

NormAUG: Normalization-guided Augmentation for Domain Generalization

Lei Qi, Hongpeng Yang, Yinghuan Shi, Xin Geng

Abstract—Deep learning has made significant advancements in supervised learning. However, models trained in this setting often face challenges due to domain shift between training and test sets, resulting in a significant drop in performance during testing. To address this issue, several domain generalization methods have been developed to learn robust and domain-invariant features from multiple training domains that can generalize well to unseen test domains. Data augmentation plays a crucial role in achieving this goal by enhancing the diversity of the training data. In this paper, inspired by the observation that normalizing an image with different statistics generated by different batches with various domains can perturb its feature, we propose a simple yet effective method called NormAUG (Normalization-guided Augmentation). Our method includes two paths: the main path and the auxiliary (augmented) path. During training, the auxiliary path includes multiple sub-paths, each corresponding to batch normalization for a single domain or a random combination of multiple domains. This introduces diverse information at the feature level and improves the generalization of the main path. Moreover, our NormAUG method effectively reduces the existing upper boundary for generalization based on theoretical perspectives. During the test stage, we leverage an ensemble strategy to combine the predictions from the auxiliary path of our model, further boosting performance. Extensive experiments are conducted on multiple benchmark datasets to validate the effectiveness of our proposed method.

Index Terms—Normalization-guided augmentation, domain generalization, domain-shift.

I. INTRODUCTION

IN the last decade, deep learning has achieved significant success in various applications, including classification [1], [2], object detection [3], and semantic segmentation [4], [5]. Most existing methods are based on the assumption of iid (independent and identically distributed) data, where the training and test data are assumed to be from the same distribution. However, in real-world applications, this assumption is often violated due to domain shift between the training (source)

and test (target) domains. For instance, a model trained on photo images may not perform well on sketch images due to distributional variations, *i.e.*, domain generalization (DG). In the DG task, the unknown data distribution of the test data poses a challenge. If the training samples lack diversity, the trained model may overfit to the training data, hindering its generalization ability. Enriching the diversity of training data can be viewed as a way to simulate the distribution of diverse data, allowing the model to capture the characteristics of unseen test data during the training stage. It is worth noting that perturbing features is a technique employed to enhance the diversity of training data, as demonstrated in various studies [6], [7], [8].

Recently, several domain generalization methods have been developed to address this domain generalization issue [9], [10], [11], [12], [13], including augmentation-based methods, meta-learning-based methods, and domain alignment-based methods. For example, Zhang *et al.* [9] propose a multi-view regularized meta-learning algorithm that utilizes multiple optimization trajectories to determine an appropriate direction for model updating. During testing, this method employs an ensemble scheme to generate the final prediction. Moreover, Ding *et al.* [11] aim to explicitly remove domain-specific features for domain generalization, effectively achieving domain alignment. In contrast, Xu *et al.* [14] introduce a novel Fourier-based perspective for domain generalization. They exploit the fact that Fourier phase information contains high-level semantics and is less affected by domain shift. Among these methods, augmentation-based methods can effectively augment the training samples, mitigating the challenge of insufficient diversity in the domain generalization task. This intuitive technique enriches the training data and addresses the limitations posed by a lack of diverse samples.

In this paper, we propose a novel method for data augmentation in domain generalization, which differs from existing methods that perform augmentation at the image or feature level [14], [12], [17] to directly enrich image style information. Instead, we indirectly conduct data augmentation from a new perspective using batch normalization (BN). To visualize this concept, we perform an experiment by combining different domains, as depicted in Fig. 1. During this experiment, we keep all model parameters fixed, except for the statistics in the normalization layers. Specifically, the statistics in each batch normalization layer are computed using the current batch. As shown in Fig. 1, each image can be perturbed by normalizing it with different statistics derived from batches of different domains. Therefore, we can leverage and explore this observation in the context of domain generalization to enhance

The work is supported by NSFC Program (Grants No. 62206052, 62125602, 62076063), Jiangsu Natural Science Foundation Project (Grant No. BK20210224), and the Xplorer Prize.

Lei Qi is with the School of Computer Science and Engineering, Southeast University, National Center of Technology Innovation for EDA, and Key Laboratory of New Generation Artificial Intelligence Technology and Its Interdisciplinary Applications (Southeast University), Ministry of Education, China, 211189 (e-mail: qilei@seu.edu.cn).

Hongpeng Yang is with the School of Cyber Science and Engineering, Southeast University, Nanjing, China, 211189 (e-mail: hp_yang@seu.edu.cn).

Yinghuan Shi is with the State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China, 210023 (e-mail: syh@nju.edu.cn).

Xin Geng is with the School of Computer Science and Engineering, Southeast University, and Key Laboratory of New Generation Artificial Intelligence Technology and Its Interdisciplinary Applications (Southeast University), Ministry of Education, China, 211189 (e-mail: xgeng@seu.edu.cn).

Corresponding author: Xin Geng.

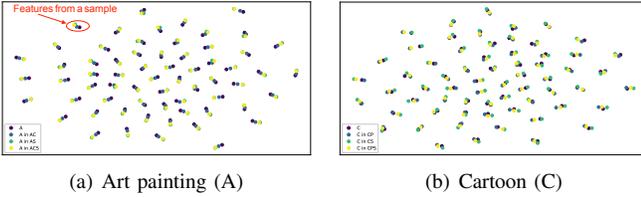


Fig. 1. Visualization of image’s features by t-SNE [15]. In this figure, “A”, “C”, “S”, and “P” represent images from the Art painting, Cartoon, Sketch, and Photo domains, respectively. For feature extraction, we utilize a ResNet-18 [1] model pre-trained on the ImageNet dataset [16]. It is important to note that we dynamically adjust the statistics (μ and σ) in all normalization layers during the feature extraction process. Each image is placed in different batches with different domains. For example, “A in ACS” indicates that images from the Art painting domain are combined with images from the Cartoon and Sketch domains for normalization during the test stage. As observed, this process perturbs the original feature representation.

the diversity of training samples.

Inspired by this observation, we propose a novel normalization-guided augmentation (NormAUG) method to enhance the model’s generalization. Our method consists of two paths: the main path and the augmented (or auxiliary) path. The main path is similar to the baseline method. The augmented path includes multiple sub-paths during training, with each sub-path representing batch normalization for a single domain or a random combination of multiple domains. This strategy effectively introduces diverse information into the feature representation. Importantly, these sub-paths in the auxiliary path can vary at each iteration as we randomly select sub-paths from a batch normalization (BN) bank. All paths and sub-paths in our method are implemented using different batch normalization layers, while other parameters are shared.

Additionally, we incorporate a classifier bank for the auxiliary path, where each classifier corresponds to a batch normalization (BN) layer in the BN bank. This method provides additional diverse information to our method. Leveraging the properties of our method, we combine the results from the auxiliary path with those from the main path to further enhance the model’s generalization. Extensive experiments demonstrate that our method outperforms state-of-the-art methods on various benchmark domain generalization datasets. Ablation analysis confirms the effectiveness of each module in our method. Furthermore, we analyze the effectiveness of NormAUG based on existing domain generalization theory, which reveals that our method achieves a lower generalization upper bound compared to the baseline method.

In this paper, our main contributions can be summarized as:

- We develop a novel Normalization-guided Augmentation (NormAUG) for domain generalization, which can effectively enhance the diversity of training data.
- We devise the BN bank and classifier bank to implement the proposed normalization-guided augmentation. This method not only enhances data diversity but also contributes to an improved ensemble prediction.
- Our method achieves state-of-the-art accuracy on multiple standard benchmark datasets, demonstrating its superiority over existing methods. We also provide an ablation study and further analysis to validate the effectiveness of our proposed method.

The structure of this paper is outlined as follows: Section II provides a literature review on relevant research. In Section III, we introduce our normalization-guided augmentation method. Section IV presents the experimental results and analysis. Finally, we conclude in Section V.

II. RELATED WORK

In this section, we review the most related domain generalization methods to our method, including data augmentation, ensemble learning and other methods. The following part presents a detailed investigation.

A. Data Augmentation

Since data augmentation can effectively enhance the diversity of training data, it has been recognized as a valuable method to improve the model’s generalization ability in unseen domains. In recent years, several methods have been developed from this perspective for the domain generalization task. For instance, Huang *et al.* [18] propose a simple training heuristic called Representation Self-Challenging (RSC) that significantly improves the generalization of convolutional neural networks (CNN) to out-of-domain data. RSC iteratively challenges the dominant features activated on the training data and encourages the network to activate remaining features that are more correlated with the labels. Another method is introduced by Xu *et al.* [14], who introduce a novel Fourier-based perspective for domain generalization. Their method leverages the assumption that the Fourier phase information contains high-level semantics and is not easily affected by domain shift, thus providing a robust representation for generalization across domains. In addition, Wang *et al.* [19] develop a style-complement module to enhance the generalization power of the model. This module synthesizes images from diverse distributions that are complementary to the source domain, thereby enriching the training data and improving the model’s ability to generalize to unseen domains.

Recently, there has been an increasing interest in augmentation methods based on Instance Normalization (IN), inspired by the AdaIN technique proposed by Huang *et al.* [20]. These methods aim to enhance the diversity and generalization of models through instance-level normalization. For example, Zhang *et al.* [17] propose a method called Exact Feature Distribution Matching (EFDM) that matches the empirical cumulative distribution functions of image features. This is achieved by applying exact histogram matching in the image feature space, enabling precise feature distribution alignment. Another method, introduced by Kang *et al.* [21], involves synthesizing novel styles continuously during training. They manage multiple queues to store observed styles and synthesize novel styles with distinct distributions compared to the styles in the queues. This method aims to enrich the style diversity and improve generalization. Additionally, MixStyle [12] explores a technique that probabilistically mixes instance-level feature statistics of training samples across source domains. By blending feature statistics, MixStyle encourages the model to learn more robust and domain-invariant representations, enhancing generalization performance.

Unlike the methods mentioned above, our method focuses on data augmentation using Batch Normalization (BN). By leveraging BN, we aim to effectively explore the diversity present in the training data. This allows us to enhance the generalization ability of our model by introducing variations in the normalization process.

B. Ensemble Learning

In the domain of domain generalization, ensemble learning has been widely utilized to improve prediction accuracy by leveraging multiple experts during the test process. For instance, Niu *et al.* [22] extend the multi-class SVM formulation to train a classifier for each class and latent domain, effectively integrating multiple classifiers to enhance generalization capability. Seo *et al.* [23] propose a simple yet effective multi-source domain generalization technique based on deep neural networks, incorporating optimized normalization layers specific to individual domains. Zhou *et al.* [10] introduce the domain adaptive ensemble learning (DAEL) framework, comprising a shared CNN feature extractor and multiple classifier heads trained to specialize in different source domains. Each classifier acts as an expert for its own domain and a non-expert for others. Segu *et al.* [24] train domain-dependent representations using ad-hoc batch normalization layers to collect independent domain statistics, enabling mapping of domains in a shared latent space. At test time, samples from an unknown domain are projected into this space to infer domain properties. Besides, Zhang *et al.* [9] employ a multi-view meta-learning scheme in the training stage and utilize an ensemble scheme by producing multiple test images for a sample to generate the final fused prediction.

In our work, our primary focus is on data augmentation rather than designing an ensemble scheme. It is important to note that the ensemble scheme serves as a supplementary component in the test stage of our method. Furthermore, unlike the method by Zhang *et al.* [9], our method does not require augmenting the test images during testing.

C. Other Methods

Besides the methods mentioned above, domain alignment or learning domain-invariant features is also crucial in domain generalization (DG) [25], [26], [27], [28], [29]. For example, Li *et al.* [25] propose learning features with domain-invariant class conditional distributions. Furthermore, meta-learning methods simulate training/test domain shift during training by synthesizing virtual test domains within each mini-batch [30], [31], [32], [33], [34]. For example, Zhang *et al.* [9] develop a multi-view regularized meta-learning algorithm employing multiple optimization trajectories. These methods aim to enhance model generalization by aligning features across all domains or using meta-learning techniques.

III. THE PROPOSED METHOD

In this paper, we are inspired by the observation depicted in Fig. 1, where the normalization of images from different domains introduces variations in their features. Motivated by this

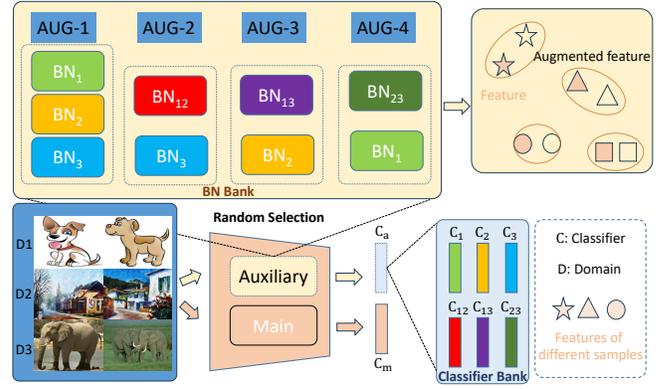


Fig. 2. The pipeline of our method in the training stage. In our training setup, we assume that there are three domains in the training set. The network architecture consists of two paths: the main path and the auxiliary path. Specifically, the data augmentation is performed through the auxiliary path. During each iteration, we randomly select a set of batch normalization (BN) layers from the BN bank (e.g., AUG-2 in the figure) and the corresponding classifier from the classifier bank to train the model. The main path and the auxiliary path share all parameters except for the BN layers. It is worth noting that BN layers with the same color in the BN bank share parameters, and all normalization layers in the auxiliary path are replaced by a BN bank. The corresponding algorithm is shown in Alg. 1.

observation, we propose a novel method called normalization-guided augmentation (NormAUG) for domain generalization, as illustrated in Fig. 2. Our method comprises two paths: the main path and the auxiliary path, which performs data augmentation using a batch normalization bank. Meanwhile, we also employ a classifier bank to better train our model. In the following sections, we will provide a detailed explanation of the background and our proposed method.

A. Background

Here, we will review the conventional batch normalization (BN) [35]. First, we define feature maps of an image $f_k \in \mathbb{R}^{C \times H \times W}$, where C are the number of channels, and H and W are the height and the width of feature maps. In general, BN leverages a global statistics of a batch to normalize all samples at each iteration, which can be defined as:

$$\text{BN}(f_k) = \gamma \frac{f_k - \mu}{\sigma} + \beta, \quad (1)$$

where $\gamma, \beta \in \mathbb{R}^C$ are learnable affine transformation parameters, and $\mu, \sigma \in \mathbb{R}^C$ (i.e., $\mu = [\mu_1, \dots, \mu_C]$ and $\sigma = [\sigma_1, \dots, \sigma_C]$) represent the channel-wise mean and standard deviation (i.e., statistics) of BN for feature maps. For statistics of the i -th channel are presented as:

$$\mu_i = \frac{1}{|\mathcal{B}|HW} \sum_{n \in \mathcal{B}} \sum_{h=1}^H \sum_{w=1}^W f[n, i, h, w], \quad (2)$$

$$\sigma_i = \sqrt{\frac{1}{|\mathcal{B}|HW} \sum_{n \in \mathcal{B}} \sum_{h=1}^H \sum_{w=1}^W (f[n, i, h, w] - \mu_i)^2 + \epsilon}, \quad (3)$$

where \mathcal{B} is a batch of samples, and $|\mathcal{B}|$ is the batch size. Besides, ϵ is a constant for numerical stability.

According to Eqs. 1 2 3, we can see that the normalized feature maps of an image are related to the statistics of a

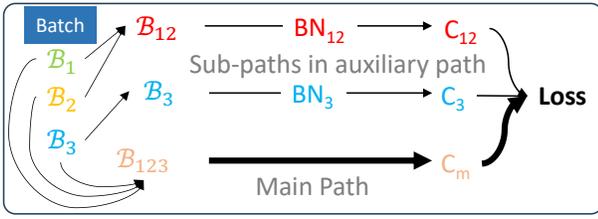


Fig. 3. The forward process of our method in the training stage. Here, we assume there are 3 source domains. This figure denotes that we randomly select the “AUG-2” in Fig. 2 from the BN bank at an iteration.

batch, meanwhile the statistics are decided by all samples of a batch. Therefore, if we randomly select one or multiple domain(s) to form a batch to perform the normalization, the diverse information can be introduced to implement data augmentation. We can also find this observation in Fig. 1.

B. NormAUG

Based on the analysis above, we propose a novel data augmentation scheme called NormAUG, which leverages the batch normalization (BN) perspective. Our model consists of two paths during the training stage: the main path and the auxiliary path. The main path serves as the baseline model, using ResNet-18 or ResNet-50 [1] in our implementation. The auxiliary path generates augmented information through the normalization-guided argumentation.

To be more specific, we randomly select an equal number of images from N domains to create a batch $\mathcal{B} = [\mathcal{B}_1; \dots; \mathcal{B}_N]$, where \mathcal{B}_i denotes the samples from the i -th domain in a batch. This batch is fed into both the main path and the auxiliary path. In the auxiliary path, we generate diverse information by randomly combining images from different domains and applying normalization. For instance, if we have 3 source domains, we can create four types of sub-batch combinations: $\{\mathcal{B}_1, \mathcal{B}_2, \mathcal{B}_3\}$, $\{\mathcal{B}_{12}, \mathcal{B}_3\}$, $\{\mathcal{B}_1, \mathcal{B}_{23}\}$, and $\{\mathcal{B}_2, \mathcal{B}_{13}\}$, for the auxiliary path. Here, \mathcal{B}_{12} is the sub-batch consisting of the samples from the first and second domains in a batch.

During the training process, our method employs a BN bank for each normalization layer in the auxiliary path, which consists of the corresponding four BN combinations: $\{\text{BN}_1, \text{BN}_2, \text{BN}_3\}$, $\{\text{BN}_{12}, \text{BN}_3\}$, $\{\text{BN}_1, \text{BN}_{23}\}$, and $\{\text{BN}_2, \text{BN}_{13}\}$, as shown in Fig. 2. When we generate the sub-batch set $\{\mathcal{B}_{12}, \mathcal{B}_3\}$, we feed it into the sub-paths using $\{\text{BN}_{12}, \text{BN}_3\}$ in each normalization layer. Besides, we feed \mathcal{B}_{123} into the main path. It is important to note that the main path and the auxiliary path share the same parameters, except for the BN layers. Therefore, our method does not introduce a large of extra parameters.

In this paper, we adopt multiple classifiers to train our model. The main path consists of an independent classifier, while in the augmented path, all BNs in the normalization bank have their respective independent classifiers. For instance, in the case of 3 source domains, we have the classifier C_m for the main path, and the classifier set $\{C_1, C_2, C_3, C_{12}, C_{13}, C_{23}\}$ for the augmented path. Each classifier corresponds to the BN set used in the respective sub-paths. The forward process of our method is illustrated in Fig. 3. The overall training loss can be defined as follows:

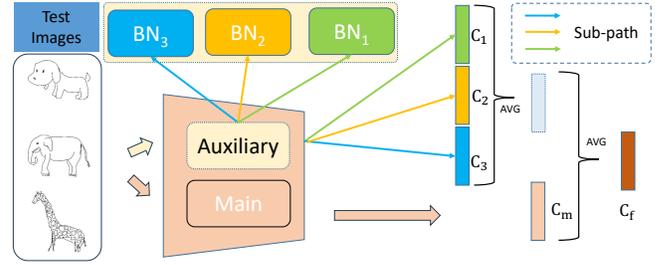


Fig. 4. The pipeline of our method in the test stage. In the test stage, we fuse the results from sub-paths (i.e., “AUG-1” in Fig. 2) to further improve the accuracy of the C_m . In this figure, “AVG” denotes the average operation.

$$\begin{aligned} \mathcal{L} = & \sum_{x \in \mathcal{B}} \sum_{c=1}^C -\log(P(c|x)) \cdot \mathbf{I}(x \in c) \\ & + \frac{1}{K} \sum_{k=1}^K \sum_{x \in \mathcal{B}_k} \sum_{c=1}^C -\log(P_k(c|x)) \cdot \mathbf{I}(x \in c), \end{aligned} \quad (4)$$

where K is the number of random sub-batches. $\mathbf{I}(x \in c)$ is equal to 1 if $x \in c$, otherwise 0. The detailed training process is described in Alg. 1.

C. Discussion

The scheme of the combination. Generally, for N domains, there are $1 + C_N^2 + C_N^3 \dots + C_N^{N-1}$ types of sub-batch combinations. Although adding BN layers does not bring a large of extra parameters, our method requires using the independent classifier for each BN, which results in the large classifier bank in the training process. Therefore, when there are N domains, we produce $N + 1$ types of sub-batch combinations. For example, we generate 5 types of sub-batch combinations: $\{\mathcal{B}_1, \mathcal{B}_2, \mathcal{B}_3, \mathcal{B}_4\}$, $\{\mathcal{B}_{123}, \mathcal{B}_4\}$, $\{\mathcal{B}_{124}, \mathcal{B}_3\}$, $\{\mathcal{B}_{134}, \mathcal{B}_2\}$, and $\{\mathcal{B}_{234}, \mathcal{B}_1\}$ for 4 domains.

The ensemble prediction. In addition, while our method aims to improve the accuracy of the main path through the auxiliary path, we can further enhance the prediction by employing a fusion scheme. By considering the diversity among the sub-paths in the auxiliary path, we select only the single-domain path to obtain an averaged result for the augmented path. The final prediction is then obtained by averaging the results from both the main path and the augmented path. This fusion process is illustrated in Fig. 4, where the combined result provides an improved prediction. The detailed test process is described in Alg. 2.

The other tricks. Instance normalization (IN) is known for its ability to remove domain-specific (style) information and improve the generalization of models. Therefore, in the main path of our method, we combine instance normalization (IN) with batch normalization (BN) to further enhance the model’s generalization. Specifically, in this paper, we utilize optimized normalization [23] in the main path of our model.

The benefits of the classifier bank. The advantages of using the classifier bank in the auxiliary path can be summarized as follows: 1) Enhanced Data Diversity: By using a single classifier for each BN in the auxiliary path, we can effectively

explore diverse information, as each BN is free to adapt to the specific domain characteristics. This flexibility allows us to avoid any underlying constraints during training and improve the model’s generalization. 2) Improved Ensemble Prediction: The independent classifiers can also be leveraged for the fused prediction in the test stage. This ensemble scheme further enhances the overall performance and will be validated in our experiments. By employing the classifier bank in the auxiliary path, our method can capitalize on these benefits, leading to superior results in domain generalization tasks.

Comparison with the most relevant methods. We compare our NormAUG method with two closely related methods, namely MixStyle [12], DSON [23], and BEN [24]. Compared to MixStyle, there are three key differences: 1) The underlying technique is different. Our NormAUG method is based on batch normalization, while MixStyle is based on instance normalization. 2) Our method employs a BN bank, which allows for the generation of more diverse samples, whereas MixStyle always mixes the styles of two images. 3) Our NormAUG method can produce an ensemble result in the test stage, thanks to its inherent properties, while MixStyle does not have this capability. In comparison to DSON, there are two main differences: 1) While DSON has independent paths for each domain, our NormAUG method includes multiple sub-paths in the auxiliary path, such as combinations of different domains. 2) The goals of the two methods differ. Our method has a main path, and our objective is to utilize the auxiliary path to enhance the generalization of the main path. On the other hand, DSON aims to learn multiple experts to improve prediction performance in unseen domains. It is worth noting that our NormAUG method, without the ensemble prediction in the test stage, can also achieve excellent performance, as demonstrated in Table VI. This is attributed to the diverse features obtained from the auxiliary path during training. In addition, BEN [24] is trained by independently evaluating domains to obtain feature embeddings of source domains. It then calculates the distance function between unknown domain samples and the source domain to linearly weight the representation of unknown domain samples during testing. Our method employs the BN bank and classifier bank for indirectly enhancing the diverse information in the training stage, and fuses different sub-paths to achieve the final prediction in the testing stage.

Explanation for NormAUG via existing theory. In this section, we use the domain generalization error bound [36] to demonstrate the effectiveness of our method. Firstly, we review the domain generalization error bound, and then we analyze our method based on it.

Theorem 1 [36], [37] (Domain generalization error bound): Let $\gamma := \min_{\pi \in \Delta_M} d_{\mathcal{H}}(\mathcal{P}_X^t, \sum_{i=1}^M \pi_i \mathcal{P}_X^i)$ ¹ with minimizer π^* being the distance of \mathcal{P}_X^t from the convex hull Λ . Let $\mathcal{P}_X^* := \sum_{i=1}^M \pi_i^* \mathcal{P}_X^i$ be the best approximator within Λ . Let $\rho := \sup_{\mathcal{P}_X', \mathcal{P}_X'' \in \Lambda} d_{\mathcal{H}}(\mathcal{P}_X', \mathcal{P}_X'')$ be the diameter of Λ . Then it holds that

¹ M is the number of source domains.

Algorithm 1: The training process of our NormAUG

Input: Training samples X_{tr} and labels Y .
Output: The trained model (θ).
1 $\theta \leftarrow$ Initialize by ResNet pre-trained on ImageNet.
// The number of epochs is T .
2 **for** $epoch \in [1, \dots, T]$ **do**
// The number of iterations in each epoch is N .
3 **for** $iteration \in [1, \dots, N]$ **do**
4 Randomly select a combination of the normalization and classifier for the auxiliary path as shown in Fig. 2.
5 Feed these input images into the main path and the auxiliary path, respectively.
6 Compute the whole loss as Eq. 4.
7 Update the model parameter θ .
8 **end**
9 **end**
10 **return**

Algorithm 2: The test process of our NormAUG

Input: Test samples X_{te} .
Output: The predicted result.
1 $\theta \leftarrow$ Initialize by the trained model θ .
// The number of iterations in the test set is N_{te} .
2 **for** $iteration \in [1, \dots, N_{te}]$ **do**
3 Feed the test images into the main path and the auxiliary path.
4 Obtain the result from the main path (p_m), and the result set from the auxiliary path (P_a).
5 Compute the average result of sub-paths in the auxiliary path as \bar{p}_a .
6 Compute the average of p_m and \bar{p}_a as p_f .
7 **end**
8 **return**

$$\epsilon^t(h) \leq \sum_{i=1}^M \pi_i^* \epsilon^i(h) + \frac{\gamma + \rho}{2} + \lambda_{\mathcal{H}}(\mathcal{P}_X^t, \mathcal{P}_X^*), \quad (5)$$

where $\lambda_{\mathcal{H}}(\mathcal{P}_X^t, \mathcal{P}_X^*)$ is the ideal joint risk across the target domain and the training domain (\mathcal{P}_X^*) with the most similar distribution to the target domain. In Theorem 1, the last item can be treated as a constant because it represents the ideal joint risk across the target domain and the training domain (\mathcal{P}_X^*) with the most similar distribution to the target domain. Besides, the first item is the empirical risk error of all source domains. Most of existing methods use the cross-entropy loss to train the model, thus it can also be viewed as a constant in the upper bound. Therefore, we primarily focus on analyzing the second item (which consist of γ and ρ).

- γ represents the discrepancy between the combination of all training domains and the target domain. In the domain generalization setting, if the test domain is far from the training domain in terms of distribution, the model’s generalization may be poor for all test samples. Particularly,

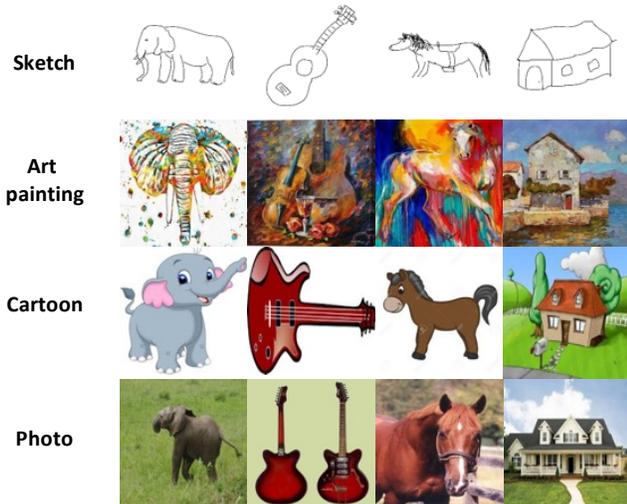


Fig. 5. Examples from PACS. We show some images from Sketch, Art painting, Cartoon, and Photo. As seen in this figure, the difference is obvious for these images with the same class from different domains.

our method utilizes the normalization based augmentation to enrich the diversity, resulting in obtaining the domain-invariant feature. We compute the divergence between the source and target domains, as shown in Tab. IX of the experimental section. From this table, we observe that the model with our augmentation generates a smaller domain gap between source and target domains than the model without our augmentation. Therefore, introducing our NormAUG can be beneficial in reducing the influence of the domain-shift between training and test sets and effectively mitigating the aforementioned risk.

- ρ indicates the maximum distance between different source domains. In our method, the normalization based augmentation scheme preserves semantic information while not introducing additional information. This indicates that *indirectly* generating diverse style information in our method does not create a large domain gap between training samples. Additionally, we can also observe this fact in Fig. 1. In summary, our method has the advantage of reducing the generalization error bound from the second item in Eq. 5.

IV. EXPERIMENTS

In this section, we begin by presenting the experimental datasets and configurations in Section IV-A. Following that, we evaluate our proposed method against the current state-of-the-art domain generalization techniques in Section IV-B. To verify the impact of different modules in our framework, we carry out ablation studies in Section IV-C. Finally, we delve deeper into the properties of our method in Section IV-D.

A. Datasets and Experimental Settings

1) *Datasets*: In this paper, we conduct the experiments to validate the effectiveness of our method on five benchmark DG datasets as follows:

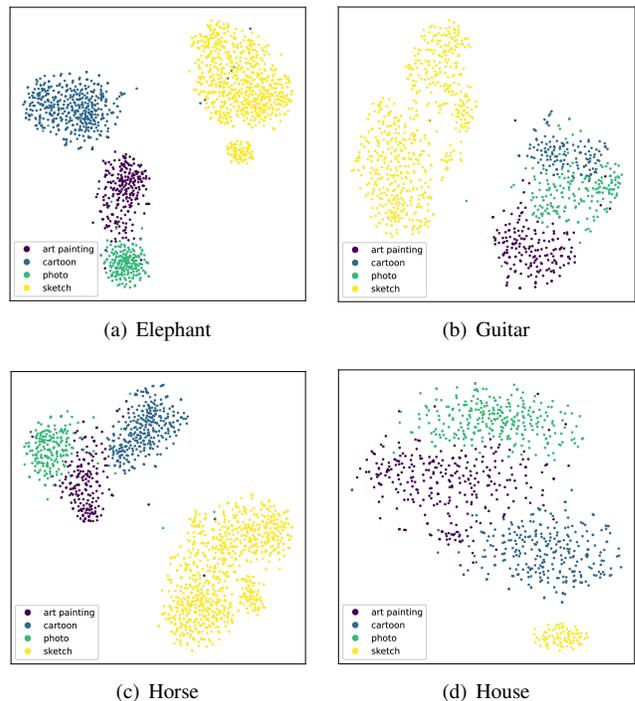


Fig. 6. Visualization of image’s features by t-SNE [15]. In this figure, we extract the image’s features using ResNet-18 [1] pre-trained on ImageNet to show the domain gap from the feature representation view. As observed in this figure, these images from the same class are in different positions, which can obviously show the domain’s difference as Fig. 5.

- **PACS** [38] consists of four different domains: Photo, Art painting, Cartoon and Sketch. It contains 9,991 images with 7 object categories in total, including Photo (1,670 images), Art (2,048 images), Cartoon (2,344 images), and Sketch (3,929 images).
- **Office-Home** [39] contains 15,588 images of 65 categories of office and home objects. It has four different domains namely Art (2,427 images), Clipart (4,365 images), Product (4,439 images) and Real World (4,357 images), which is originally introduced for UDA but is also applicable in the DG setting.
- **DomainNet** [40] is an extensive domain generalization dataset, comprising 596,010 images distributed across 345 categories from 6 distinct domains: Clipart (48,837 images), Infograph (53,201 images), Painting (75,759 images), Quickdraw (172,500 images), Real (175,327 images), and Sketch (70,386 images).
- **mini-DomainNet** [10] takes a subset of DomainNet [40]. mini-DomainNet includes four domains and 126 classes. As a result, mini-DomainNet contains 18,703 images of Clipart, 31,202 images of Painting, 65,609 images of Real and 24,492 images of Sketch.
- **DigitsDG** [41] is a digit recognition benchmark consisting of four classical datasets MNIST [42], MNIST-M [43], SVHN [44], SYN [43]. The four datasets mainly differ in font style, background and image quality. We use the original train/validation split in [41] with 600 images per class per dataset.

We show some examples from four different domains on PACS Fig. 5. As seen, there is an obvious difference among

different domains. Besides, we also visualize the features of four categories on PACS by t-SNE [15], as illustrated in Fig. 6. In this figure, different colors denote different domains. We observe that different domains appear in different spaces, validating that there exists the domain shift in the training set.

2) *Implementation Details*: In this study, we utilize ResNet-18 [1] and ResNet-50 [1] models pretrained on the ImageNet [16] dataset as the backbone for our framework. All images are resized to dimensions of 224×224 . We randomly sample 16 images from each domain to form a batch of data for input to our network. During training, we apply data augmentations including horizontal flipping, random cropping, color jittering. We use the SGD optimizer for both the classifier and backbone networks. By default, we set the initial learning rate (lr_c) for the classifier as 0.01, and the initial learning rate (lr_b) for the backbone to 0.003. Especially, for the Office-Home dataset, we use lr_c and lr_b as 0.005 and 0.001, respectively. We follow the standard data splits and adopt the leave-one-domain-out evaluation protocol as used in [45]. For evaluating the performance on the target domain, we select the model from the final training epoch. The accuracy on the target domain is reported and averaged over three runs. We employ the same settings for experiments on all datasets.

B. Comparison with State-of-the-art Methods

In this section, we compare our NormAUG method with several state-of-the-art (SOTA) methods on four benchmark datasets: PACS, Office-Home, mini-DomainNet, Digits-DG, and DomainNet. The experimental results are reported in Tabs. I, II, III, IV, and V. In the following part, we will give the detailed analysis.

Results on PACS. We compare our NormAUG method with several augmentation-based methods, including EFDMix [17], FACT [14], RSC [18], and STNP [21]. The experimental results are shown in Tab. I. As observed, our NormAUG consistently outperforms these methods on both ResNet-18 and ResNet-50. Additionally, we compare our method with ensemble methods such as DAEL [10] and DSON [23], demonstrating the superiority of our method. Furthermore, we compare our method with other state-of-the-art (SOTA) methods. For instance, our NormAUG outperforms MVDG by +0.48% (87.04 vs. 86.56) on ResNet-18. It is worth noting that MVDG generates multiple images through data augmentation for each test image and combines all the results for final prediction, which is also an ensemble scheme in the test stage.

Results on Office-Home. We also present the experimental results on the Office-Home dataset, as shown in Tab. II. As observed, our NormAUG method consistently outperforms the compared augmentation-based methods, including RSC [23], MixStyle [12], L2A-OT [53], FACT [14], DSU [54], and STNP [21]. Additionally, we compare our method with the recent method DCG [34], further highlighting the effectiveness of our method.

Results on mini-DomainNet. We further evaluate the effectiveness of our method on the mini-DomainNet dataset, which has a larger number of categories compared to PACS and Office-Home. In Tab. III, we report the experimental

TABLE I
DOMAIN GENERALIZATION ACCURACY (%) ON PACS DATASET WITH RESNET-18 (TOP) AND RESNET-50 (BOTTOM) BACKBONE. THE BEST PERFORMANCE IS MARKED AS **BOLD**, AND THE UNDERLINE IS THE SECOND BEST RESULT.

Methods	A	C	P	S	Avg.
COMEN [46]	82.60	81.00	94.60	84.50	85.68
CIRL [47]	86.08	80.59	95.93	82.67	86.32
I ² -ADR [28]	82.90	80.80	95.00	83.50	85.55
XDED [29]	85.60	84.20	<u>96.50</u>	79.10	86.35
LRDG [11]	81.88	80.20	95.21	84.65	85.48
DAEL [10]	84.60	74.40	95.60	78.90	83.40
DSON [23]	84.67	77.65	95.87	82.23	85.11
BEN [24]	78.80	78.90	94.80	79.70	83.10
MVDG [9]	<u>85.62</u>	79.98	95.54	85.08	<u>86.56</u>
MixStyle [12]	84.10	78.80	96.10	75.90	83.70
EFDMix [17]	83.90	79.40	96.80	75.00	83.78
FACT [14]	85.37	78.38	95.15	79.15	84.51
RSC [18]	83.43	80.31	95.99	80.85	85.15
STNP [21]	84.41	79.25	94.93	83.27	85.47
NormAUG (ours)	85.60	81.85	95.70	<u>85.00</u>	87.04
LRDG [11]	86.57	<u>85.78</u>	95.57	86.59	88.63
Fishr [48]	88.40	78.70	97.00	77.80	85.50
mDSDI [49]	87.70	80.40	98.10	78.40	86.20
GIN [50]	89.00	81.50	<u>98.00</u>	80.20	87.20
CACE-D [51]	89.20	82.10	<u>98.00</u>	80.50	87.50
I ² -ADR [28]	88.50	83.20	95.20	85.80	88.18
SWAD [52]	89.30	83.40	97.30	82.50	88.10
DSON [23]	87.04	80.62	95.99	82.90	86.64
MVDG [9]	89.31	84.22	97.43	<u>86.36</u>	89.33
RSC [18]	87.89	82.16	97.92	83.35	87.83
FACT [14]	89.63	81.77	96.75	84.46	88.15
STNP [21]	<u>90.35</u>	84.20	96.73	85.18	89.12
EFDMix [17]	90.60	82.50	98.10	76.40	86.90
NormAUG (ours)	88.95	86.00	97.15	85.95	89.51

TABLE II
DOMAIN GENERALIZATION ACCURACY (%) ON OFFICE-HOME. THE BEST PERFORMANCE IS MARKED AS **BOLD**, AND THE UNDERLINE IS THE SECOND BEST RESULT.

Methods	A	C	P	R	Avg.
RSC [23]	58.42	47.90	71.63	74.54	63.12
MixStyle [12]	58.70	53.40	74.20	75.90	65.55
L2A-OT [53]	60.60	50.10	74.80	77.00	65.63
FACT [14]	60.34	54.85	74.48	76.55	66.56
DSU [54]	60.20	54.80	74.10	75.10	66.05
STNP [21]	59.55	55.01	73.57	75.52	65.91
DAEL [10]	59.40	55.10	74.00	75.70	66.10
DCG [34]	60.67	<u>55.46</u>	75.26	76.82	<u>67.05</u>
NormAUG (ours)	61.25	58.00	<u>75.25</u>	<u>76.65</u>	67.79

results and compare our method with various methods, including meta-learning, domain-invariant, augmentation, and other methods. As observed, our NormAUG method exhibits a clear advantage over all the compared methods, demonstrating its effectiveness on datasets with multiple classes.

Results on Digits-DG. We also evaluate the performance of our method on the Digits-DG dataset, which is a widely used benchmark dataset for domain generalization in digit recognition. The experimental results are presented in Tab. IV. As observed in the table, our method consistently outperforms the other methods, highlighting its effectiveness in the domain generalization task for digit recognition. The superior performance of our method on this dataset further validates its robustness and generalization capability.

Results on DomainNet. We evaluate the performance of our method on the DomainNet dataset, which has more

TABLE III

DOMAIN GENERALIZATION ACCURACY (%) ON MINI-DOMAINNET. THE BEST PERFORMANCE IS MARKED AS **BOLD**, AND THE UNDERLINE IS THE SECOND BEST RESULT.

Methods	C	P	R	S	Avg.
MLDG [32]	65.70	57.00	63.70	58.10	61.12
DAEL [10]	69.95	55.13	66.11	55.72	61.73
MMD [55]	65.00	58.00	63.80	58.40	61.30
Mixup [56]	67.10	59.10	64.30	59.20	62.43
SagNet [57]	65.00	58.10	64.20	58.10	61.35
CORAL [40]	66.50	59.50	66.00	59.50	62.87
MTL [58]	65.30	59.00	65.60	58.50	62.10
DCG [34]	69.38	61.79	66.34	<u>63.21</u>	<u>65.18</u>
NormAUG (ours)	70.20	66.90	71.20	63.40	67.93

TABLE IV

DOMAIN GENERALIZATION ACCURACY (%) ON DIGITS-DG. THE BEST PERFORMANCE IS MARKED AS **BOLD**, AND THE UNDERLINE IS THE SECOND BEST RESULT.

Methods	MN	MN-M	SV	SY	Avg.
FACT [14]	97.90	65.60	72.40	90.30	81.55
COMEN [46]	97.10	67.60	75.10	91.30	82.78
CIRL [47]	96.08	69.87	<u>76.17</u>	87.68	82.45
STEAM [59]	96.80	67.50	76.00	<u>92.20</u>	<u>83.13</u>
NormAUG (ours)	<u>97.50</u>	<u>67.65</u>	80.10	96.85	85.53

domains compared to the other datasets in our experiment. The results are reported in Tab. V. We can observe that on large dataset with multiple source domains, our method exhibits a remarkably substantial improvement in performance.

C. Ablation Study

We conducted an ablation study on the PACS and mini-DomainNet datasets to analyze the effectiveness of different modules. The experimental results are presented in Tab. VI. In this table, “ON” represents the optimized normalization [23] module in the main path, “AUG” represents the normalization-based augmentation module in the training stage, and “EP” represents the ensemble prediction based on the normalization-based augmentation module.

As observed in the table, the “ON” module is effective for the domain generalization task as it incorporates the instance normalization scheme within the normalization layer. It contributes to improved performance. Furthermore, the addition of the “AUG” module enhances the diversity of samples and further boosts the performance. Finally, the combination of the ensemble prediction (“EP”) in the test stage leads to additional performance improvement.

Overall, the results demonstrate the effectiveness of each module and the benefits of combining them for improved domain generalization performance.

D. Further Analysis

Impact on different numbers of classifiers. As mentioned in Section III, we explore different configurations for the classifiers in the training process. Specifically, we consider two cases: one shared classifier and two classifiers. In the case of two classifiers, one is assigned to the main path, while the other is assigned to the auxiliary path (*i.e.*, all sub-paths in the auxiliary path share a classifier). It is important to note that the

TABLE V

DOMAIN GENERALIZATION ACCURACY (%) ON DOMAINNET WITH RESNET-18 (TOP) AND RESNET-50 (BOTTOM) BACKBONE. THE BEST PERFORMANCE IS MARKED AS **BOLD**, AND THE UNDERLINE IS THE SECOND BEST RESULT.

Methods	C	I	P	Q	R	S	Avg.
MetaReg [31]	53.70	<u>21.10</u>	<u>45.30</u>	10.60	58.50	42.30	38.58
DMG [27]	60.10	18.80	44.50	<u>14.20</u>	54.70	41.70	39.00
I ² -ADR [28]	57.30	15.20	44.10	12.10	53.90	<u>46.70</u>	38.22
ITTA [60]	50.70	13.90	39.40	11.90	50.20	43.50	34.90
NormAUG (ours)	<u>57.40</u>	22.70	49.00	14.60	<u>58.30</u>	48.70	41.78
RSC [18]	55.00	18.30	44.40	12.20	55.70	47.80	38.90
DMG [27]	65.20	<u>22.20</u>	50.00	15.70	59.60	49.00	43.62
SagNet [57]	57.70	19.00	45.30	12.70	58.10	48.80	40.27
SelfReg [61]	60.70	21.60	49.40	12.70	60.70	51.70	42.80
I ² -ADR [28]	64.40	20.20	49.20	15.00	61.60	53.30	43.95
PTE [62]	62.40	21.00	50.50	13.80	64.60	52.40	44.12
SAGM [63]	<u>64.90</u>	21.10	<u>51.50</u>	14.80	<u>64.10</u>	<u>53.60</u>	<u>45.00</u>
NormAUG (ours)	63.10	27.30	54.30	17.30	62.00	54.80	46.47

TABLE VI

ABLATION STUDIES ON PACS AND MINI-DOMAINNET. IN THIS TABLE, “DA” DENOTES THE “DEEPALL” MODEL.

Methods	ON	AUG	EP	PACS				
				A	C	P	S	Avg.
DA [34]	-	-	-	77.63	76.77	95.85	69.50	79.94
Model-1	✓	-	-	80.64	81.68	95.30	79.34	84.24
Model-2	✓	✓	-	83.85	82.55	94.95	83.85	86.30
Ours	✓	✓	✓	85.60	81.85	95.70	85.00	87.04
Method	ON	AUG	EP	mini-DomainNet				
DA [34]	-	-	-	65.30	58.40	64.70	59.00	61.86
Model-1	✓	-	-	67.95	63.10	70.00	57.10	64.54
Model-2	✓	✓	-	69.80	65.80	69.90	61.70	66.80
Ours	✓	✓	✓	70.20	66.90	71.20	63.40	67.93

experimental setup and the fused scheme used in the test stage are consistent across all variant methods. The experimental results are summarized in Fig. 7. As observed in the figure, using independent classifiers for each normalization combination yields the best performance among the variant methods. This can be attributed to the increased diversity introduced by the multiple different classifiers. Therefore, our method adopts the use of independent classifiers for all paths. It is important to emphasize that all methods in Fig. 7 are evaluated using the same fused scheme in the test stage. Furthermore, we investigate the impact of the number of classifiers for the main path, as reported in Tab. VII. On the PACS dataset, using one classifier (“1CLS”) and two classifiers (“2CLS”) for the main path results in a slight performance decrease of -0.22% and -0.26% , respectively, compared to our method. Similarly, on the mini-DomainNet dataset, “1CLS” and “2CLS” lead to a performance decrease of -0.41% and -0.00% compared to our method in the main path. These results indicate that employing different classifiers for each BN in the BN bank can also slightly improve the performance of the main path.

Evaluation on different fusion schemes on prediction. In the test phase, we combine the results from the main path and the independent domain paths. In our method, recognizing the importance of the main path, we adopt a two-step averaging method. First, we average the results from the independent domain paths, and then we average this averaged result with the result from the main path. Additionally, we also evaluate a direct averaging method where we average all results

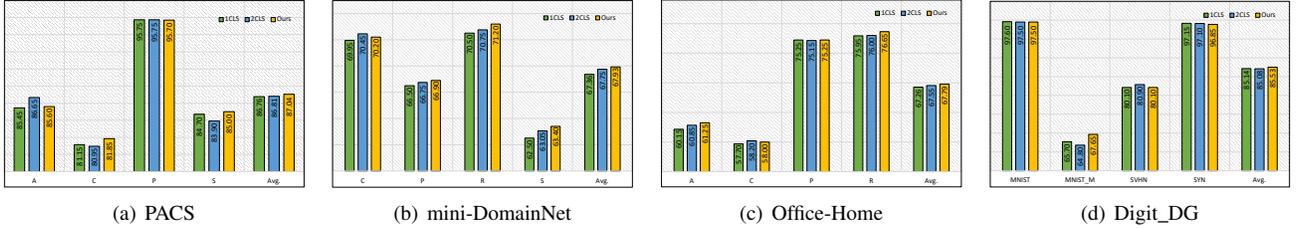


Fig. 7. Experimental results of different numbers of classifiers in our method on four datasets. In this figure, “1CLS” denotes that the main path and the auxiliary path share the same classifier. “2CLS” indicates that the main path and the auxiliary path use different classifiers, and these sub-paths in the auxiliary path share a classifiers. It is worth noting that all method use the same fused scheme in the test stage.

TABLE VII

EXPERIMENTAL RESULTS OF THE MAIN PATH. THE MODEL USES DIFFERENT NUMBERS OF CLASSIFIERS, WHICH IS CORRESPONDING WITH FIG 7 (A) AND (B). “DROP” IN THIS TABLE DENOTES THE VALUE OF “1CLS OR 2CLS-OURS”.

PACS						
Methods	A	C	P	S	Avg.	Drop
Ours	83.85	82.55	94.95	83.85	86.30	-
1CLS	83.25	82.15	95.05	83.85	86.08	-0.22
2CLS	84.20	81.85	95.00	83.10	86.04	-0.26
mini-DomainNet						
Methods	C	P	R	S	Avg.	Drop
Ours	69.80	65.80	69.90	61.70	66.80	-
1CLS	69.50	65.45	69.35	61.25	66.39	-0.41
2CLS	70.05	65.70	69.55	61.90	66.80	-0.00

from the main path and the independent domain paths. The experimental results are presented in Tab. VIII. As observed in the table, using the two-step averaging scheme leads to a decrease in performance compared to the direct averaging method. This suggests that the two-step averaging scheme reduces the influence of the main path’s result in the final ensemble result.

TABLE VIII

EXPERIMENTAL RESULTS OF DIFFERENT FUSION SCHEMES ON PREDICTION.

Methods	PACS	Office-Home	mini-DomainNet	Digit_DG
MeanAll	86.21	67.70	67.88	85.33
Ours	87.04	67.79	67.93	85.53

Analysis of the divergence. In this section, we evaluate the effectiveness of our NormAUG method in reducing the domain gap between source and source (D_{s2s}) and between source and target (D_{s2t}). To compute these divergences, we calculate the averaged feature of the d -th domain, denoted as f_d^s , the averaged feature of the source domains, denoted as f^s , and the averaged feature of the target domain, denoted as f^t . By utilizing these features, we can compute the divergence between source and source as $D_{s2s} = \frac{1}{D} \sum_{d=1}^D \|f^s - f_d^s\|$ and the divergence between source and target as $D_{s2t} = \|f^s - f^t\|$, where D is the number of source domains. The experimental results are presented in Tab. IX. In this table, “Model-1” is the same as “Model-1” in Tab. VI. As observed, our NormAUG method effectively reduces the domain gap between source and source as well as between source and target. This demonstrates that our method is capable of learning domain-invariant features, thus enhancing the generalization capability across different domains.

Evaluation on “AUG-1” in the BN bank. As discussed

TABLE IX

DATA-DISTRIBUTION DISTANCE BETWEEN SOURCE AND SOURCE (S2S) AND BETWEEN SOURCE AND TARGET (S2T) ON PACS. HERE, A SMALLER DOMAIN GAP IS PREFERABLE.

Tasks	S2T		S2S	
	Ours	Model-1	Ours	Model-1
CPS→A	1.82	2.49	1.83	1.91
APS→C	1.88	2.40	1.79	1.88
ACS→P	3.49	3.76	1.59	1.84
ACP→S	4.62	4.79	1.65	1.67

in Section III, our NormAUG method incorporates a random combination scheme for conducting the normalization operation in the auxiliary at each iteration. To evaluate the effectiveness of this scheme, we always perform experiments using a single domain to independently conduct the normalization, which is the same as the test scheme in Fig. 4 or “AUG-1” in Fig. 2. The experimental results are presented in Tab. X. From the results, we can observe that the random combination scheme significantly improves the model’s generalization compared to using a single domain for normalization. This improvement can be attributed to the increased diversity introduced by the random combination scheme. The random combination allows the model to explore different combinations of multiple domains during the normalization process, leading to a more comprehensive representation of the data and enhancing the model’s ability to generalize across domains.

TABLE X

COMPARISON BETWEEN THE SINGLE-DOMAIN AUGMENTATION AND RANDOM COMBINATION AUGMENTATION.

Methods	PACS				
	A	C	P	S	Avg.
Single	85.28	79.54	96.66	83.10	86.15
Ours	85.60	81.85	95.70	85.00	87.04
Methods	Office-Home				
	A	C	P	R	Avg.
Single	60.60	56.40	73.60	76.20	66.70
Ours	61.25	58.00	75.25	76.65	67.79
Methods	mini-DomainNet				
	C	P	R	S	Avg.
Single	68.40	65.90	70.90	61.90	66.78
Ours	70.20	66.90	71.20	63.40	67.93
Methods	Digit_DG				
	MN	MN-M	SV	SY	Avg.
Single	97.45	67.10	79.68	96.40	85.16
Ours	97.50	67.65	80.10	96.85	85.53

Result of using all sub-paths in the test stage. We also conduct experiments using all sub-paths of the auxiliary path in the test stage. The experimental results are presented in

Tab. XI. From this table, we observe that using all sub-paths of the auxiliary path (*i.e.*, “All” in Tab. XI) can slightly improve the performance compared to using the sub-paths for the independent domains alone, especially on the PACS and Office-Home datasets. However, it is important to note that using all sub-paths requires additional computation cost in the test stage. In our NormAUG method, we primarily focus on utilizing the sub-paths for the independent domain, which has shown significant improvements in performance, as demonstrated in Tab. VI. This method strikes a balance between performance improvement and computational efficiency. By utilizing the sub-paths for the independent domain, we can effectively enhance the model’s generalization while maintaining a reasonable computational cost in the test stage.

TABLE XI

EXPERIMENTAL RESULTS WITH DIFFERENT SUB-PATHS OF THE AUXILIARY PATH IN THE TEST STAGE. “ALL” IS USING ALL SUB-PATHS, AND “INDEPENDENT” IS ONLY USING THESE SUB-PATHS FOR THE INDEPENDENT DOMAIN IN THE TEST STAGE.

PACS					
Methods	A	C	P	S	Avg.
All	86.05	82.25	95.80	84.80	87.23
Independent (ours)	85.60	81.85	95.70	85.00	87.04
Office-Home					
Methods	A	C	P	R	Avg.
All	61.20	58.15	75.30	76.75	67.85
Independent (ours)	61.25	58.00	75.25	76.65	67.79
mini-DomainNet					
Methods	C	P	R	S	Avg.
All	70.40	66.80	71.20	63.20	67.90
Independent (ours)	70.20	66.90	71.20	63.40	67.93

Stability of the training process. We demonstrate the stability of our method during the training process, as illustrated in Fig. 8. From these figures, it is evident that our method exhibits a consistent improvement in performance over time, specifically in the unseen target domain. This stability further emphasizes the effectiveness and robustness of our method in enhancing the model’s generalization capabilities. Different from multiple existing methods that select the best model for evaluation [28], [9], we leverage the model from the last one epoch for evaluation in all experiments.

Evaluation of the test time. Our method introduces a Batch Normalization (BN) bank and multiple classifiers to enhance the generalization performance of features. To ensure a balance between improved performance and computational efficiency, we retain only a subset of paths. The time consumption during testing in the Art Painting domain of the PACS on a single RTX 3090 is presented in Tab. XII. It is evident that our method incurs only a slight increase in time compared to DeepAll [34], while achieving state-of-the-art performance. Besides, as the classifier is solely a fully-connected layer, reflected in the results presented in Tab. XII, our method experiences a slight increment in parameter count.

Further evaluation of the different fusion scheme. We conduct experiments to investigate the influence of various prediction’s strategies on generalization performance. We employed average and maximum fusion strategies for the independent domain paths and the main path, respectively. As shown in Tab. XIII, it is evident that our strategy achieves

TABLE XII
PARAMETER AND TIME RESULTS FOR VARIOUS METHODS TESTED ON PACS.

Methods	Parameters	Time per Image (ms/pic)	Accuracy
DeepAll [34]	11.18M	1.46	79.94
Mixstyle [12]	11.18M	1.49	83.70
I ² -ADR [28]	12.23M	1.44	85.55
Ours	11.21M	1.63	87.04

the highest average performance on four benchmark datasets. It shows that the averaging fusion strategy can improve generalization performance.

TABLE XIII

EXPERIMENTAL RESULTS OF DIFFERENT FUSION SCHEMES ON PREDICTION ON PACS, OFFICE-HOME (OH), MINI-DOMAINNET (MD), AND DIGIT_DG (DD). “M” MEANS THE MAIN PATH AND “I” MEANS THE INDEPENDENT DOMAIN PATHS.

Fusion Strategy	PACS	OH	mD	DD	Avg.
M	86.30	65.45	66.80	84.48	75.76
Max(I)	84.40	66.68	65.78	84.45	75.33
Mean(I)	84.02	66.21	67.40	84.38	75.50
Mean(I, M)	86.21	67.70	67.88	85.33	76.78
Max(I, M)	86.45	66.95	66.63	85.31	76.33
Max(Mean(I), M)	86.58	66.45	67.28	85.34	76.41
Mean(Max(I), M)	87.18	67.79	67.25	85.48	76.93
Ours Mean(Mean(I), M)	87.04	67.79	67.93	85.53	77.07

Visualization of the activation map. In this section, we present the activation maps of our method, as depicted in Fig. 9. As observed, when compared to the baseline, our method exhibits a more focused activation on the key regions of each object. This highlights the effectiveness of our NormAUG method in improving the model’s generalization in the target domain by capturing and highlighting the relevant features and regions.

V. CONCLUSION

In this paper, we introduce a novel data augmentation method based on the normalization perspective for domain generalization. Our method consists of two paths during the training stage: the main path, which serves as the baseline in the domain generalization task, and the auxiliary path, which incorporates our proposed augmentation method. Unlike traditional augmentation schemes, we adopt a random combination of domains for the normalization operation, thus enriching the diversity of the training data. Additionally, our method allows for the fusion of auxiliary results to enhance the model’s performance in the test stage. Experimental results on various benchmark datasets demonstrate the effectiveness of our normalization-guided augmentation.

In our method, the normalization-based augmentation scheme does not introduce additional information, indicating that indirectly generating diverse style information in our method does not create a large domain gap between training samples, as shown in the theoretical analysis. However, the generated diverse data from our NormAUG is not the richest training set for each specific task. Therefore, in future work, we will explore introducing slight extra information to further enrich data’s diversity while keeping the semantic information.

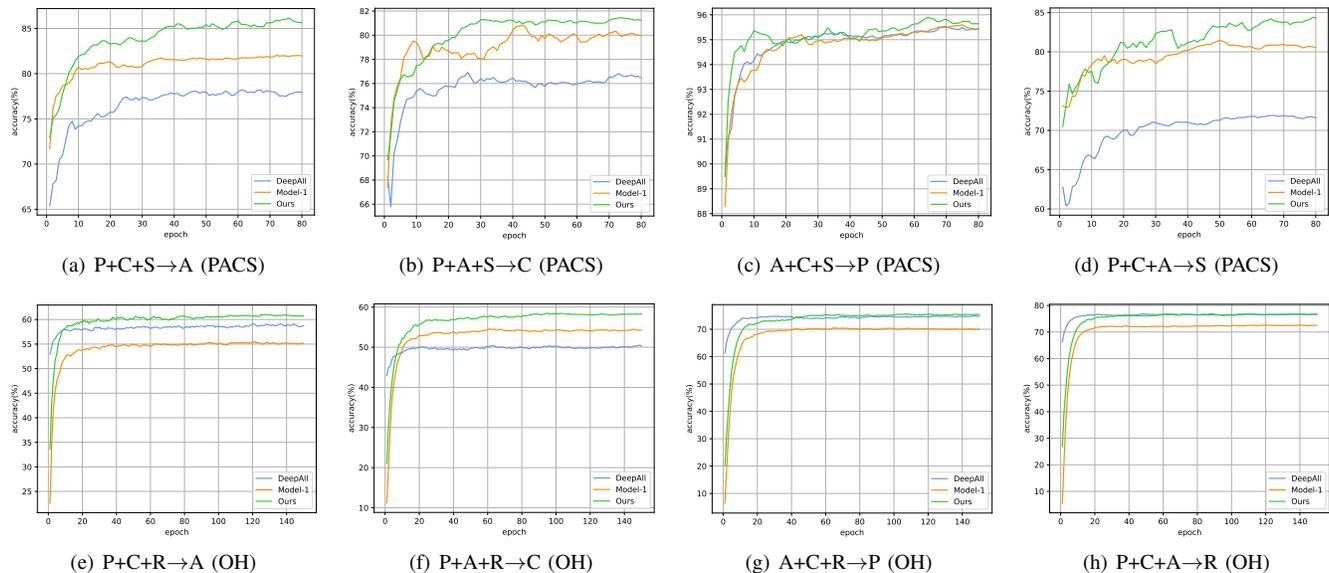


Fig. 8. The stability of the training process on PACS and Office-Home (OH).

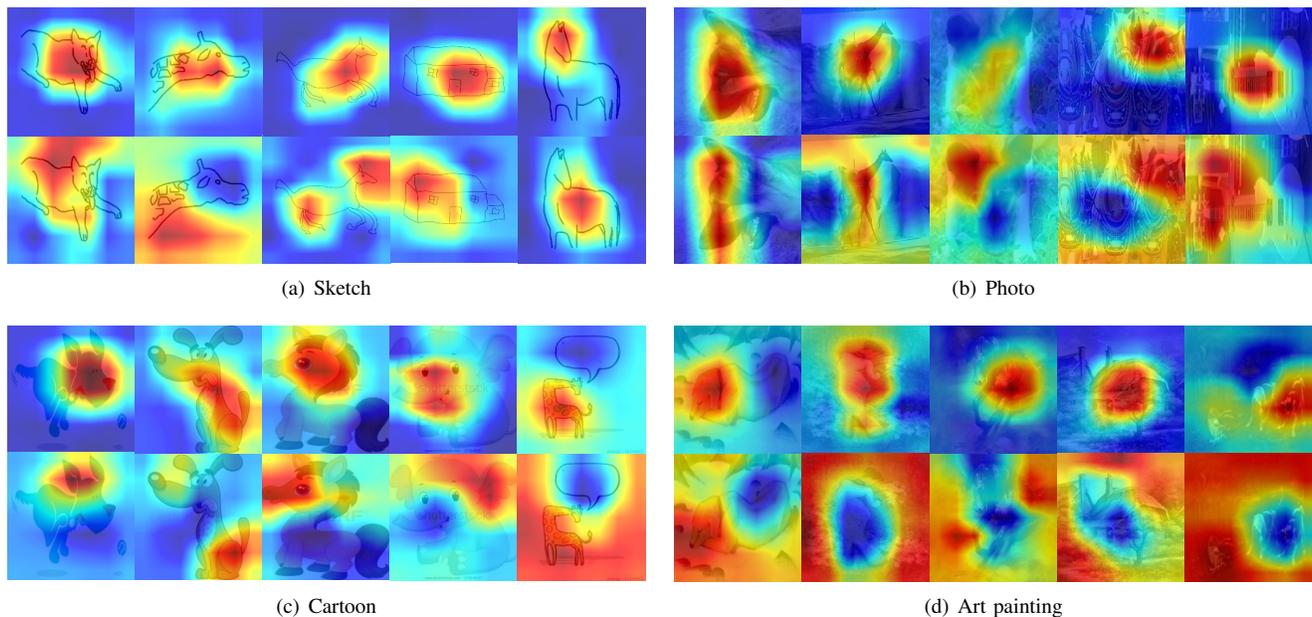


Fig. 9. The activation maps of our method (top) and the baseline (bottom) on PACS. In this figure, the redder area indicates the more attention.

REFERENCES

- [1] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [2] R. Aversa, P. Coronica, C. De Nobili, and S. Cozzini, “Deep learning, feature learning, and clustering analysis for sem image classification,” *Data Intelligence (DI)*, vol. 2, no. 4, pp. 513–528, 2020.
- [3] X. Wang, T. E. Huang, B. Liu, F. Yu, X. Wang, J. E. Gonzalez, and T. Darrell, “Robust object detection via instance-level temporal cycle confusion,” in *IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 9123–9132.
- [4] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015, pp. 234–241.
- [5] S. Lingwal, K. K. Bhatia, and M. Singh, “Semantic segmentation of landcover for cropland mapping and area estimation using machine learning techniques,” *Data Intelligence (DI)*, vol. 5, no. 2, pp. 370–387, 2023.
- [6] P. Li, D. Li, W. Li, S. Gong, Y. Fu, and T. M. Hospedales, “A simple feature augmentation for domain generalization,” in *International Conference on Computer Vision (ICCV)*, 2021, pp. 8886–8895.
- [7] C. Li, D. Zhang, W. Huang, and J. Zhang, “Cross contrasting feature perturbation for domain generalization,” in *International Conference on Computer Vision (ICCV)*, 2023, pp. 1327–1337.
- [8] Y. Wang, L. Qi, Y. Shi, and Y. Gao, “Feature-based style randomization for domain generalization,” *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)*, vol. 32, no. 8, pp. 5495–5509, 2022.
- [9] J. Zhang, L. Qi, Y. Shi, and Y. Gao, “MVDG: A unified multi-view framework for domain generalization,” in *European Conference on Computer Vision (ECCV)*, 2022, pp. 161–177.
- [10] K. Zhou, Y. Yang, Y. Qiao, and T. Xiang, “Domain adaptive ensemble learning,” *IEEE Transactions on Image Processing (TIP)*, vol. 30, pp. 8008–8018, 2021.
- [11] Y. Ding, L. Wang, B. Liang, S. Liang, Y. Wang, and F. Chen, “Domain generalization by learning and removing domain-specific features,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [12] K. Zhou, Y. Yang, Y. Qiao, and T. Xiang, “Domain generalization

- with mixstyle,” in *International Conference on Learning Representations (ICLR)*, 2021.
- [13] Z. Ding and Y. Fu, “Deep domain generalization with structured low-rank constraint,” *IEEE Transactions on Image Processing (TIP)*, vol. 27, no. 1, pp. 304–313, 2018.
- [14] Q. Xu, R. Zhang, Y. Zhang, Y. Wang, and Q. Tian, “A fourier-based framework for domain generalization,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 14 383–14 392.
- [15] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of Machine Learning Research (JMLR)*, vol. 9, no. 11, 2008.
- [16] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 248–255.
- [17] Y. Zhang, M. Li, R. Li, K. Jia, and L. Zhang, “Exact feature distribution matching for arbitrary style transfer and domain generalization,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 8025–8035.
- [18] Z. Huang, H. Wang, E. P. Xing, and D. Huang, “Self-challenging improves cross-domain generalization,” in *European Conference on Computer Vision (ECCV)*, 2020, pp. 124–140.
- [19] Z. Wang, Y. Luo, R. Qiu, Z. Huang, and M. Baktashmotlagh, “Learning to diversify for single domain generalization,” in *IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 814–823.
- [20] X. Huang and S. J. Belongie, “Arbitrary style transfer in real-time with adaptive instance normalization,” in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 1510–1519.
- [21] J. Kang, S. Lee, N. Kim, and S. Kwak, “Style neophile: Constantly seeking novel styles for domain generalization,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 7120–7130.
- [22] L. Niu, W. Li, and D. Xu, “Visual recognition by learning from web data: A weakly supervised domain generalization approach,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 2774–2783.
- [23] S. Seo, Y. Suh, D. Kim, G. Kim, J. Han, and B. Han, “Learning to optimize domain specific normalization for domain generalization,” in *European Conference on Computer Vision (ECCV)*, 2020, pp. 68–83.
- [24] M. Segù, A. Tonioni, and F. Tombari, “Batch normalization embeddings for deep domain generalization,” *Pattern Recognition (PR)*, vol. 135, p. 109115, 2023.
- [25] Y. Li, M. Gong, X. Tian, T. Liu, and D. Tao, “Domain generalization via conditional invariant representations,” in *AAAI Conference on Artificial Intelligence (AAAI)*, 2018, pp. 3579–3587.
- [26] T. Matsuura and T. Harada, “Domain generalization using a mixture of multiple latent domains,” in *AAAI Conference on Artificial Intelligence (AAAI)*, 2020, pp. 11 749–11 756.
- [27] P. Chattopadhyay, Y. Balaji, and J. Hoffman, “Learning to balance specificity and invariance for in and out of domain generalization,” in *European Conference on Computer Vision (ECCV)*, 2020, pp. 301–318.
- [28] R. Meng, X. Li, W. Chen, S. Yang, J. Song, X. Wang, L. Zhang, M. Song, D. Xie, and S. Pu, “Attention diversification for domain generalization,” in *European Conference on Computer Vision (ECCV)*, 2022, pp. 322–340.
- [29] K. Lee, S. Kim, and S. Kwak, “Cross-domain ensemble distillation for domain generalization,” in *European Conference on Computer Vision (ECCV)*, 2022, pp. 1–20.
- [30] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *International Conference on Machine Learning (ICML)*, pp. 1126–1135.
- [31] Y. Balaji, S. Sankaranarayanan, and R. Chellappa, “Metareg: Towards domain generalization using meta-regularization,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp. 1006–1016.
- [32] D. Li, Y. Yang, Y. Song, and T. M. Hospedales, “Learning to generalize: Meta-learning for domain generalization,” in *AAAI Conference on Artificial Intelligence (AAAI)*, 2018, pp. 3490–3497.
- [33] Q. Dou, D. C. de Castro, K. Kamnitsas, and B. Glocker, “Domain generalization via model-agnostic learning of semantic features,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 6447–6458.
- [34] F. Lv, J. Liang, S. Li, J. Zhang, and D. Liu, “Improving generalization with domain convex game,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 24 315–24 324.
- [35] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *International Conference on Machine Learning (ICML)*, 2015, pp. 448–456.
- [36] I. Albuquerque, J. Monteiro, M. Darvishi, T. H. Falk, and I. Mitliagkas, “Generalizing to unseen domains via distribution matching,” *arXiv preprint arXiv:1911.00804*, 2019.
- [37] J. Wang, C. Lan, C. Liu, Y. Ouyang, and T. Qin, “Generalizing to unseen domains: A survey on domain generalization,” in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2021, pp. 4627–4635.
- [38] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales, “Deeper, broader and artier domain generalization,” in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 5543–5551.
- [39] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan, “Deep hashing network for unsupervised domain adaptation,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 5018–5027.
- [40] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang, “Moment matching for multi-source domain adaptation,” in *IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 1406–1415.
- [41] K. Zhou, Y. Yang, T. Hospedales, and T. Xiang, “Deep domain-adversarial image generation for domain generalisation,” in *AAAI Conference on Artificial Intelligence (AAAI)*, vol. 34, no. 07, 2020, pp. 13 025–13 032.
- [42] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [43] Y. Ganin and V. Lempitsky, “Unsupervised domain adaptation by back-propagation,” in *International conference on machine learning (ICML)*, 2015, pp. 1180–1189.
- [44] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, “Reading digits in natural images with unsupervised feature learning,” 2011.
- [45] F. M. Carlucci, A. D’Innocente, S. Bucci, B. Caputo, and T. Tommasi, “Domain generalization by solving jigsaw puzzles,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 2229–2238.
- [46] C. Chen, J. Li, X. Han, X. Liu, and Y. Yu, “Compound domain generalization via meta-knowledge encoding,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 7109–7119.
- [47] F. Lv, J. Liang, S. Li, B. Zang, C. H. Liu, Z. Wang, and D. Liu, “Causality inspired representation learning for domain generalization,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 8036–8046.
- [48] A. Ramé, C. Dancette, and M. Cord, “Fishr: Invariant gradient variances for out-of-distribution generalization,” in *International Conference on Machine Learning (ICML)*, 2022, pp. 18 347–18 377.
- [49] M. Bui, T. Tran, A. Tran, and D. Q. Phung, “Exploiting domain-specific features to enhance domain generalization,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021, pp. 21 189–21 201.
- [50] H. Xia, T. Jing, and Z. Ding, “Generative inference network for imbalanced domain generalization,” *IEEE Transactions on Image Processing (TIP)*, vol. 32, pp. 1694–1704, 2023.
- [51] Y. Wang, F. Liu, Z. Chen, Y. Wu, J. Hao, G. Chen, and P. Heng, “Contrastive-acc: Domain generalization through alignment of causal mechanisms,” *IEEE Transactions on Image Processing (TIP)*, vol. 32, pp. 235–250, 2023.
- [52] J. Cha, S. Chun, K. Lee, H. Cho, S. Park, Y. Lee, and S. Park, “SWAD: domain generalization by seeking flat minima,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021, pp. 22 405–22 418.
- [53] K. Zhou, Y. Yang, T. M. Hospedales, and T. Xiang, “Learning to generate novel domains for domain generalization,” in *European Conference on Computer Vision (ECCV)*, 2020, pp. 561–578.
- [54] X. Li, Y. Dai, Y. Ge, J. Liu, Y. Shan, and L. Duan, “Uncertainty modeling for out-of-distribution generalization,” in *International Conference on Learning Representations (ICLR)*, 2022.
- [55] C. Lee, T. Batra, M. H. Baig, and D. Ulbricht, “Sliced wasserstein discrepancy for unsupervised domain adaptation,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 10 285–10 295.
- [56] H. Zhang, M. Cissé, Y. N. Dauphin, and D. Lopez-Paz, “mixup: Beyond empirical risk minimization,” in *International Conference on Learning Representations (ICLR)*, 2018.
- [57] H. Nam, H. Lee, J. Park, W. Yoon, and D. Yoo, “Reducing domain gap by reducing style bias,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 8690–8699.
- [58] G. Blanchard, A. A. Deshmukh, Ü. Dogan, G. Lee, and C. Scott, “Domain generalization by marginal transfer learning,” *Journal of Machine Learning Research (JMLR)*, vol. 22, pp. 2:1–2:55, 2021.

- [59] Y. Chen, Y. Wang, Y. Pan, T. Yao, X. Tian, and T. Mei, "A style and semantic memory mechanism for domain generalization," in *IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 9144–9153.
- [60] L. Chen, Y. Zhang, Y. Song, Y. Shan, and L. Liu, "Improved test-time adaptation for domain generalization," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 24 172–24 182.
- [61] D. Kim, Y. Yoo, S. Park, J. Kim, and J. Lee, "Selfreg: Self-supervised contrastive regularization for domain generalization," in *IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 9599–9608.
- [62] S. Min, N. Park, S. Kim, S. Park, and J. Kim, "Grounding visual representations with texts for domain generalization," in *European Conference on Computer Vision (ECCV)*, 2022, pp. 37–53.
- [63] P. Wang, Z. Zhang, Z. Lei, and L. Zhang, "Sharpness-aware gradient matching for domain generalization," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 3769–3778.