to the dense prediction task, evaluation on semantic segmentation is addressed here. Following the mainstream protocol, we train models on a synthetic dataset and evaluate them on several datasets which mainly belong to real-world. Specifically, we adopt GTA5 [61] as a source dataset which consists of 24,966 images. For target datasets, Cityscapes [13], BDD [81], and Mapillary [57] are real-world datasets whose image sizes are 5,000, 10,000, and 25,000, respectively. Lastly, SYNTHIA [63] has 9,400 images. Note that ResNet-50 is used as the backbone and the common 19 classes are used across all datasets.

Results. We remark that XDED constructs an ensemble by simply averaging all the logits from the pixels whose gt is the same in a mini-batch. As shown in Table 6, ours outperforms the competing methods overall, even if those are dedicated to this task only. We show that our method can be extended to the pixel-wise classification with little modification on XDED. Also, the results support our claim that our method is simple yet effective in a wide range of tasks.

4.4 In-depth Analysis

Ablation study. To investigate the impact of each component in our method, we conduct an ablation study which is summarized in Table 7. The result reveals that two components are complementary and consistently help the model improve the generalization ability. For image classification, XDED contributes most to the performance, and UniStyle boosts the effect of XDED. On both

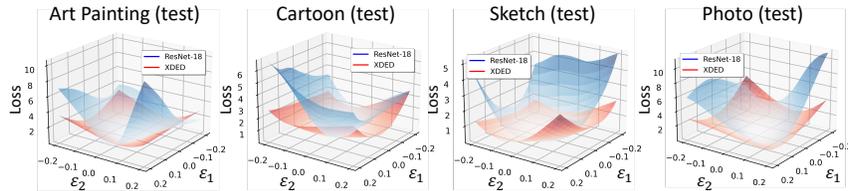


Fig. 4: Visualization results of the loss landscapes incorporating the vanilla method and XDED on the PACS dataset. Note that each loss landscape is visualized on the data of source domains, not the data of the marked target domain. Blue and red surfaces are from the vanilla method and XDED, respectively.

Table 8: Multi-source domain generalization accuracy (%) on Photo of PACS before and after applying given adversarial attacks.

Methods	Photo	w/ FGSM	w/ PGD
ResNet-18	96.0	39.6	16.3
Label smoothing [69]	95.6	43.5	20.2
Mixup [86]	95.8	46.5	21.9
Manifold mixup [73]	93.5	46.6	23.8
MixStyle [95]	96.1	41.4	22.7
Ours	96.5	55.4	30.4

domains, XDED uniformly improves the vanilla method by about 6%, whereas UniStyle shows different degrees of improvement. It is because the image style discrepancy between domains in OfficeHome is less severe than that in PACS. Interestingly, for the task of person re-ID, UniStyle reveals more impact than does XDED. Due to the inherent characteristics of the task itself, the effect of XDED on collecting meaningful knowledge of the same pedestrian from different cameras may become less significant.

Loss surface visualization. To further illustrate how XDED leads to flat minima in the loss landscapes, we provide qualitative results that visualize the loss landscapes. Following [9], we plot the loss landscapes on data of source domains per each case by perturbing the model parameters across the first and second Hessian eigenvectors which are provided by PyHessian [80] which is a framework for Hessian-based analysis of neural networks. As shown in Fig. 4, we observe that the loss landscapes incorporating XDED clearly become flatter than those incorporating the vanilla method for all cases. We argue that these qualitative results also consistently support that XDED promotes flat minima.

Robustness to adversarial examples. Recent studies have demonstrated that convergence on flat minima strengthens the adversarial robustness [76,67]. To revalidate that our method promotes flat minima, we evaluate the adversarial robustness of learned models. Specifically, we trained models on source domains and added adversarial perturbations on images of the unseen target domain by using existing adversarial attack methods: FGSM [25] and PGD [51]. Table 8 shows that our method outperforms other regularization methods in terms of

Table 9: Average classification error (%) on the corruption benchmarks.

Methods	CIFAR-10-C	CIFAR-100-C
40-2 WRN [85]	26.9	53.3
Cutout [15]	26.8	53.5
Mixup [86]	22.3	50.4
CutMix [83]	27.1	52.9
AutoAug [14]	23.9	49.6
AugMix [29]	11.2	35.9
Ours	<u>18.5</u>	<u>46.6</u>

robustness against both unseen data and adversarial attacks. Considering that adversarial attacks are made to maximize the loss value, we argue that our superiority in adversarial robustness is also attributed to the capability of promoting flat minima as desired, even though our method has no direct connection to adversarial training.

Results on corruption benchmarks. We further measure the resilience of learned models to image corruptions. Following the protocol provided by [28], we trained models on the original training dataset, and evaluated them on the test dataset constructed by corrupting the original test dataset through predefined corruption types. Table. 9 shows that our method outperforms all regularization methods except AugMix [30]. Considering AugMix is a state of the art that is dedicated to corruption robustness while ours is not, we argue that our method still shows its significant robustness against image corruptions.

5 Conclusion

We have presented a simple yet effective framework for domain generalization. XDED first generates an ensemble of output distributions for the data with the same label but from different domains, and then penalizes each output distribution for the mismatch with the ensemble in the form of self-knowledge distillation. With this approach, our model can learn domain-invariant features and also easily converges to flat minima. Besides, the proposed UniStyle suppresses domain-specific style bias to boost the effect of XDED and encourage style-consistent predictions. Furthermore, we empirically validate the generalization ability of the proposed method from the perspective of flat minima and reduced divergence between source and target. Through extensive experimental results, we demonstrate the superiority of the proposed framework.

Acknowledgement. This work was supported by the NRF grant and the IITP grant funded by Ministry of Science and ICT, Korea (NRF-2021R1A2C3012728, IITP-2019-0-01906, IITP-2022-0-00926, IITP-2022-0-00290).

References

1. Arjovsky, M., Bottou, L., Gulrajani, I., Lopez-Paz, D.: Invariant risk minimization. arXiv preprint arXiv:1907.02893 (2019) [11](#)
2. Balaji, Y., Sankaranarayanan, S., Chellappa, R.: Metareg: Towards domain generalization using meta-regularization. In: Proc. Neural Information Processing Systems (NeurIPS) (2018) [3](#)
3. Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., Vaughan, J.W.: A theory of learning from different domains. Machine learning (2010) [6](#), [8](#)
4. Ben-David, S., Blitzer, J., Crammer, K., Pereira, F., et al.: Analysis of representations for domain adaptation. In: Proc. Neural Information Processing Systems (NeurIPS) (2007) [6](#)
5. Blanchard, G., Deshmukh, A.A., Dogan, U., Lee, G., Scott, C.: Domain generalization by marginal transfer learning. Journal of Machine Learning Research (JMLR) (2021) [11](#)
6. Brendel, W., Bethge, M.: Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. In: Proc. International Conference on Learning Representations (ICLR) (2019) [4](#)
7. Carlucci, F.M., D’Innocente, A., Bucci, S., Caputo, B., Tommasi, T.: Domain generalization by solving jigsaw puzzles. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019) [10](#)
8. Cha, J., Chun, S., Lee, K., Cho, H.C., Park, S., Lee, Y., Park, S.: Swad: Domain generalization by seeking flat minima. In: Proc. Neural Information Processing Systems (NeurIPS) (2021) [3](#), [4](#), [7](#), [8](#)
9. Cha, S., Hsu, H., Hwang, T., Calmon, F.P., Moon, T.: Cpr: Classifier-projection regularization for continual learning. In: Proc. International Conference on Learning Representations (ICLR) (2021) [2](#), [4](#), [7](#), [8](#), [13](#)
10. Chaudhari, P., Choromanska, A., Soatto, S., LeCun, Y., Baldassi, C., Borgs, C., Chayes, J., Sagun, L., Zecchina, R.: Entropy-sgd: Biasing gradient descent into wide valleys. In: Proc. International Conference on Learning Representations (ICLR) (2017) [2](#), [4](#)
11. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Proc. European Conference on Computer Vision (ECCV) (2018) [12](#)
12. Choi, S., Jung, S., Yun, H., Kim, J.T., Kim, S., Choo, J.: Robustnet: Improving domain generalization in urban-scene segmentation via instance selective whitening. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2021) [4](#), [12](#)
13. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016) [12](#)
14. Cubuk, E.D., Zoph, B., Mane, D., Vasudevan, V., Le, Q.V.: Autoaugment: Learning augmentation policies from data (2019) [14](#)
15. DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. arXiv preprint arXiv:1708.04552 (2017) [14](#)
16. Dou, Q., Coelho de Castro, D., Kamnitsas, K., Glocker, B.: Domain generalization via model-agnostic learning of semantic features. In: Proc. Neural Information Processing Systems (NeurIPS) (2019) [3](#), [10](#)

17. Dumoulin, V., Shlens, J., Kudlur, M.: A learned representation for artistic style. In: Proc. International Conference on Learning Representations (ICLR) (2017) [6](#)
18. Dziugaite, G.K., Roy, D.M.: Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In: Proc. The Conference on Uncertainty in Artificial Intelligence (UAI) (2017) [2](#), [4](#)
19. Feng, H.Z., You, Z., Chen, M., Zhang, T., Zhu, M., Wu, F., Wu, C., Chen, W.: Kd3a: Unsupervised multi-source decentralized domain adaptation via knowledge distillation. In: Proc. International Conference on Machine Learning (ICML) (2021) [4](#)
20. Foret, P., Kleiner, A., Mobahi, H., Neyshabur, B.: Sharpness-aware minimization for efficiently improving generalization. In: Proc. International Conference on Learning Representations (ICLR) (2021) [2](#), [4](#)
21. Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., Lempitsky, V.: Domain-adversarial training of neural networks. *Journal of Machine Learning Research (JMLR)* (2016) [11](#)
22. Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F.A., Brendel, W.: Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In: Proc. International Conference on Learning Representations (ICLR) (2019) [4](#)
23. Ghiasi, G., Lin, T.Y., Le, Q.V.: Dropblock: A regularization method for convolutional networks. In: Proc. Neural Information Processing Systems (NeurIPS) (2018) [11](#)
24. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Proc. Neural Information Processing Systems (NeurIPS) (2014) [1](#)
25. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. In: Proc. International Conference on Learning Representations (ICLR) (2014) [13](#)
26. Gulrajani, I., Lopez-Paz, D.: In search of lost domain generalization. In: Proc. International Conference on Learning Representations (ICLR) (2021) [3](#), [9](#)
27. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2016) [9](#), [11](#)
28. Hendrycks, D., Dietterich, T.: Benchmarking neural network robustness to common corruptions and perturbations. In: Proc. International Conference on Learning Representations (ICLR) (2019) [14](#)
29. Hendrycks, D., Mu, N., Cubuk, E.D., Zoph, B., Gilmer, J., Lakshminarayanan, B.: Augmix: A simple data processing method to improve robustness and uncertainty (2020) [14](#)
30. Hendrycks, D., Mu, N., Cubuk, E.D., Zoph, B., Gilmer, J., Lakshminarayanan, B.: Augmix: A simple data processing method to improve robustness and uncertainty. In: Proc. International Conference on Learning Representations (ICLR) (2020) [14](#)
31. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015) [3](#), [4](#)
32. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2017) [6](#)
33. Huang, Z., Wang, H., Xing, E.P., Huang, D.: Self-challenging improves cross-domain generalization. In: Proc. European Conference on Computer Vision (ECCV) (2020) [10](#), [11](#)

34. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: Proc. International Conference on Machine Learning (ICML) (2015) [3](#)
35. Izmailov, P., Podoprikin, D., Garipov, T., Vetrov, D., Wilson, A.G.: Averaging weights leads to wider optima and better generalization. In: Proc. The Conference on Uncertainty in Artificial Intelligence (UAI) (2018) [4](#)
36. Jiang, Y., Neyshabur, B., Mobahi, H., Krishnan, D., Bengio, S.: Fantastic generalization measures and where to find them. In: Proc. International Conference on Learning Representations (ICLR) (2020) [2](#)
37. Kang, J., Lee, S., Kim, N., Kwak, S.: Style neophile: Constantly seeking novel styles for domain generalization. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2022) [3, 4, 10, 11](#)
38. Keskar, N.S., Mudigere, D., Nocedal, J., Smelyanskiy, M., Tang, P.T.P.: On large-batch training for deep learning: Generalization gap and sharp minima. In: Proc. International Conference on Learning Representations (ICLR) (2017) [2, 4](#)
39. Kifer, D., Ben-David, S., Gehrke, J.: Detecting change in data streams. In: In Very Large Databases (VLDB) (2004) [8](#)
40. Kim, D., Park, S., Kim, J., Lee, J.: Selfreg: Self-supervised contrastive regularization for domain generalization. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2021) [1, 10](#)
41. Kim, N., Son, T., Lan, C., Zeng, W., Kwak, S.: Wedge: Web-image assisted domain generalization for semantic segmentation. arXiv preprint arXiv:2109.14196 (2021) [3](#)
42. Kim, S., Kim, D., Cho, M., Kwak, S.: Embedding transfer with label relaxation for improved metric learning. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2021) [3](#)
43. Krizhevsky, A., Sutskever, I., Hinton, G.E.: ImageNet classification with deep convolutional neural networks. In: Proc. Neural Information Processing Systems (NeurIPS) (2012) [1](#)
44. Krueger, D., Caballero, E., Jacobsen, J.H., Zhang, A., Binas, J., Zhang, D., Le Priol, R., Courville, A.: Out-of-distribution generalization via risk extrapolation (rex). In: Proc. International Conference on Machine Learning (ICML) (2021) [11](#)
45. Li, D., Yang, Y., Song, Y.Z., Hospedales, T.M.: Deeper, broader and artier domain generalization. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2017) [9](#)
46. Li, D., Yang, Y., Song, Y.Z., Hospedales, T.M.: Learning to generalize: Meta-learning for domain generalization. In: Proc. AAAI Conference on Artificial Intelligence (AAAI) (2018) [11](#)
47. Li, D., Zhang, J., Yang, Y., Liu, C., Song, Y.Z., Hospedales, T.M.: Episodic training for domain generalization. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2019) [3, 10](#)
48. Li, H., Pan, S.J., Wang, S., Kot, A.C.: Domain generalization with adversarial feature learning. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018) [1, 3, 10, 11](#)
49. Li, Y., Tian, X., Gong, M., Liu, Y., Liu, T., Zhang, K., Tao, D.: Deep domain generalization via conditional invariant adversarial networks. In: Proc. European Conference on Computer Vision (ECCV) (2018) [1, 3, 11](#)
50. Long, M., Cao, Y., Wang, J., Jordan, M.: Learning transferable features with deep adaptation networks. In: Proc. International Conference on Machine Learning (ICML) (2015) [9](#)

51. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: Proc. International Conference on Learning Representations (ICLR) (2018) [13](#)
52. Meng, Z., Li, J., Gong, Y., Juang, B.H.: Adversarial teacher-student learning for unsupervised domain adaptation. In: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (2018) [4](#)
53. Miyato, T., Maeda, S.i., Koyama, M., Ishii, S.: Virtual adversarial training: a regularization method for supervised and semi-supervised learning. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) (2018) [5](#)
54. Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M.: Playing atari with deep reinforcement learning. In: NeurIPS Deep Learning Workshop (2013) [1](#)
55. Muandet, K., Balduzzi, D., Schölkopf, B.: Domain generalization via invariant feature representation. In: Proc. International Conference on Machine Learning (ICML) (2013) [1](#), [3](#)
56. Nam, H., Lee, H., Park, J., Yoon, W., Yoo, D.: Reducing domain gap by reducing style bias. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2021) [4](#), [9](#), [10](#), [11](#)
57. Neuhold, G., Ollmann, T., Rota Bulò, S., Kotschieder, P.: The mapillary vistas dataset for semantic understanding of street scenes. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2017) [12](#)
58. Pan, X., Zhan, X., Shi, J., Tang, X., Luo, P.: Switchable whitening for deep representation learning. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2019) [12](#)
59. Park, W., Kim, D., Lu, Y., Cho, M.: Relational knowledge distillation. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019) [3](#)
60. Pereyra, G., Tucker, G., Chorowski, J., Kaiser, L., Hinton, G.: Regularizing neural networks by penalizing confident output distributions. In: ICLR Workshop (2017) [2](#), [4](#)
61. Richter, S.R., Vineet, V., Roth, S., Koltun, V.: Playing for data: Ground truth from computer games. In: Proc. European Conference on Computer Vision (ECCV) (2016) [12](#)
62. Ristani, E., Solera, F., Zou, R., Cucchiara, R., Tomasi, C.: Performance measures and a data set for multi-target, multi-camera tracking. In: Proc. European Conference on Computer Vision (ECCV) (2016) [11](#)
63. Ros, G., Sellart, L., Materzynska, J., Vazquez, D., Lopez, A.M.: The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016) [12](#)
64. Sagawa, S., Koh, P.W., Hashimoto, T.B., Liang, P.: Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In: Proc. International Conference on Learning Representations (ICLR) (2020) [11](#)
65. Seo, S., Suh, Y., Kim, D., Kim, G., Han, J., Han, B.: Learning to optimize domain specific normalization for domain generalization. In: Proc. European Conference on Computer Vision (ECCV) (2020) [3](#), [10](#)
66. Shankar, S., Piratla, V., Chakrabarti, S., Chaudhuri, S., Jyothi, P., Sarawagi, S.: Generalizing across domains via cross-gradient training. In: Proc. International Conference on Learning Representations (ICLR) (2018) [3](#), [10](#)

67. Stutz, D., Hein, M., Schiele, B.: Relating adversarially robust generalization to flat minima. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2021) [13](#)
68. Sun, B., Saenko, K.: Deep coral: Correlation alignment for deep domain adaptation. In: Proc. European Conference on Computer Vision (ECCV) (2016) [11](#)
69. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016) [13](#)
70. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Instance normalization: The missing ingredient for fast stylization. arXiv preprint arXiv:1607.08022 (2016) [6](#)
71. Vapnik, V.: Statistical learning theory. NY: Wiley (1998) [11](#)
72. Venkateswara, H., Eusebio, J., Chakraborty, S., Panchanathan, S.: Deep hashing network for unsupervised domain adaptation. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017) [9](#)
73. Verma, V., Lamb, A., Beckham, C., Najafi, A., Mitliagkas, I., Lopez-Paz, D., Bengio, Y.: Manifold mixup: Better representations by interpolating hidden states. In: Proc. International Conference on Machine Learning (ICML) (2019) [13](#)
74. Wang, S., Yu, L., Li, C., Fu, C.W., Heng, P.A.: Learning from extrinsic and intrinsic supervisions for domain generalization. In: Proc. European Conference on Computer Vision (ECCV) (2020) [10](#)
75. Wang, Z., Luo, Y., Qiu, R., Huang, Z., Baktashmotlagh, M.: Learning to diversify for single domain generalization. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2021) [10](#)
76. Wu, D., Xia, S.T., Wang, Y.: Adversarial weight perturbation helps robust generalization. In: Proc. Neural Information Processing Systems (NeurIPS) (2020) [13](#)
77. Xu, Q., Zhang, R., Zhang, Y., Wang, Y., Tian, Q.: A fourier-based framework for domain generalization. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2021) [3](#), [10](#)
78. Xu, T.B., Liu, C.L.: Data-distortion guided self-distillation for deep neural networks. In: Proc. AAAI Conference on Artificial Intelligence (AAAI) (2019) [4](#)
79. Xu, Z., Li, W., Niu, L., Xu, D.: Exploiting low-rank structure from latent domains for domain generalization. In: Proc. European Conference on Computer Vision (ECCV) (2014) [3](#)
80. Yao, Z., Gholami, A., Keutzer, K., Mahoney, M.W.: Pyhessian: Neural networks through the lens of the hessian. In: 2020 IEEE International Conference on Big Data (Big Data) (2020) [13](#)
81. Yu, F., Xian, W., Chen, Y., Liu, F., Liao, M., Madhavan, V., Darrell, T.: Bdd100k: A diverse driving video database with scalable annotation tooling. arXiv preprint arXiv:1805.04687 (2018) [12](#)
82. Yue, X., Zhang, Y., Zhao, S., Sangiovanni-Vincentelli, A., Keutzer, K., Gong, B.: Domain randomization and pyramid consistency: Simulation-to-real generalization without accessing target domain data. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2019) [12](#)
83. Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2019) [14](#)
84. Yun, S., Park, J., Lee, K., Shin, J.: Regularizing class-wise predictions via self-knowledge distillation. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020) [4](#), [5](#)
85. Zagoruyko, S., Komodakis, N.: Wide residual networks. In: Proc. British Machine Vision Conference (BMVC) (2016) [14](#)

86. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. In: Proc. International Conference on Learning Representations (ICLR) (2018) [11](#), [13](#), [14](#)
87. Zhang, L., Song, J., Gao, A., Chen, J., Bao, C., Ma, K.: Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2019) [2](#), [4](#), [7](#), [8](#)
88. Zhang, M.M., Marklund, H., Dhawan, N., Gupta, A., Levine, S., Finn, C.: Adaptive risk minimization: A meta-learning approach for tackling group shift. In: Proc. Neural Information Processing Systems (NeurIPS) (2021) [11](#)
89. Zhang, Y., Xiang, T., Hospedales, T.M., Lu, H.: Deep mutual learning. In: Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018) [2](#), [4](#), [7](#), [8](#)
90. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2015) [3](#), [11](#)
91. Zheng, Z., Zheng, L., Yang, Y.: Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In: Proc. IEEE International Conference on Computer Vision (ICCV) (2017) [3](#), [11](#)
92. Zhong, Z., Zheng, L., Kang, G., Li, S., Yang, Y.: Random erasing data augmentation. In: Proc. AAAI Conference on Artificial Intelligence (AAAI) (2020) [11](#)
93. Zhou, K., Yang, Y., Hospedales, T., Xiang, T.: Learning to generate novel domains for domain generalization. In: Proc. European Conference on Computer Vision (ECCV) (2020) [3](#), [10](#)
94. Zhou, K., Yang, Y., Qiao, Y., Xiang, T.: Domain adaptive ensemble learning. IEEE Transactions on Image Processing (TIP) (2021) [3](#)
95. Zhou, K., Yang, Y., Qiao, Y., Xiang, T.: Domain generalization with mixstyle. In: Proc. International Conference on Learning Representations (ICLR) (2021) [4](#), [7](#), [10](#), [11](#), [13](#)