

PLACE dropout: A Progressive Layer-wise and Channel-wise Dropout for Domain Generalization

JINTAO GUO, The State Key Laboratory for Novel Software Technology, Nanjing University, China
 LEI QI*, 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

YINGHUAN SHI*, The State Key Laboratory for Novel Software Technology, Nanjing University, China
 YANG GAO, The State Key Laboratory for Novel Software Technology, Nanjing University, China

Domain generalization (DG) aims to learn a generic model from multiple observed source domains that generalizes well to arbitrary unseen target domains without further training. The major challenge in DG is that the model inevitably faces a severe overfitting issue due to the domain gap between source and target domains. To mitigate this problem, some dropout-based methods have been proposed to resist overfitting by discarding part of the representation of the intermediate layers. However, we observe that most of these methods only conduct the dropout operation in some specific layers, leading to an insufficient regularization effect on the model. We argue that applying dropout at multiple layers can produce stronger regularization effects, which could alleviate the overfitting problem on source domains more adequately than previous layer-specific dropout methods. In this paper, we develop a novel layer-wise and channel-wise dropout for DG, which randomly selects one layer and then randomly selects its channels to conduct dropout. Particularly, the proposed method can generate a variety of data variants to better deal with the overfitting issue. We also provide theoretical analysis for our dropout method and prove that it can effectively reduce the generalization error bound. Besides, we leverage the progressive scheme to increase the dropout ratio with the training progress, which can gradually boost the difficulty of training the model to enhance its robustness. Extensive experiments on three standard benchmark datasets have demonstrated that our method outperforms several state-of-the-art DG methods. Our code is available at <https://github.com/lingeringlight/PLACEdropout>.

CCS Concepts: • **Computing methodologies** → **Object recognition**.

Additional Key Words and Phrases: Domain generalization, dropout regularization, overfitting problem, distribution shift

1 INTRODUCTION

Deep learning has achieved tremendous progress in various tasks over the last few years. Under the assumption that training and test data come from similar data distributions, deep convolutional neural networks have shown remarkable ability in a wide range of visual applications [21, 38, 56].

*Corresponding authors: Yinghuan Shi and Lei Qi.

Authors' addresses: Jintao Guo, The State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China, guojintao@smail.nju.edu.cn; Lei Qi, 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, Nanjing, China, qilei.cs@gmail.com; Yinghuan Shi, The State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China, syh@nju.edu.cn; Yang Gao, The State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China, gaoy@nju.edu.cn.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

1551-6857/2023/9-ART1 \$15.00
<https://doi.org/10.1145/3624015>

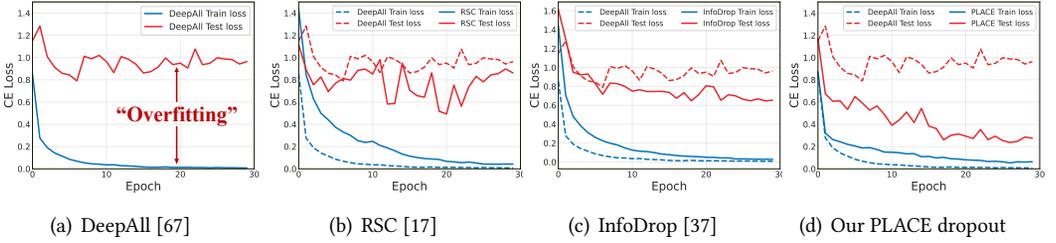


Fig. 1. The difference between the training loss and the test loss indicates the overfitting degree. The larger the difference, the more severe the overfitting issue of the model on source domains. We experiment on the PACS dataset with ResNet-18 as the backbone architecture. We select the conventional DG method DeepAll [67] and two representative layer-specific dropout methods, *i.e.*, RSC [17] and InfoDrop [37]. Our PLACE dropout can produce regularization effects on multiple layers of the network, which alleviates overfitting effectively with a decreased gap between source train loss and target test loss as shown in column 4.

However, their trained models often overfit the training data, leading to the inferior performance on out-of-distribution data [19]. Such catastrophic performance degradation caused by distribution discrepancy (*i.e.*, domain shift [30]) hinders the applications of deep neural networks, as in reality training and test data are often from different distributions [23, 54]. To address this issue, domain adaptation (DA) has been proposed to narrow the potential distribution discrepancy by leveraging labeled source domains and unlabeled target domain to jointly train the model [53, 55].

Unfortunately, despite their success, DA methods could not guarantee the model performance on unknown target domains that have not been seen during training [39, 51]. Since target data cannot always be available in real-world scenarios, collecting data from all potential target domains is expensive and impractical. Moreover, even if data from the target domain can be obtained, DA methods need to re-train the model on the new target domain, which is hard to achieve in reality. Aware of this fact, domain generalization (DG) has been proposed as a more challenging yet practical problem setting, which aims to train a model with multiple different but related source domains that can generalize well to arbitrary unseen target domains without re-training [46, 66]. The DG problem has aroused wide attention from researchers and many existing works have shown promising results by utilizing domain-invariant learning [6, 10, 24], meta-learning [1, 8, 48, 61], data augmentation [44, 50, 67], regularization strategies [17, 33, 37], *etc.*

One of the main challenges in DG is that the model is prone to overfitting the source domains due to the domain gap between source domains and unknown target domain, which greatly impairs the ability of the model to generalize to the target domain. One way to identify whether overfitting occurs is to examine if the difference between the training loss and the test loss is increasing [15]. As shown in Fig. 1(a), the conventional DG method, *i.e.*, Deepall [67], overfits the observed source domains soon after training begins due to the large distribution discrepancy between source and target domains. To alleviate the overfitting, some dropout-based regularization methods [17, 37] have shown their promising performance in DG task. Without introducing extra parameters, these methods can still achieve remarkable improvements over other DG methods by discarding potentially overfitting-related features during training. However, exiting dropout-based DG methods are mainly designed for specific layers, *e.g.*, InfoDrop [37] only works in the shallowest layer to reduce the texture-bias of the model, and RSC [17] only fits in the fully connected layer to reduce the model dependence on the over-dominate features. As a result, these methods could not produce a sufficiently regularizing effect on the model, leading to the relatively poor generalization

ability to unseen target domain if the model suffers a severe overfitting issue on source domains as shown in 1(b) and 1(c). Therefore, different from previous layer-specific dropout methods, we design a multiple-layers dropout method for DG to alleviate the overfitting issue adequately.

To be specific, we propose a novel Progressive LAYer-wise and ChannEl-wise (**PLACE** in short) dropout for domain generalization. The PLACE dropout consists of two components, *i.e.*, 1) the layer-wise and channel-wise dropout for perturbing the feature maps in multiple layers, and 2) the progressive scheme to gradually increase the difficulty of training the model. 1) *Layer-wise and channel-wise dropout*. Supported by previous dropout-related theoretical works [58, 60], we argue that conducting dropout in multiple layers could produce a stronger regularization effect than the single-layer dropout. However, conducting dropout in multiple layers simultaneously suffers a high risk of losing excessive information, which will most likely hinder model training and slow down the learning speed of the lower layers [31]. To overcome these issues, we design a simple *layer-wise and channel-wise dropout* which randomly selects a layer and then randomly drops its channels at each iteration. Our method can effectively resist the overfitting problem by generating diverse data variants in multiple layers during training. We also theoretically analyze the property of *layer-wise and channel-wise dropout* and prove that it can achieve a small generalization error on target domains. 2) *Progressive scheme*. Considering the risk of overfitting is small at the beginning of training but increases as training progresses [26], we introduce a training strategy that enhances the dropout ratio with a progressive scheme, which gradually raises the regularization effect to better tackle the overfitting issue. Extensive experiments validate the effectiveness of both the proposed dropout and the progressive scheme. We also conduct an in-depth analysis of the impact of our method on the cross-domain gap, which demonstrates that PLACE dropout can narrow the gap between source and target domains by reducing the impact of overfitting on source domains.

Our contributions can be summarized as follows:

- We propose a novel regularization method for DG, namely *LAYer-wise and ChannEl-wise dropout*, which can generate diverse variants of data in multiple layers to fight against the overfitting issue of the model on source domains. We also provide theoretical analysis that our method can generate a tight generalization error bound.
- We design a simple yet effective *progressive scheme* that gradually boosts the ratio of the layer-wise and channel-wise dropout to continually raise the difficulty of training the model, which can further improve the generalization ability of the model.
- Without introducing additional network parameters, our method achieves state-of-the-art accuracy on multiple benchmarks and outperforms all established dropout-based methods. We also build a strong baseline for DG with multiple augmentation methods, with which our PLACE dropout can achieve new SOTA performances, *e.g.*, 89.03% on PACS with ResNet-50.

2 RELATED WORK

2.1 Domain Generalization

Domain generalization (DG) aims to learn a robust model from multiple distinct but related domains that can generalize well to arbitrary unseen target domains. Existing DG methods can be roughly divided into four categories, including domain-invariant learning, meta-learning methods, data augmentation and regularization methods. Several works for domain-invariant learning have been explored, including domain-adversarial learning [10] and feature disentanglement [6]. Another popular way to address the DG problem is meta-learning, which simulates the domain shift by splitting the source domains into meta-train and meta-test domains [8, 52, 61, 62]. Data augmentation is also an important technology to empower the model generalization by enriching the diversity of existing training data from the image or feature level. The image-level augmentation methods

mainly generate virtual images via gradient-based adversarial attacks [44], domain-adversarial generation [39] or learnable augmentation networks [67]. And the feature-level augmentation methods can diversify image style by mixing or perturbing the feature statistics, thus generating data in different styles from source domains to boost model generalization [20, 50, 63]. Recently, some regularization methods have also been proposed to address the DG task via dropout strategies [17, 37], shape-biased learning [28] and self-supervise methods [4, 47].

Our work is most relevant to the dropout-based regularization methods [17, 37], which have shown promising performance for significantly enhancing the model generalization ability with introducing no extra parameters. Concretely, Shi *et al.* [37] find less informative regions contain texture-biased representation, thus proposing to mask them during training to reduce the model's texture bias. Huang *et al.* [17] assume that representations with the highest gradients are over-dominant and hinder the model from generalizing, thus designing a self-challenging algorithm to discard them during training. However, these assumption-based dropouts only apply to specific layers, which is insufficient to regularize the model if there exists a severe overfitting issue on source domains [31, 40]. In contrast to all the methods above, we investigate the role of dropout in model training and propose an assumption-free dropout method with layer-wise and channel-wise dropout and a progressive scheme to tackle the DG issue. Our method can be applied to multiple layers of the network to generate various data variants during training, which can adequately mitigate the overfitting problem of the model on source domains.

2.2 Dropout Regularization

Dropout [40] is one of the most widely employed regularization techniques to enhance the generalization capability of deep neural networks. Over the years, dropout has been extended in both channel-wise and space-wise manners. Tompson *et al.* [41] introduced Spatial Dropout, which drops entire channels in feature maps, while Ghiasi *et al.* [11] proposed DropBlock, a method that randomly masks contiguous feature regions. Subsequently, various novel dropout methods have been investigated, broadly categorized into task-auxiliary dropout and structure-information dropout. Task-auxiliary dropout methods aim to incorporate auxiliary tasks to help the model focus on relevant information for the primary task [27, 35], *e.g.*, Nagpal *et al.* [27] employ an auxiliary network to drop filters responsible for encoding given sensitive attributes. On the other hand, structure-information dropout methods leverage feature-level structure information to guide the dropout operations during training for diverse tasks [14, 22, 59].

However, these methods are primarily designed for supervised and semi-supervised learning settings, where training and test data typically share similar distributions. Consequently, they may not be well-suited for Domain Generalization (DG) tasks due to their insufficient regularization effect in combating the severe overfitting issues caused by the domain gap between source and target domains. Moreover, since the target domain is unavailable during training, unsuitable guidance for dropout might potentially damage the generalization ability of the model [59]. To address these challenges, we propose a novel dropout-based framework for DG, which can effectively mitigate the overfitting on source domains, reducing the cross-domain gap, and facilitating the model to generalize effectively to unseen target domains.

3 PROPOSED METHOD

In this section, we introduce the proposed method, Progressive LAYER-wise and CHANNEL-wise (PLACE in short) dropout. We first present the necessary preliminaries of our work. Then, we provide the details of the proposed method, including the layer-wise and channel-wise dropout and the progressive scheme. Finally, we outline the whole training algorithm of our PLACE dropout.

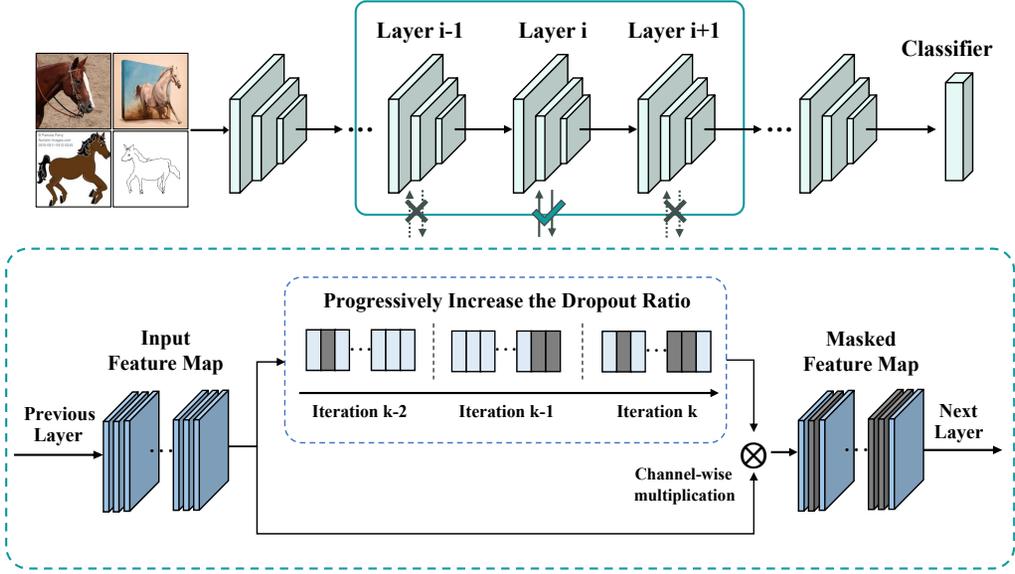


Fig. 2. Illustration of our PLACE dropout. Our method contains two primary components, including the layer-wise and channel-wise dropout and the progressive scheme. The figure shows the model is trained at the k -th iteration and conducted dropout in the i -th Layer. Detailed procedure is discussed in Section 3.3.

3.1 Preliminaries

Given a training set of multiple observed source domains $\mathcal{D}_S = \{D_1, D_2, \dots, D_K\}$ with N_k labelled samples $\{(x_i^k, y_i^k)\}_{i=1}^{N_k}$ in the k -th domain S_k , where K is the number of total source domains, x_i^k and y_i^k denote the samples and labels, respectively. For simplicity, we use (\mathbf{x}, \mathbf{y}) to replace the sample-label pair (x_i^k, y_i^k) in the following. The goal of domain generalization is to learn a domain-agnostic model $f(\cdot; \theta)$ on source domains \mathcal{D}_S that can generalize well to any unseen target domain \mathcal{D}_T without extra training. θ denotes the network parameters of $f(\cdot; \theta)$.

With the training data of source domains \mathcal{D}_S , our method minimizes the standard loss function:

$$\mathcal{L}_{ent} = \sum_{\langle \mathbf{x}, \mathbf{y} \rangle \sim \mathcal{D}_s} \ell(f(\mathbf{x}; \theta), \mathbf{y}). \quad (1)$$

where $\ell(\cdot, \cdot)$ is the cross entropy loss function. However, unlike previous DG methods, PLACE dropout focuses on generating noise by applying dropout to multiple layers for boosting the model robustness and minimizing the above loss function. Our method consists of two components, *i.e.*, *layer-wise and channel-wise dropout* to generate diverse data variants in multiple network layers, and *progressive scheme* for increasing dropout ratio to gradually raise the difficulty of training model. Concretely, at each iteration, PLACE dropout first randomly selects one middle layer, then randomly mutes its channels by a proportion that gradually increases as the training progresses, and finally updates the entire model. The overview of our method is illustrated in Fig. 2. We present the proposed method in detail in the following parts.

3.2 Progressive Layer- and Channel-wise Dropout

Channel-wise Dropout. We utilize channel-wise dropout [41], which randomly discards entire channels in feature maps, rather than the standard dropout [40] or the spatial-wise dropout [11]. The reason behind this choice lies in the high similarity between adjacent units within feature maps. Standard dropout, which drops individual neurons, does not exhibit a marked effect in convolution layers [31]. On the other hand, spatial-wise dropout, which masks contiguous regions, faces challenges in its regularization effect, as it is easily influenced by the location and size of the dropped regions, *i.e.*, over-dropping if the informative regions (*e.g.*, foreground regions) are masked or under-dropping if the irrelevant regions (*e.g.*, background regions) are discarded [11, 59]. In contrast, channel-wise dropout operates at the level of individual channels, which are basic units of feature maps corresponding to specific patterns in the input image. As a result, channel-wise dropout has been demonstrated to effectively ensure the dropout’s impact in convolution layers [14, 41]. Besides, dropping channels in feature maps reduces co-adaptations among different channels [40, 41], which can effectively mitigate the overfitting issue on source domains and implicitly encourage the model to learn comprehensive feature patterns.

Layer-wise and Channel-wise Dropout. Different from most previous dropout-based DG methods by conducting the dropout on specific layers [17, 37], we propose layer-wise dropout that randomly selects a middle layer of the network and then randomly discards its channels at each iteration. We first analyze the effect of random channel-wise dropout in different layers on the inter-domain discrepancy between source and target domains, which potentially reflects the overfitting degree of the model to the observed source domains. The discrepancy is calculated as [50]:

$$d^i = \frac{1}{K} \sum_{s=1}^K \|\text{GAP}(\bar{\mathbf{z}}_s^i) - \text{GAP}(\bar{\mathbf{z}}_t^i)\|_2, \quad (2)$$

where $\bar{\mathbf{z}}_s^i$ is the averaged feature maps of all samples from the i -th layer in the s -th source domain, and $\bar{\mathbf{z}}_t^i$ is the mean feature maps from the i -th layer in the target domain. K is the number of source domains. $\text{GAP}(\cdot)$ is the global average pooling operation. We calculate the domain discrepancy $d_{drop_j}^i$ according to features from the i -th layer during inserting dropout into the j -th layer, and d_{bal}^i for the baseline model. Then we calculate $d_{drop_j}^i - d_{bal}^i$ to investigate the effectiveness of dropout, where the “negative” value indicates that the dropout can reduce the domain difference.

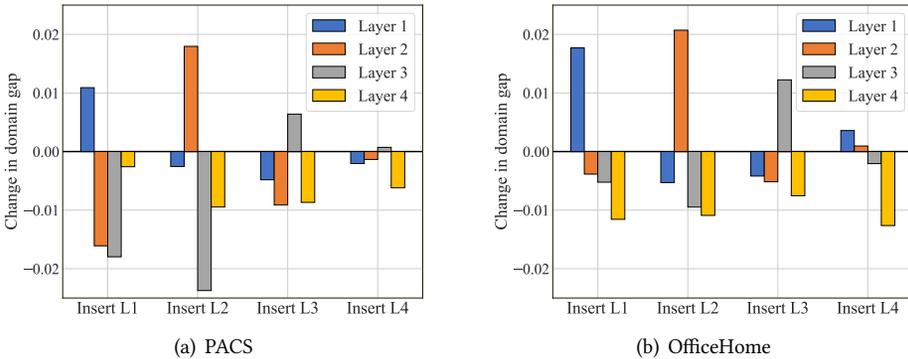


Fig. 3. The effectiveness of dropout in each layer of ResNet-18 when dropout is inserted in different layers. The horizontal axis denotes the position of inserting channel-wise dropout (*i.e.*, L_i is the i -th layer), and the vertical axis is the magnitude of change in domain gap from the i -th layer, which is calculated by $d_{drop}^i - d_{bal}^i$.

The results are presented in Fig. 3. Specifically, Fig. 3(a) presents the results on PACS when considering Art as the target domain and using the other three domains to train the model, and Fig. 3(b) presents the results on OfficeHome when taking RealWorld as the target domain. From Fig. 3, we observe that regardless of the layer where dropout is inserted, it consistently leads to a lower domain gap between the source and target domains at the 4-th layer, which is nearest to the prediction layer and its output can directly influence the prediction performance of the model. Furthermore, the results suggest that placing dropout in different layers has different effects on the domain gap between the source and target domains in each layer, *i.e.*, leading to a larger domain gap in the inserted layer but a smaller domain gap in the adjacent layers. This observation could be interpreted as that dropout introduces noise to the feature maps of the inserted layer, compelling the adjacent layers to learn general and robust features to resist the noise. In conclusion, *the feature representations in different layers exhibit significant diversity, leading to various effects on model training when dropout is placed in different layers.*

However, conducting dropout in multiple layers simultaneously could lead to conflicts, *e.g.*, if dropout is set in the 2-nd and 3-rd layers, dropout in the 2-nd layer will reduce the domain discrepancy on 3-rd layer, but setting dropout in 3-rd layer will increase the domain discrepancy itself. Moreover, the simultaneous usage of dropout in multiple layers could result in the absence of excessive information during training and hinder the model convergence [31]. To address these issues, we propose a novel approach called *layer-wise and channel-wise dropout*. Instead of applying dropout to all layers at once, our method randomly selects a middle layer of the network and performs dropout on its feature maps at each iteration. This approach introduces noise at multiple scales to the model, generating various data variants to enhance model robustness. Furthermore, our proposed PLACE dropout can effectively integrate the effects of dropout in different layers while avoiding conflicts in the inter-domain discrepancy across layers, thus producing a strong regularization effect to mitigate the overfitting issue on source domains.

Progressive Scheme. We design a time-adaptive regularization method, namely *progressive scheme*, to mitigate the overfitting issue adequately while ensuring the stability of model convergence. It is motivated by the intuition that the risk of overfitting on source domains is low at the beginning of training, and increases as the training proceeds [26]. To address this issue, we propose to increase the dropout ratio progressively, which can gradually enhance the regularization effect of dropout [45] and help the model better tackle the overfitting problem on source domains. Given the current iteration t , we compute the dropout ratio p and the number of dropped channels γ :

$$\gamma = C \times p = C \times P_{max} \times G(t), \quad (3)$$

$$G(t) = \frac{2}{\pi} \arctan\left(\frac{t}{V}\right), \quad (4)$$

where P_{max} denotes the maximum of channel-wise dropout ratio, C represents the number of channels, and V is a hyper-parameter to control the increasing speed of p , respectively. $G(\cdot)$ could be an appropriate monotonically increasing and strongly-convex function with t as input and a value from 0 and 1 as output, *i.e.*, the growth rate is relatively high at the start of training but gradually decreases as the training processes. Such function is also utilized in curriculum learning and related progressive settings [2, 26]. We empirically adopt the arctangent function (*i.e.*, $\arctan(\cdot)$) for $G(\cdot)$, which is formulated as Eq. (4), in the following experiments.

3.3 Training Algorithm

As introduced before, at each iteration, we first randomly select an index l of a middle layer from the network. Then we obtain the output of the l -th layer by forward propagation, which is denoted as $\mathbf{F}^l \in \mathbb{R}^{C \times H \times W}$, and generate an all-one mask matrix $\mathbf{M}^l \in \mathbb{R}^{C \times H \times W}$ with the same size as \mathbf{F}^l . C is

Algorithm 1: PLACE dropout Algorithm

Input: Batch size N , learning rate η , source data $\{(x_k, y_k)\}$, candidate layer set $\{l_1, \dots, l_n\}$, updated network f_{θ_u} .

Output: Trained network $f_{\theta_u^*}$.

for sampled mini-batch $\{(x_k, y_k)\}_{k=1}^N$ **do**

Sample l from $\{l_1, \dots, l_n\}$; // Select one middle layer from the network.

$\mathbf{F}^l \in \mathbb{R}^{C \times H \times W} \leftarrow$ Computed by forward propagation;

p and $\gamma \leftarrow$ Computed by Eq. (3); // Compute the proportion of dropout.

for each sample (x_i, y_i) **do**

Randomly sample γ channels $\{r_1, r_2, \dots, r_\gamma\}$ from F_i^l ;

$\mathbf{M}_i^l \leftarrow$ Computed by Eq. (5); // Compute the channel-wise mask matrix.

$\widehat{\mathbf{F}}^l \leftarrow$ Computed by Eq. (6); // Compute the masked feature map.

end

Update $\mathbf{F}^l \leftarrow \widehat{\mathbf{F}}^l$; // Update current feature map to the masked version.

Compute the output of the succeeding layers; // Forward propagation.

Compute $L \leftarrow -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log(p_{ik})$; // Compute the prediction loss.

Update $\theta_u \leftarrow \theta_u - \eta \nabla_{\theta_u} L$; // Update the parameters of the network.

end

the number of channels, H and W denote the dimension of height and width, respectively. Next, we calculate the dropout proportion p according to current iteration t by Eq. (3) and randomly sample γ distinct channel indexes $\{r_1, r_2, \dots, r_\gamma\}$ from the C channels of \mathbf{M}^l . The element in \mathbf{M}^l is set to 0 if its corresponding index is one of the selected channel indexes, and set to 1 otherwise:

$$\mathbf{M}_{i,j,k}^l = \begin{cases} 0, & i \in \{r_1, r_2, \dots, r_\gamma\} \\ 1, & \text{otherwise} \end{cases}. \quad (5)$$

We then obtain the masked output features map $\widehat{\mathbf{F}}^l$ as Eq. (6) by computing element-wise multiplication (denoted as \odot) of the output \mathbf{F}^l and the mask matrix \mathbf{M}^l :

$$\widehat{\mathbf{F}}^l = \mathbf{F}^l \odot \mathbf{M}^l. \quad (6)$$

Finally, $\widehat{\mathbf{F}}^l$ is forward propagated to the following part of the network to compute the prediction loss and update the parameters of the entire network.

The overall training procedure of PLACE dropout is summarized in Algorithm 1. In the training stage, our method can generate diverse data variants in multiple layers for tackling the overfitting issue on source domains. During the inference process, our PLACE dropout is closed as the conventional dropout [40]. It's worth noting that the PLACE dropout only comprises a few simple operations, such as random selection and channel-level product, *i.e.*, our method introduces no additional parameters and only incurs negligible computational cost during training. It also does not consume any extra computation time in the inference stage.

4 THEORETICAL ANALYSIS

In this section, we provide the theoretical analysis for layer-wise and channel-wise dropout and prove that our method can effectively reduce the generalization error bound of the model. We first analyze the domain generalization error bound for the dropout-based model, which can help us analyze the relationship between multi-layers dropout and model generalization. Then we

demonstrate that under mild assumptions, 1) *layer-wise and channel-wise dropout* can reduce the model sensitivity to the perturbation caused by dropout. 2) *layer-wise and channel-wise dropout* can generate more diverse augmented data than single-layer dropout; The results indicate that *layer-wise and channel-wise dropout* can lower the upper bound of generalization error and help the model generalize well to arbitrary unseen target domains.

Notations. Given the sample-label pair (\mathbf{x}, \mathbf{y}) , we use \mathbf{z} to denote the feature representation of (\mathbf{x}, \mathbf{y}) learned by the model. We define the task component of the model as $h: \mathcal{Z} \rightarrow \mathcal{Y}$ such that $h \in \mathcal{H}$, where \mathcal{H} is a set of candidate hypothesis. h takes \mathbf{z} as input and outputs the corresponding predict label. For simplicity, we take S to represent the set of feature-label pairs $\{(\mathbf{z}, \mathbf{y})\}$. $\tilde{\mathbf{z}}$ is the perturbed version of \mathbf{z} generated by dropout. \tilde{S} is the set containing all possible perturbed versions $\tilde{\mathbf{z}}$ of representations \mathbf{z} in S . Given a hypothesis h , the empirical risk $R[h]$ on domain \mathcal{D} is defined:

$$R_D[h] = \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim \mathcal{D}} \ell[h(\mathbf{z}), \mathbf{y}], \quad (7)$$

where the loss: $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow R_+$ quantifies the difference between $h(\mathbf{z})$ and \mathbf{y} for a data pair (\mathbf{z}, \mathbf{y}) .

We first introduce the generalization error bound of the dropout-based model under the DG set. Assuming that samples in the unknown target domain T could be regarded as perturbed versions of samples in source domains S , the empirical risk on the target domain is upper-bounded [17] by:

Lemma 1 [17]. *Let $L(\mathbf{z}_S, h)$ be the loss of a given hypothesis h on the source domains S , defined as:*

$$L(\mathbf{z}_S, h) = \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim S} \ell[h(\mathbf{z}), \mathbf{y}]. \quad (8)$$

$\xi(h)$ is a function of h that can reflect the sensitivity of the model to the perturbation of dropout:

$$\xi(h) = \sup_{\tilde{\mathbf{z}}_S \in \tilde{S}} |L(\mathbf{z}_S, h) - L(\tilde{\mathbf{z}}_S, h)|, \quad (9)$$

where $\tilde{\mathbf{z}}$ is the perturbed version of \mathbf{z} by dropout. Let N be the dataset size. Then for any $h \in \mathcal{H}$ and $\delta \in (0, 1)$, the following inequality holds with the probability at least $1 - \delta$:

$$R_T[h] \leq R_S[h] + \xi(h) \sqrt{\frac{\ln |\mathcal{H}| + \ln(2/\delta)}{2N}}. \quad (10)$$

The bound at the right side of Eq. 10 contains two terms: (1) the first is the empirical risk over all source domains; (2) the second indicates that the generalization bound of the model is positively related to the sensitivity $\xi(h)$ of the model to dropout, and negatively related to the dataset size N of source domains. We then analyze the effectiveness of layer-wise and channel-wise dropout by presenting that it reduces the sensitivity $\xi(h)$ while increasing the dataset size N .

Proposition 1. *Let $\hat{\xi}_t(h)$ be the estimated value of $\xi(h)$ at iteration t , which is defined as: $\hat{\xi}_t(h) = |L(\mathbf{z}_t, h_t) - L(\tilde{\mathbf{z}}_t, h_t)|$. Given a sufficiently small learning rate η , if discarding structural features will increase the empirical loss at the current iteration t , it holds that the layer-wise and channel-wise dropout can continually decrease $\hat{\xi}_t(h)$ and lead it to be a small number at the end of training.*

Proof. See in the Appendix.

Proposition 2. *Dropout in different layers perform as different data augmentations in the input space, thus the layer-wise dropout can generate more diverse augmented data than the single-layer dropout, i.e., increasing the dataset size N .*

Proof. See in the Appendix.

According to Propositions 1 and 2, we can conclude that *the layer-wise and channel-wise dropout can enrich the diversity of training data and reduce the sensitivity of model to the perturbation by dropout, thus generating a tight generalization error bound.* The theoretical results justify our

motivation and formally verify the efficacy of our method in terms of improving the generalization error bound. In the following, we will experimentally demonstrate the superiority of our method.

5 EXPERIMENTS AND RESULTS

5.1 Datasets and Settings

Datasets. We evaluate our method on three popularly-used DG benchmark datasets as follows:

- **PACS** [19] consists of images from 4 domains with a large discrepancy in image styles: Photo, Art Painting, Cartoon, and Sketch, including 7 object categories and 9,991 images total. We adopt the official split provided by [19] for training and test;
- **VLCS** [42] comprises of 5 categories, which are selected from 4 photo datasets, *i.e.*, VOC 2007 (Pascal), LabelMe, Caltech and Sun datasets. We utilize the same experimental setup as [4] and divide the dataset into the training and test sets based on 7 : 3;
- **Office-Home** [43] is an object recognition benchmark that contains around 15,500 images of 65 categories from 4 domains: Artistic, Clipart, Product and Real-World. Following [4], we randomly split each domain into 90% for training and 10% for test.

Experimental Settings. We apply the leave-one-domain-out protocol for all benchmarks, *i.e.*, we train the model on source domains and test the model on the remaining domain. Due to the progressive scheme, we select the model of the last epoch as the final model and report the top-1 classification accuracy. All the reported results are the averaged value over five runs.

5.2 Implementation Details

For the PACS dataset, we use the ImageNet pre-trained ResNet-18 and ResNet-50 as backbones following [12] and adopt the same hyper-parameters as [4]. For the VLCS dataset, We employ the experimental protocol as mentioned in [64] and use ResNet-18 as the backbone. For the OfficeHome dataset, we follow the same experimental setup as [9]. We train the network using SGD with a momentum of 0.9 and weight decay of 5×10^{-4} for a total of 30 epochs. The initial learning rate of the network is 4×10^{-3} and decayed by 0.1 at 80% of the total epochs.

Our method consists of two primary components, *i.e.*, the progressive scheme and the layer-wise and channel-wise dropout. For the progressive scheme, the dropout proportion is related to two hyperparameters, *i.e.*, the maximum dropout ratio P_{max} , and the progressive rate V . The maximum dropout ratio P_{max} is set as 0.33 for PACS and VLCS, and 0.25 for OfficeHome, respectively. The progressive rate V is set to 4 in all datasets. For the layer-wise and channel-wise dropout, we select the 1st, 2nd, and 3rd residual layers as the candidate set for both ResNet-18 and ResNet-50, which is determined experimentally in Tab. 8. At each iteration, we randomly select a residual layer from the candidate set and perform the progressive channel-wise dropout on its feature maps. Practically, we design a strong baseline model, denoted as DeepAll⁺⁺, for domain generalization with image-level and feature-level augmentation methods, *i.e.*, we perform random augmentation [7] on the images before training, and then randomly swap the style statistics of features in the 3-rd layer for any two samples during training [3], which is motivated by AdaIN [16]. The strength of random augmentation is controlled by both the number of transformations and the magnitude of distortion, which are set to 8 and 4, respectively.

5.3 Comparison with State-of-the-Art Methods

We compare our method with other recent state-of-the-art (SOTA) domain generalization methods on three public standard benchmark datasets, *i.e.*, PACS, VLCS, and OfficeHome.

Results on PACS. We compare our PLACE dropout with SOTA methods of both dropout and domain generalization on the PACS dataset. Tab. 1 and Tab. 2 present the results using ResNet-18

Table 1. Comparison of performance (%) among different methods using ResNet-18 on PACS [19]. The best and second-best are **bolded** and underlined, respectively.

Method	Venue	Art	Cartoon	Sketch	Photo	Avg.
C-Drop [26] [†]	ICCV 2017	79.64	76.49	72.37	95.93	81.11
WCD [14] [†]	AAAI 2019	81.56	78.24	75.53	94.99	82.58
RSC [17] [†]	ECCV 2020	82.03	77.39	75.64	95.63	82.67
I-Drop [37] [†]	ICML 2020	80.27	76.54	76.38	96.11	82.33
DSON [36]	ECCV 2020	84.67	77.65	82.23	95.87	85.11
SFA-A [20]	ICCV 2021	81.20	77.80	73.70	93.90	81.70
MixStyle [68]	ICLR 2021	84.10	78.80	75.90	96.10	83.70
pAdaIN [29]	CVPR 2021	81.74	76.91	75.13	96.29	82.51
SagNet [28]	CVPR 2021	83.58	77.66	76.30	95.47	83.25
FACT [54]	CVPR 2021	<u>85.37</u>	78.38	79.15	95.15	84.51
StableNet [64]	CVPR 2021	81.74	<u>79.91</u>	80.50	<u>96.53</u>	84.69
EFDMix [65]	CVPR 2022	83.90	79.40	75.00	96.80	83.90
StyleNeophile [18]	CVPR 2022	84.41	79.25	83.27	94.93	85.47
IRMCon-IPW [32]	ECCV 2022	81.10	77.30	76.60	95.40	82.60
I ² -ADR [25]	ECCV 2022	82.90	80.80	<u>83.50</u>	95.00	<u>85.60</u>
DeepAll [67] (<i>our imple.</i>)	AAAI 2020	80.19	77.19	73.48	95.71	81.64
+ PLACE dropout	Ours	82.60	78.33	81.47	95.65	84.51
DeepAll ⁺⁺	Ours	83.64	78.28	81.41	96.34	84.92
+ PLACE dropout	Ours	85.40	79.69	83.97	96.23	86.32

and ResNet-50 as backbones, respectively. In Tab. 1, PLACE dropout significantly outperforms the SOTA dropout-based method (RSC [17]) by 1.84% (84.51% vs. 82.67%), which develops a self-challenging algorithm by discarding over-dominant features with large gradients to encourage the model to rely more on the remaining features. The superior performance of PLACE dropout compared to previous dropout-based methods highlights the importance of utilizing dropout in multiple layers. Furthermore, when compared to SOTA DG methods, PLACE dropout achieves competitive performance and outperforms the baseline by a significant margin of 2.87% (84.51% vs. 81.64%). While integrated with our proposed DeepAll⁺⁺, which incorporates existing image-level and feature-level augmentation methods, PLACE dropout further improves the model generalization and outperforms the SOTA DG method I²-ADR [25] by 0.72% (86.32% vs. 85.60%). It is noteworthy that our method achieves excellent performance with introducing no extra parameters and little training time, proving its efficiency and superiority. Tab. 2 presents the results on PACS using ResNet-50 as the backbone. As observed, PLACE dropout continues to achieve substantial improvements over DeepAll by 3.15% (87.83% vs. 84.68%) and DeepAll⁺⁺ by 1.93% (89.03% vs. 87.10%). We notice that on DeepAll⁺⁺ in Tab. 2, PLACE dropout slightly degrades the performance on Cartoon, which could be attributed to the excessive perturbations by the combination of PLACE dropout and the two augmentation methods employed. Considering that different intensities of perturbations could lead to variations in the model performance across different domains (as evidenced by Tab. 9), the performance of our method on cartoon could be improved by fine-tuning perturbation strength

Table 2. Comparison of performance (%) among different methods using ResNet-50 on PACS [19]. The best and second-best are **bolded** and underlined, respectively.

Method	Venue	Art	Cartoon	Sketch	Photo	Avg.
RSC [17]	ECCV 2020	84.13	79.83	82.32	95.35	85.41
EISNet [47]	ECCV 2020	86.64	81.53	78.07	97.11	85.84
DSON [36]	ECCV 2020	87.04	80.62	82.90	95.99	86.64
MDGHybrid [23]	ICML 2021	86.74	82.32	82.66	98.36	87.52
FACT [54]	CVPR 2021	89.63	81.77	84.46	96.75	88.15
PCL [57]	CVPR 2022	<u>90.20</u>	83.90	82.60	<u>98.10</u>	<u>88.70</u>
EFDMix [65]	CVPR 2022	90.60	82.50	76.40	<u>98.10</u>	86.90
I ² -ADR [25]	ECCV 2022	88.50	83.20	<u>85.80</u>	95.20	88.20
DeepAll [67] (<i>our imple.</i>)	AAAI 2020	85.55	79.82	76.69	96.66	84.68
+ PLACE dropout	Ours	87.55	83.11	83.48	97.19	87.83
DeepAll ⁺⁺	Ours	85.11	<u>83.23</u>	83.58	96.47	87.10
+ PLACE dropout	Ours	89.03	83.04	86.82	97.22	89.03

accordingly (*i.e.*, decreasing the dropout rate). Moreover, our approach surpasses the second-best method PCL [57] by 0.33% (89.03% vs. 88.70%), demonstrating the stability and effectiveness of PLACE dropout even when incorporated into a larger network like ResNet-50. Noting that certain methods exhibit slightly better performance than ours on specific domains, , likely due to the utilization of specialized structures or assumptions tailored for those particular domains. In contrast, our method is an assumption-free technique, which is orthogonal to existing methods. The results in Tab. 5 further prove that combining existing methods with PLACE dropout can lead to enhanced generalization performance.

Results on VLCS. As shown in Tab. 3, our PLACE dropout outperforms the baseline DeepAll by a significant margin of 4.96% (77.44% vs. 72.48%) and achieves competitive performance with the SOTA method StableNet [64] while introducing no additional training parameters. Furthermore, based on the DeepAll⁺⁺, our PLACE dropout surpasses the SOTA method StableNet [64] by a margin of 0.11% (77.76% vs. 77.65%) on average, demonstrating the effectiveness of our method. Specifically, our method outperforms StableNet on the Caltech and Pascal domains with considerable improvements of 1.56% (98.23% vs. 96.67%) and 2.71% (76.30% vs. 73.59%), respectively, but slightly underperforms on the other two domains. The results could be attributed to the different strategies employed by our method and StableNet. Our method randomly drops task-relevant features to mitigate the overfitting issue, while StableNet discards the task-irrelevant features for stable learning, causing their different effects on the generalization ability of the model.

Results on OfficeHome. We conducted experiments on OfficeHome and reported the results in Tab. 4. OfficeHome is a more challenging benchmark for DG due to relatively smaller domain shifts and a larger number of categories compared to PACS and VLCS datasets. Despite these challenges, our method still performs competitively with the latest DG methods on this benchmark, which improves the baseline performance by 5.77% (66.04% vs. 60.27%) without extra computational cost. We observe that PLACE dropout performs well on the Artistic, Product, and Realworld tasks,

[†]The dropout-based regularization methods.

Table 3. Comparison of performance (%) among different methods using ResNet-18 on VLCS [42]. The best and second-best are **bolded** and underlined, respectively.

Method	Venue	Caltech	LableMe	Pascal	Sun	Avg.
JiGen [4]	CVPR 2019	96.17	62.06	70.93	71.40	75.14
MMLD [24]	AAAI 2020	97.01	62.20	73.01	<u>72.49</u>	76.18
RSC [17]	ECCV 2020	95.83	63.74	71.86	72.12	75.89
StableNet [64]	CVPR 2021	96.67	65.36	73.59	74.97	<u>77.65</u>
DeepAll [67] (<i>our imple.</i>)	AAAI 2020	91.86	61.81	67.48	68.77	72.48
+ PLACE dropout	Ours	98.14	64.12	<u>75.45</u>	72.06	77.44
DeepAll ⁺⁺	Ours	98.27	62.02	74.09	71.29	76.42
+ PLACE dropout	Ours	<u>98.23</u>	<u>64.38</u>	76.30	72.12	77.76

Table 4. Comparison of performance (%) among different methods using ResNet-18 on Office-Home [43]. The best and second-best are **bolded** and underlined, respectively.

Method	Venue	Artistic	Clipart	Product	Real	Avg.
RSC [17]	ECCV 2020	57.70	48.58	72.59	74.17	63.26
MixStyle [68]	ICLR 2021	58.70	53.40	74.20	75.90	65.50
SagNet [28]	CVPR 2021	<u>60.20</u>	45.38	70.42	73.38	62.34
FACT [54]	CVPR 2021	60.34	54.85	<u>74.48</u>	<u>76.55</u>	<u>66.56</u>
StyleNeophile [18]	CVPR 2022	59.55	55.01	73.57	75.52	65.89
COMEN [5]	CVPR 2022	57.60	<u>55.80</u>	75.50	76.90	66.50
DeepAll [67] (<i>our imple.</i>)	AAAI 2020	52.06	46.12	70.45	72.45	60.27
+ PLACE dropout	Ours	<u>60.20</u>	54.02	73.80	76.13	66.04
DeepAll ⁺⁺	Ours	56.83	55.49	70.12	72.28	63.68
+ PLACE dropout	Ours	59.99	58.40	74.14	76.29	67.20

but achieves relatively mediocre results on the most difficult task, Clipart. The possible reason for this difference in performance is that the Clipart dataset contains more noise compared to the other three domains, which negatively impacts the model performance. As shown in Tab. 4, some data augmentation methods can combat this noise and perform better on the Clipart domain compared to other methods, *e.g.*, MixStyle [68] and FACT [54]. By combining PLACE dropout with the augmentation-based DeepAll⁺⁺, our method achieves a significant improvement over the baseline, *i.e.*, 12.28% (58.40% vs. 46.12%), and surpasses the nearest competitor COMEN [5] by 2.60% (58.40% vs. 55.80%) on the Clipart dataset. Furthermore, our method outperforms SOTA DG methods, *e.g.*, StyleNeophile [18] by 1.31% (67.20% vs. 65.89%) and FACT [54] by 0.64% (67.20% vs. 66.56%). The results further support the effectiveness of PLACE dropout in challenging DG tasks.

Comparisons with Dropout-based Methods. To comprehensively prove the effectiveness of PLACE dropout, we compare it with the SOTA dropout-based methods, including methods designed for domain generalization (RSC [17] [‡] and I-Drop [37]), supervised learning (C-Drop [26] and WCD

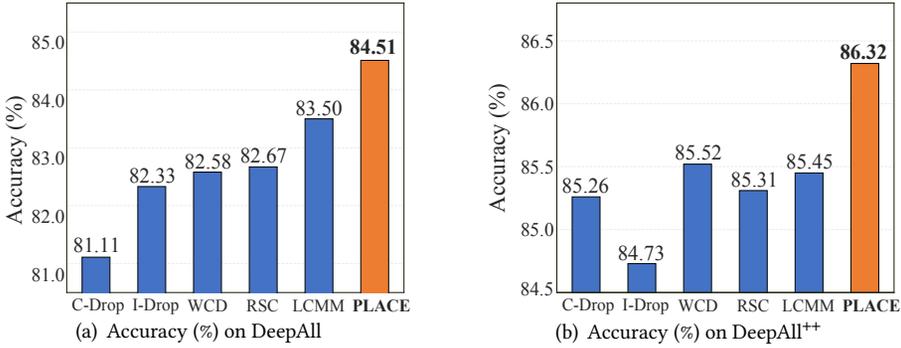


Fig. 4. Comparison of PLACE dropout and other dropout-based methods on PACS using ResNet-18.

Table 5. Effect (%) of incorporating PLACE dropout on other SOTA methods on the PACS dataset. For the CNN-based methods, *i.e.*, MixStyle [68] and FACT [54], we use the ResNet-18 as the backbone. We also validate the effectiveness of PLACE dropout on two SOTA MLP-like networks, *i.e.*, GFNet [34] and ViP [13].

Method	Art	Cartoon	Sketch	Photo	Avg.
MixStyle [68] (ICLR'20)	84.10	78.80	75.90	96.10	83.70
+ PLACE dropout	84.62	79.78	83.66	95.69	85.94
FACT [54] (CVPR'21)	85.37	78.38	79.15	95.15	84.51
+ PLACE dropout	85.50	82.08	83.18	94.91	86.42
GFNet [34] (NeurIPS'21)	89.37	84.74	79.01	97.94	87.76
+ PLACE dropout	92.19	86.86	84.22	98.44	90.43
ViP [13] (TPAMI'22)	88.09	84.22	82.41	98.38	88.27
+ PLACE dropout	91.66	85.47	85.62	99.04	90.45

[14]), and image compression (LMMD [22]). We conduct comparative experiments on both DeepAll and DeepAll⁺⁺, and the results are presented in Fig. 4. As shown in Fig. 4(a), our PLACE dropout achieves a significant improvement of 1.01% (84.51% vs. 83.50%) compared to the second-best approach LCMM, which utilizes a learnable vector to explicitly discard redundant channels. This enhancement is attributed to our method's ability to generate diverse data variants across multiple layers, effectively mitigating overfitting on source domains and promoting model generalization. Furthermore, in Fig. 4(b), when applied to DeepAll⁺⁺, our PLACE dropout still achieves SOTA performance and outperforms all other dropout-based methods, highlighting its superiority in promoting model generalization performance. Besides, it is worth noting that our DeepAll⁺⁺ significantly benefits most of the other dropout-based methods, *e.g.*, the performance of WCD [14] is boosted by 2.94% (85.52% vs. 82.58%) on DeepAll⁺⁺ compared to that on DeepAll. We also observe a performance decline of I-Drop on DeepAll⁺⁺, possibly due to that on the strong baseline, the regions with relatively less information could still contain important semantic information. In conclusion, PLACE dropout surpasses the SOTA dropout-based methods of DG on both DeepAll

[‡]We rerun the official codes of RSC with the same hyperparameters as mentioned in [17] but fail to reproduce the reported results, which may be due to the different hardware environment. This problem is also reported in [29, 49, 54]. To be fair, we compare our method with RSC in our environment.

Table 6. Effect (%) of each component in PLACE dropout on PACS with ResNet-18 or ResNet-50 as the backbone architecture. For simplicity, we denote the channel-wise dropout as C, the layer-wise dropout as L, and the progressive scheme as P. The best and second-best are **bolded** and underlined, respectively.

Method	C	L	P	Art	Cartoon	Sketch	Photo	Avg.
DeepAll								
Baseline	-	-	-	80.19	77.19	73.48	95.71	81.64
Variant 1	✓	-	-	79.02	76.44	79.93	92.02	81.85
Variant 2	✓	-	✓	80.34	78.07	<u>80.91</u>	94.17	83.38
Variant 3	-	✓	-	79.95	76.82	75.90	95.28	81.99
Variant 4	-	✓	✓	81.03	77.18	76.04	95.25	82.38
Variant 5	✓	✓	-	<u>82.23</u>	78.34	80.40	95.23	<u>84.05</u>
PLACE dropout	✓	✓	✓	82.60	<u>78.33</u>	81.47	<u>95.65</u>	84.51
DeepAll ⁺⁺								
Baseline	-	-	-	83.64	78.28	81.41	<u>96.34</u>	84.92
Variant 1	✓	-	-	81.71	76.01	82.42	94.33	83.62
Variant 2	✓	-	✓	82.63	77.56	83.27	95.09	84.64
Variant 3	-	✓	-	83.69	77.94	82.92	96.65	85.30
Variant 4	-	✓	✓	84.47	<u>79.25</u>	83.22	96.11	85.76
Variant 5	✓	✓	-	<u>84.57</u>	79.18	84.16	96.23	<u>86.03</u>
PLACE dropout	✓	✓	✓	85.40	79.69	<u>83.97</u>	96.23	86.32

and DeepAll⁺⁺, demonstrating the effectiveness of our method in improving the generalization ability of the model. Additionally, we provide preliminary evidence that dropout-based methods can cooperate with other data augmentation techniques to enhance model performance in practical scenarios [66], which we hope will inspire future research in DG.

Incorporating PLACE dropout with other SOTA DG methods. To validate the generalization of PLACE dropout, we conducted experiments by integrating it with other SOTA DG methods, including CNN-based methods MixStyle [68] and FACT [54], as well as MLP-like networks GFNet [34] and ViP [13]. As shown in Tab. 5, when combined with the SOTA CNN-based methods, PLACE dropout consistently achieves substantial performance improvements, with a margin of 2.24% (85.94% vs. 83.70%) for MixStyle [68] and 1.91% (86.42% vs. 84.51%) for FACT [54]. Furthermore, when applied to GFNet [34], a novel network exploring feature dependencies in frequency space, PLACE dropout effectively enhances the performance with an improvement of 2.67% (90.43% vs. 87.76%). Besides, our method outperforms ViP [13], which learns long-range dependencies in both height and width directions through linear projections, by 2.18% (90.45% vs. 88.27%). The results indicate that PLACE dropout complements DG methods effectively, confirming that mitigating overfitting on source domains can further enhance model generalization ability.

5.4 Ablation Studies

Impact of different components. We conducted experiments to investigate the contribution of each component in PLACE dropout on PACS using ResNet-18 as the backbone. The results are presented in Tab. 6, and each component has demonstrated significant importance in our method.

Starting with the baseline, Variant 1 represents Channel-wise dropout, which applies dropout in the channel dimension to multiple layers simultaneously. It achieves a slight improvement on DeepAll but results in performance degradation on DeepAll⁺⁺, indicating that the simultaneous use of dropout in multiple layers may cause excessive information loss and hinder model training. Nevertheless, the results also show that Channel-wise dropout can generate strong regularization effects, as evidenced by its notable improvement on Sketch. On the other hand, Layer-wise dropout (Variant 3) involves randomly selecting a network layer (Layers 1-3) at each iteration to apply dropout in the pixel dimension instead of the channel dimension. Pure Layer-wise dropout exhibited a slight improvement of 0.35% over the Baseline (81.99% vs. 81.64%), indicating the limited regularization strength of dropout in the pixel dimension. To address these issues, we developed Layer-wise and Channel-wise dropout (Variant 5), which significantly improved performance by 2.41% on DeepAll (84.05% vs. 81.64%) and by 1.11% on DeepAll⁺⁺ (86.03% vs. 84.92%). The results prove that Layer-wise and Channel-wise dropout can effectively mitigate the detrimental effects of noise on model training while providing strong regularization to combat overfitting. Besides, we extended Variants 1, 3, and 5 to Variants 2, 4, and PLACE dropout with the progressive scheme to verify its effectiveness. Based on Variant 1, the progressive scheme achieved a relatively significant improvement on both DeepAll (by 1.47%) and DeepAll⁺⁺ (by 1.02%), owing to it can reduce the hindrance of multi-layers dropout to model convergence at the beginning of training. On other variants, the progressive scheme also consistently improves performance by about 0.5% on average, verifying its effectiveness for adequately mitigating overfitting. Specifically, we observe that the Baseline achieves remarkable performance on Photo. It could be due to the similarity between Photo and the pre-trained dataset ImageNet [17, 54], which allows the model performance to easily saturate on this domain. However, when applying dropout during training, it would generate additional noise and lead to performance fluctuations, *e.g.*, a slight performance decrease on Photo. The issue could be alleviated by adjusting the random seed or the dropout rate. Nevertheless, our method enhances the accuracy on other challenging domains and achieves the best overall performance, proving that the three modules are indispensable for superior generalization ability.

Different dimensions to conduct dropout. We investigate different variants of PLACE dropout that apply dropout in various dimensions, including pixel-wise dropout (PLACE-E), spatial-wise dropout (PLACE-S), and channel-wise dropout (PLACE-C). We also conducted experiments to combine PLACE-S and PLACE-C (denoted as PLACE-S/C), which randomly selects either PLACE-S or PLACE-C to perform dropout with a probability of 50% at each iteration. As shown in Tab. 7, all the variants show performance gains, with PLACE-S and PLACE-C achieving larger improvements compared to PLACE-E. The results suggest that structural dropouts, *i.e.*, spatial- and channel-wise dropout, contribute to better resistance against overfitting. Moreover, PLACE-C shows the best performance on both DeepAll and DeepAll⁺⁺, confirming the discussion in Section 3.2 that channel-wise dropout can effectively perturb the patterns learned by the model, thereby providing a strong regularization effect. Besides, we observed that the combined usage of PLACE-C and PLACE-S (PLACE-S&C) yields improved results compared to using PLACE-S alone, but slightly reduces the accuracy compared to PLACE-C. A possible reason for the observation is that the regularization effects of PLACE-S and PLACE-C are not orthogonal, as both methods provide regularization by discarding a subset of feature maps. Since numerous regions within the spatial dimension contain redundant information, as already discussed in Section 3.2, the regularization effect of PLACE-S could not be as effective as that of PLACE-C, thus slightly weakening the regularization effect when combined with PLACE-C. To effectively combat the overfitting on source domains, we apply PLACE dropout in the channel dimension to further improve the model generalization.

Where to apply PLACE dropout. To investigate which layers to apply PLACE dropout, we conducted the experiments on PACS using the ResNet-18 architecture. Given that a standard ResNet

Table 7. Comparison of different variants of PLACE dropout that conduct dropout in different dimensions, *i.e.*, pixel dimension (PLACE-E), spatial dimension (PLACE-S), channel dimension (PLACE-C), and spatial or channel dimension (PLACE-S&C). The experiments are conducted on PACS using ResNet-18 as the backbone architecture. The best and second-best are **bolded** and underlined, respectively.

Method	Art	Cartoon	Sketch	Photo	Avg.
DeepAll	80.19	77.19	73.48	95.71	81.64
+ PLACE-E	81.20	<u>78.11</u>	73.26	95.40	81.99
+ PLACE-S	82.08	77.90	78.18	95.07	83.31
+ PLACE-C	82.60	78.33	81.47	<u>95.65</u>	84.51
+ PLACE-S/C	<u>82.13</u>	77.98	<u>78.95</u>	95.39	<u>83.61</u>
DeepAll ⁺⁺	83.64	78.28	81.41	96.34	84.92
+ PLACE-E	83.77	79.02	82.51	96.15	85.36
+ PLACE-S	84.67	79.27	82.82	96.07	85.71
+ PLACE-C	85.40	79.69	83.97	<u>96.23</u>	86.32
+ PLACE-S/C	<u>84.97</u>	<u>79.34</u>	<u>83.38</u>	96.13	<u>85.96</u>

Table 8. Effect (%) of different layers where PLACE dropout is inserted on the PACS dataset with ResNet-18 or ResNet-50 as the backbone architecture. For notation, L1 means PLACE dropout is applied after the first residual block; {L1, L2} means PLACE dropout is applied after the first or the second blocks randomly. The best and second-best are **bolded** and underlined, respectively.

Method	Position				Baseline	
	L1	L2	L3	L4	DeepAll	DeepAll ⁺⁺
Baseline	-	-	-	-	81.64 ± 0.39	84.92 ± 0.16
Variant 1	✓	-	-	-	82.81 ± 0.03	85.67 ± 0.13
Variant 2	-	✓	-	-	83.63 ± 0.15	85.66 ± 0.04
Variant 3	-	-	✓	-	83.25 ± 0.13	85.91 ± 0.12
Variant 4	-	-	-	✓	81.48 ± 0.40	84.83 ± 0.34
Variant 5	✓	✓	-	-	83.58 ± 0.15	85.69 ± 0.31
Variant 6	✓	✓	✓	-	84.51 ± 0.09	86.32 ± 0.15
Variant 7	✓	✓	✓	✓	<u>83.85 ± 0.12</u>	<u>86.06 ± 0.20</u>

model has four residual blocks denoted by L1-4, we trained various models with PLACE dropout applied to different layers. The results are presented in Tab. 8. We first tested the performance of the models with PLACE dropout in one single layer, *i.e.*, Variant 1-4. The model with PLACE dropout in the first three layers achieves consistent improvement, indicating that PLACE dropout can alleviate the overfitting issue in every network layer. However, inserting dropout in the 4-th layer reduces the model accuracy on both DeepAll and DeepAll⁺⁺, which might be due to L4 being near the prediction layer and tending to capture label-related information, thus randomly dropout is likely to break the inherent label space and hamper the model generalization ability. We further tested the model performance with PLACE dropout in multiple layers. By comparing Variant 6 with 7, we observed that dropout in the 4-th layer reduces the performance, which is likely related to the randomness of dropout and could be overcome by carefully guiding the dropout [14, 17]. Finally,

Table 9. Effect(%) of PLACE dropout with different P_{\max} to model performance on the PACS dataset with ResNet-18 as the backbone architecture. The best and second-best are **bolded** and underlined, respectively.

	P_{\max}	Art	Cartoon	Sketch	Photo	Avg.
DeepAll	20%	81.51	77.68	78.71	95.67	83.39
	25%	<u>82.31</u>	78.09	79.98	95.64	84.01
	33%	82.60	78.33	81.47	<u>95.65</u>	84.51
	40%	82.14	<u>78.27</u>	<u>80.81</u>	95.29	<u>84.13</u>
DeepAll ⁺⁺	20%	83.94	79.14	83.13	<u>96.71</u>	85.73
	25%	<u>84.63</u>	<u>79.89</u>	82.85	96.70	86.02
	33%	85.40	79.69	83.97	96.23	86.32
	40%	83.94	79.91	<u>83.56</u>	96.77	<u>86.05</u>

Variant 6 with PLACE dropout in the 1-st, 2-nd, and 3-rd layers achieves the best performance and yields the baseline 2.87% (84.51% vs. 81.64%) on DeepAll and 1.14% (86.06% vs. 84.92%) on DeepAll⁺⁺, aligning well with the discussion for the layer-wise dropout in Section 3.2.

Sensitivity to the maximum dropout rate. To investigate the optimal value of the hyperparameter P_{\max} , which indicates the maximum proportion of dropping channels, we conducted the experiments on PACS with ResNet-18 as the backbone. The experimental results on both DeepAll and DeepAll⁺⁺ are presented in Tab. 9. Remarkably, the model performance remains relatively consistent across different values of P_{\max} , and it achieves excellent performance when P_{\max} is set between 25% and 40% on both baselines. Particularly, our method achieves the highest average accuracy on both DeepAll and DeepAll⁺⁺ when P_{\max} is set to 33%, indicating that the baseline does not significantly influence the selection of P_{\max} . Considering that the maximum dropout rate P_{\max} influences the strength of the regularization effect during training, the optimal value should not be excessively large or small. A too large P_{\max} could lead to excessive information loss during training, hindering the model from learning discriminative information from source domains. Conversely, a too small P_{\max} might fail to effectively address the overfitting issue on source domains. Based on the experimental results, We adopt 33% as the default setting in all experiments if not specified.

5.5 Further analysis

Layer-wise dropout helps reduce domain gap. To analyze the effect of layer-wise dropout on the model generalization, we conducted experiments to measure the inter-domain discrepancy of extracted features from intermediate layers with PLACE dropout on both DeepAll and DeepAll⁺⁺. The discrepancy was computed using Eq. (2). We also compared Layer-wise dropout (LD) with Multi-layers dropout (MD), where dropout with lower dropping rates is simultaneously applied to L1 – 3 of the model. As shown in Fig. 5, layer-wise dropout effectively reduces the inter-domain distances across all network layers, leading to a significant improvement in model generalization, regardless of whether applied to DeepAll or DeepAll⁺⁺. Specifically, when we employ MD instead of LD, we observe that the values of $d_{drop}^i - d_{bal}^i$ in the 2-th layer are positive, indicating an increase in inter-domain discrepancy on that specific layer. This aligns well with the discussion in Sec. 3.2, indicating that simultaneous dropout across multiple layers could lead to cumulative noise that hinders some layers from fully learning, causing an increase in the inter-domain gap on certain layers, *e.g.*, the 2-th layer. To tackle this issue, we design the Layer-wise dropout, where we randomly select one layer to adopt dropout at each iteration, thereby mitigating excessive noise that could impede the learning process. As demonstrated in "LD0.33" of Fig. 5, Layer-wise

Table 10. Comparison of layer-wise dropout and multi-layers dropout on PACS with ResNet-18 (Best in **bold**). For simplicity, We denote multi-layers dropout as MD and layer-wise dropout as LD, *e.g.*, MD0.11 is multi-layers dropout with P_{\max} as 0.11. The best and second-best are **bolded** and underlined, respectively.

	Method	Art	Cartoon	Sketch	Photo	Avg.
DeepAll	MD0.11	81.91	78.08	78.86	<u>95.64</u>	83.62
	MD0.22	<u>82.42</u>	<u>78.44</u>	80.27	94.37	83.88
	MD0.33	80.25	78.50	81.64	94.23	83.65
	LD0.33	82.60	78.33	<u>81.47</u>	95.65	84.51
DeepAll++	MD0.11	84.18	<u>78.75</u>	83.33	<u>95.93</u>	85.55
	MD0.22	<u>84.25</u>	78.63	84.28	95.91	<u>85.77</u>
	MD0.33	82.45	77.22	83.33	94.97	84.49
	LD0.33	85.40	79.69	<u>83.97</u>	96.23	86.32

dropout achieves a smaller domain gap than multi-layers dropout, especially for the 2-th layer, indicating that it can narrow the domain gap between source and target domains more effectively than simply utilizing dropout for multiple layers with a lower dropping rate. Furthermore, we provide detailed comparison results on both DeepAll and DeepAll++ in Tab. 10. The model with layer-wise dropout consistently achieves superior performance compared to the model with multi-layers dropout. These results verify that simply utilizing dropout for multiple layers with a lower dropping rate cannot lead to the same regularization as our layer-wise dropout. Additionally, the results demonstrate that layer-wise dropout can effectively reduce the domain gap between source and target domains, thereby improving the model generalization ability.

PLACE dropout indicates tighter generalization bound. We investigate the effect of PLACE dropout on the model generalization bound. As is shown in Fig. 6(a), InfoDrop, RSC, and PLACE dropout can all increase the difficulty of model training, but layer-wise and channel-wise dropout exhibits the strongest regularization effect in the late stages of training, and the progressive scheme ensures stable model convergence in the early training stages. Fig.6(b) and Fig.6(c) illustrate the accuracy and loss on the test set with an increasing number of epochs for each method. Our proposed method consistently reduces the model prediction loss and enhances accuracy more effectively

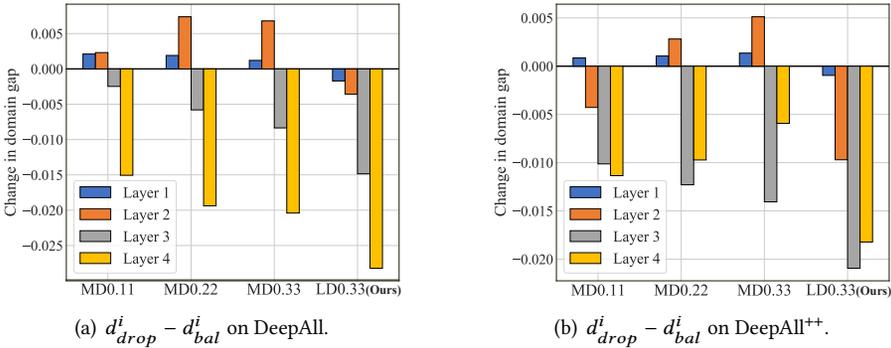


Fig. 5. Comparison of inter-domain discrepancy between multi-layers dropout and our layer-wise dropout. We denote multi-layers dropout as MD and layer-wise dropout as LD, *e.g.*, MD0.11 is multi-layers dropout with drop ratio p as 0.11. (a) and (b) are the results on DeepAll and DeepAll++, respectively. The smaller $d_{drop}^i - d_{bal}^i$ indicates the larger reduction of the domain gap and the better generalization of the model.

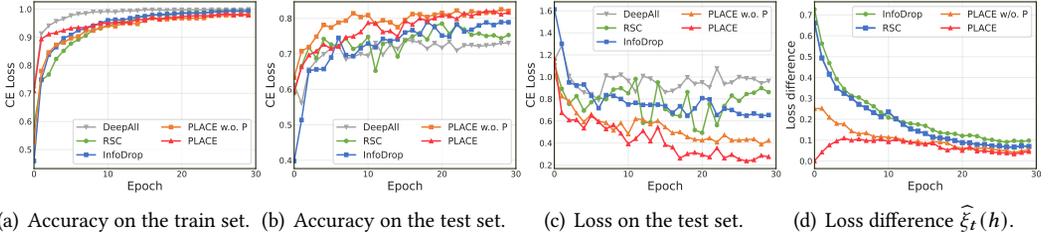


Fig. 6. The influence of PLACE dropout on the procedure of model training and test. We compare PLACE dropout with the baseline DeepAll and RSC [17]. We also investigate the effect of the progressive scheme. (a) presents the accuracy of each method on the train set, (b) and (c) are the accuracy and loss of each method on the test set, and (d) shows the change of loss difference $\hat{\xi}_t(h)$ as the training progresses.

Table 11. The distribution discrepancy ($\times 10$) of inter-domain (across all source domains) and intra-class (for all classes in source domains) on PACS. Note that the lower the value indicates the better performance.

Method	PACS		VLCS		OfficeHome	
	Inter-domain	Intra-class	Inter-domain	Intra-class	Inter-domain	Intra-class
DeepAll	2.43	2.95	1.58	2.41	0.76	1.49
+ PLACE	1.96	2.48	1.54	2.35	0.71	1.44
DeepAll ⁺⁺	1.63	2.13	1.47	2.34	0.69	1.40
+ PLACE	1.53	2.04	1.38	2.27	0.63	1.36

on the test set compared to other methods. Furthermore, the progressive scheme, when combined with layer-wise and channel-wise dropout, further improves model performance by protecting early-stage convergence. By employing the progressive scheme, the model learns "smarter" weights, avoiding learning potentially harmful information [26] and effectively addressing overfitting. Fig. 6(d) presents the changes of $\hat{\xi}_t(h)$ during training, reflecting the evolution of the generalization error bound. PLACE dropout enhances the model robustness against perturbations caused by dropout, *i.e.*, reducing the empirical loss difference $\hat{\xi}_t(h)$, which aligns well with the theoretical discussion that layer-wise and channel-wise dropout can yield a tight generalization error bound. The empirical results also validate the effectiveness of the progressive scheme.

Distribution Discrepancy of Feature Representations To investigate the influence of our PLACE dropout, we examine the inter-domain (across all source domains) and intra-class (inter-domain distance for the same-class samples) discrepancy. The inter-domain discrepancy and the intra-class discrepancy are computed by the following formulations:

$$d_{inter-domain} = \frac{2}{K(K-1)} \sum_{m=1}^K \sum_{n=m+1}^K \|\bar{\mathbf{z}}_m - \bar{\mathbf{z}}_n\|_2, \quad (11)$$

$$d_{intra-class} = \frac{1}{CK} \sum_{c=1}^C \sum_{k=1}^K \|\bar{\mathbf{z}}_{k,c} - \bar{\mathbf{z}}_k\|_2, \quad (12)$$

where $\bar{\mathbf{z}}_k, \bar{\mathbf{z}}_m, \bar{\mathbf{z}}_n$ are the feature means of all classes in the k -th, m -th, n -th source domain, respectively. $\bar{\mathbf{z}}_{k,c}$ is the averaged representation of the c -th class in the k -th source domain, K denotes the number of source domains and C is the number of categories. As shown in Tab. 11, our method can effectively reduce the domain gap of source domains on both DeepAll and DeepAll⁺⁺, especially

on PACS which shows a large domain gap. Even on the more challenging datasets VLCS and OfficeHome, our method can still effectively close the domain gap. Besides, PLACE dropout can also reduce intra-class discrepancy, which indicates that it can effectively capture discriminative features. The results demonstrate the effectiveness of our method in mitigating the overfitting issue, which is beneficial to reducing the intra-class and inter-domain distances among source domains.

6 CONCLUSIONS

In this paper, we propose a simple yet effective dropout-based method for domain generalization to mitigate the overfitting problem on source domains and help the model generalize well to unseen target domains. We name our method Progressive Layer-wise and Channel-wise (PLACE) dropout. At each iteration, PLACE dropout randomly selects one middle network layer and then randomly mutes its channels by a proportion that gradually increases as the training progresses. We also provide theoretical analysis for layer-wise and channel-wise dropout and prove its effectiveness in reducing the generalization error bound. Extensive experiments on various datasets demonstrate that our method outperforms other related SOTA DG methods with no extra network parameters and almost no increment in computing cost. In future work, we will further analyze the overfitting problem in the DG task and develop advanced methods that enable the model to adaptively learn suitable layers and dropout rates based on input data.

ACKNOWLEDGMENTS

The work is supported by NSFC Program (62222604, 62206052, 62192783) and Jiangsu Natural Science Foundation Project (BK20210224).

REFERENCES

- [1] Yogesh Balaji, Swami Sankaranarayanan, and Rama Chellappa. 2018. Metareg: Towards domain generalization using meta-regularization. In *NeurIPS*.
- [2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *ICML*.
- [3] Francesco Cappio Borlino, Antonio D’Innocente, and Tatiana Tommasi. 2021. Rethinking domain generalization baselines. In *ICPR*.
- [4] Fabio M Carlucci, Antonio D’Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. 2019. Domain generalization by solving jigsaw puzzles. In *CVPR*.
- [5] Chaoqi Chen, Jiongcheng Li, Xiaoguang Han, Xiaoqing Liu, and Yizhou Yu. 2022. Compound Domain Generalization via Meta-Knowledge Encoding. In *CVPR*.
- [6] Yang Chen, Yu Wang, Yingwei Pan, Ting Yao, Xinmei Tian, and Tao Mei. 2021. A Style and Semantic Memory Mechanism for Domain Generalization. In *ICCV*.
- [7] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. 2020. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPRW*.
- [8] Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. 2019. Domain generalization via model-agnostic learning of semantic features. In *NeurIPS*.
- [9] Antonio D’Innocente and Barbara Caputo. 2018. Domain generalization with domain-specific aggregation modules. In *GCPR*.
- [10] Xinjie Fan, Qifei Wang, Junjie Ke, Feng Yang, Boqing Gong, and Mingyuan Zhou. 2021. Adversarially Adaptive Normalization for Single Domain Generalization. In *CVPR*.
- [11] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. 2018. Dropblock: A regularization method for convolutional networks. In *NeurIPS*.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*.
- [13] Qibin Hou, Zihang Jiang, Li Yuan, Ming-Ming Cheng, Shuicheng Yan, and Jiashi Feng. 2022. Vision permutator: A permutable mlp-like architecture for visual recognition. *TPAMI* (2022).
- [14] Saihui Hou and Zilei Wang. 2019. Weighted channel dropout for regularization of deep convolutional neural network. In *AAAI*.

- [15] Jiaxing Huang, Dayan Guan, Aoran Xiao, and Shijian Lu. 2021. Rda: Robust domain adaptation via fourier adversarial attacking. In *ICCV*.
- [16] Xun Huang and Serge Belongie. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*.
- [17] Zeyi Huang, Haohan Wang, Eric P Xing, and Dong Huang. 2020. Self-challenging improves cross-domain generalization. In *ECCV*.
- [18] Juwon Kang, Sohyun Lee, Namyup Kim, and Suha Kwak. 2022. Style Neophile: Constantly Seeking Novel Styles for Domain Generalization. In *CVPR*.
- [19] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. 2017. Deeper, broader and artier domain generalization. In *ICCV*.
- [20] Pan Li, Da Li, Wei Li, Shaogang Gong, Yanwei Fu, and Timothy M Hospedales. 2021. A Simple Feature Augmentation for Domain Generalization. In *ICCV*.
- [21] Zekun Li, Lei Qi, Yinghuan Shi, and Yang Gao. 2023. IOMatch: Simplifying Open-Set Semi-Supervised Learning with Joint Inliers and Outliers Utilization. In *ICCV*.
- [22] Lin Liu, Mingming Zhao, Shanxin Yuan, Wenlong Lyu, Wengang Zhou, Houqiang Li, Yanfeng Wang, and Qi Tian. 2023. Exploring Effective Mask Sampling Modeling for Neural Image Compression. *arXiv preprint arXiv:2306.05704* (2023).
- [23] Divyat Mahajan, Shruti Tople, and Amit Sharma. 2021. Domain generalization using causal matching. In *ICML*.
- [24] Toshihiko Matsuura and Tatsuya Harada. 2020. Domain generalization using a mixture of multiple latent domains. In *AAAI*.
- [25] Rang Meng, Xianfeng Li, Weijie Chen, Shicai Yang, Jie Song, Xinchao Wang, Lei Zhang, Mingli Song, Di Xie, and Shiliang Pu. 2022. Attention diversification for domain generalization. In *ECCV*.
- [26] Pietro Morerio, Jacopo Cavazza, Riccardo Volpi, René Vidal, and Vittorio Murino. 2017. Curriculum dropout. In *ICCV*.
- [27] Shruti Nagpal, Maneet Singh, Richa Singh, and Mayank Vatsa. 2020. Attribute aware filter-drop for bias invariant classification. In *CVPRW*.
- [28] Hyeonseob Nam, HyunJae Lee, Jongchan Park, Wonjun Yoon, and Donggeun Yoo. 2021. Reducing Domain Gap by Reducing Style Bias. In *CVPR*.
- [29] Oren Nuriel, Sagie Benaim, and Lior Wolf. 2021. Permuted AdaIN: reducing the bias towards global statistics in image classification. In *CVPR*.
- [30] Sinno Jialin Pan and Qiang Yang. 2009. A survey on transfer learning. *TKDE* 22, 10 (2009), 1345–1359.
- [31] Sungheon Park and Nojun Kwak. 2016. Analysis on the dropout effect in convolutional neural networks. In *ACCV*.
- [32] Jiaxin Qi, Kaihua Tang, Qianru Sun, Xian-Sheng Hua, and Hanwang Zhang. 2022. Class is invariant to context and vice versa: on learning invariance for out-of-distribution generalization. In *ECCV*.
- [33] Lei Qi, Jiaqi Liu, Lei Wang, Yinghuan Shi, and Xin Geng. 2023. Unsupervised generalizable multi-source person re-identification: A Domain-specific adaptive framework. *PR* (2023).
- [34] Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. 2021. Global filter networks for image classification. In *NeurIPS*.
- [35] Christian Schreckenberger, Christian Bartelt, and Heiner Stuckenschmidt. 2019. iDropout: Leveraging deep taylor decomposition for the robustness of deep neural networks. In *OTM*.
- [36] Seonguk Seo, Yumin Suh, Dongwan Kim, Geeho Kim, Jongwoo Han, and Bohyung Han. 2020. Learning to optimize domain specific normalization for domain generalization. In *ECCV*.
- [37] Baifeng Shi, Dinghuai Zhang, Qi Dai, Zhanxing Zhu, Yadong Mu, and Jingdong Wang. 2020. Informative dropout for robust representation learning: A shape-bias perspective. In *ICML*.
- [38] Qinghongya Shi, Hong-Bo Zhang, Zhe Li, Ji-Xiang Du, Qing Lei, and Jing-Hua Liu. 2022. Shuffle-invariant Network for Action Recognition in Videos. *TOMM* 18, 3 (2022), 1–18.
- [39] Yang Shu, Zhangjie Cao, Chenyu Wang, Jianmin Wang, and Mingsheng Long. 2021. Open Domain Generalization with Domain-Augmented Meta-Learning. In *CVPR*.
- [40] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *JMLR* 15, 1 (2014), 1929–1958.
- [41] Jonathan Tompson, Ross Goroshin, Arjun Jain, Yann LeCun, and Christoph Bregler. 2015. Efficient object localization using convolutional networks. In *CVPR*.
- [42] Antonio Torralba and Alexei A Efros. 2011. Unbiased look at dataset bias. In *CVPR*.
- [43] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. 2017. Deep hashing network for unsupervised domain adaptation. In *CVPR*.
- [44] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John Duchi, Vittorio Murino, and Silvio Savarese. 2018. Generalizing to unseen domains via adversarial data augmentation. In *NeurIPS*.
- [45] Stefan Wager, Sida Wang, and Percy S Liang. 2013. Dropout training as adaptive regularization. In *NeurIPS*.

- [46] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and Philip Yu. 2022. Generalizing to unseen domains: A survey on domain generalization. *TKDE* (2022).
- [47] Shujun Wang, Lequan Yu, Caizi Li, Chi-Wing Fu, and Pheng-Ann Heng. 2020. Learning from extrinsic and intrinsic supervisions for domain generalization. In *ECCV*.
- [48] Xiran Wang, Jian Zhang, Lei Qi, and Yinghuan Shi. 2023. Generalizable Decision Boundaries: Dualistic Meta-Learning for Open Set Domain Generalization. In *ICCV*.
- [49] Yufei Wang, Haoliang Li, Lap-pui Chau, and Alex C Kot. 2021. Embracing the Dark Knowledge: Domain Generalization Using Regularized Knowledge Distillation. In *ACM MM*.
- [50] Yue Wang, Lei Qi, Yinghuan Shi, and Yang Gao. 2022. Feature-based Style Randomization for Domain Generalization. *TCSVT* 32, 8 (2022), 54595–5509.
- [51] Zijian Wang, Yadan Luo, Ruihong Qiu, Zi Huang, and Mahsa Baktashmotlagh. 2021. Learning to diversify for single domain generalization. In *ICCV*.
- [52] Guoqiang Wei, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. 2021. MetaAlign: Coordinating Domain Alignment and Classification for Unsupervised Domain Adaptation. In *CVPR*.
- [53] Lei Wu, Hefei Ling, Yuxuan Shi, and Baiyan Zhang. 2022. Instance Correlation Graph for Unsupervised Domain Adaptation. *TOMM* 18, 1s (2022), 1–23.
- [54] Qinwei Xu, Ruipeng Zhang, Ya Zhang, Yanfeng Wang, and Qi Tian. 2021. A Fourier-based Framework for Domain Generalization. In *CVPR*.
- [55] Yifan Xu, Kekai Sheng, Weiming Dong, Baoyuan Wu, Changsheng Xu, and Bao-Gang Hu. 2022. Towards Corruption-Agnostic Robust Domain Adaptation. *TOMM* 18, 4 (2022), 1–16.
- [56] Zhenzhen Yang, Pengfei Xu, Yongpeng Yang, and Bing-Kun Bao. 2021. A densely connected network based on U-Net for medical image segmentation. *TOMM* 17, 3 (2021), 1–14.
- [57] Xufeng Yao, Yang Bai, Xinyun Zhang, Yuechen Zhang, Qi Sun, Ran Chen, Ruiyu Li, and Bei Yu. 2022. PCL: Proxy-based Contrastive Learning for Domain Generalization. In *CVPR*.
- [58] Matthew D Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *ECCV*.
- [59] Yuyuan Zeng, Tao Dai, Bin Chen, Shu-Tao Xia, and Jian Lu. 2021. Correlation-based structural dropout for convolutional neural networks. *PR* (2021).
- [60] Dinghui Zhang, Kartik Ahuja, Yilun Xu, Yisen Wang, and Aaron Courville. 2021. Can Subnetwork Structure be the Key to Out-of-Distribution Generalization?. In *ICML*.
- [61] Jian Zhang, Lei Qi, Yinghuan Shi, and Yang Gao. 2022. More is Better: A Novel Multi-view Framework for Domain Generalization. In *ECCV*.
- [62] Jian Zhang, Lei Qi, Yinghuan Shi, and Yang Gao. 2022. MVDG: A Unified Multi-view Framework for Domain Generalization. In *ECCV*.
- [63] Jian Zhang, Lei Qi, Yinghuan Shi, and Yang Gao. 2023. DomainAdaptor: A Novel Approach to Test-time Adaptation. In *ICCV*.
- [64] Xingxuan Zhang, Peng Cui, Renzhe Xu, Linjun Zhou, Yue He, and Zheyang Shen. 2021. Deep Stable Learning for Out-Of-Distribution Generalization. In *CVPR*.
- [65] Yabin Zhang, Minghan Li, Ruihuang Li, Kui Jia, and Lei Zhang. 2022. Exact feature distribution matching for arbitrary style transfer and domain generalization. In *CVPR*.
- [66] Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. 2021. Domain generalization in vision: A survey. *arXiv preprint arXiv:2103.02503* (2021).
- [67] Kaiyang Zhou, Yongxin Yang, Timothy Hospedales, and Tao Xiang. 2020. Deep domain-adversarial image generation for domain generalisation. In *AAAI*.
- [68] Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. 2020. Domain Generalization with MixStyle. In *ICLR*.

A PROOF OF THEORETICAL ANALYSIS

A.1 Proof of Proposition 1

Proposition 1. Let $\widehat{\xi}_t(h)$ be the estimated value of $\xi(h)$ at iteration t , which is defined as: $\widehat{\xi}_t(h) = |L(\mathbf{z}_t, h_t) - L(\widetilde{\mathbf{z}}_t, h_t)|$. Given a sufficiently small learning rate η , if discarding structural features will increase the empirical loss at the current iteration t , it holds that the layer-wise and channel-wise dropout can continually decrease $\widehat{\xi}_t(h)$ and lead it to be a small number at the end of training.

Proof. To investigate the trend of $\widehat{\xi}_t(h)$ over iterations, we first utilize the Taylor series formula to derive the expansion of the empirical loss $L(\mathbf{z}_{t+1}, h)$ of the model at iteration $t + 1$:

$$L(\mathbf{z}_{t+1}, h_{t+1}) = L(\mathbf{z}_t, h_t) + \frac{\partial L(\mathbf{z}_t, h_t)}{\partial h_t} (h_{t+1} - h_t) + \frac{1}{2} \frac{\partial^2 L(\mathbf{z}_t, h_t)}{\partial^2 h_t} \|h_{t+1} - h_t\|^2 + \dots \quad (13)$$

Note that the update step size of the hypothesis h_{t+1} can be denoted as:

$$h_{t+1} - h_t = -\eta \widetilde{\mathbf{g}}_t, \quad \text{where } \widetilde{\mathbf{g}}_t = \frac{\partial L(\widetilde{\mathbf{z}}_t, h_t)}{\partial h_t}. \quad (14)$$

The assumption that discarding structural features will increase the empirical loss could be formulated as: $L(\widetilde{\mathbf{z}}_t, h_t) = \gamma_t L(\mathbf{z}_t, h_t)$, $\gamma_t \geq 1$. The small learning rate η allows us to discard the higher order terms in the Taylor expansion of empirical loss $L(\mathbf{z}_{t+1}, h)$ in Eq. (13), then we can obtain:

$$L(\mathbf{z}_{t+1}, h_{t+1}) = L(\mathbf{z}_t, h_t) - \frac{1}{\gamma_t} \eta \widetilde{\mathbf{g}}_t^2. \quad (15)$$

We can also get an approximation of $L(\widetilde{\mathbf{z}}_{t+1}, h_{t+1})$ through a similar process as above, which is formulated as $L(\widetilde{\mathbf{z}}_{t+1}, h_{t+1}) = L(\widetilde{\mathbf{z}}_t, h_t) - \eta \widetilde{\mathbf{g}}_t^2$. Finally, we can derive the update step size of $\widehat{\xi}_t(h)$:

$$\begin{aligned} \widehat{\xi}_{t+1}(h) - \widehat{\xi}_t(h) &= |L(\widetilde{\mathbf{z}}_{t+1}, h_{t+1}) - L(\mathbf{z}_{t+1}, h_{t+1})| - |L(\widetilde{\mathbf{z}}_t, h_t) - L(\mathbf{z}_t, h_t)| \\ &= L(\widetilde{\mathbf{z}}_{t+1}, h_{t+1}) - L(\widetilde{\mathbf{z}}_t, h_t) - (L(\mathbf{z}_{t+1}, h_{t+1}) - L(\mathbf{z}_t, h_t)) \\ &= -\eta \widetilde{\mathbf{g}}_t^2 - \left(-\frac{1}{\gamma_t} \eta \widetilde{\mathbf{g}}_t^2\right) = -(1 - \frac{1}{\gamma_t}) \eta \widetilde{\mathbf{g}}_t^2 < 0. \end{aligned} \quad (16)$$

The result indicates that layer-wise and channel-wise dropout can decrease $\widehat{\xi}_t(h)$ at every iteration, thus forcing it to be a small number at the end of training. This discussion is also experimentally verified in Section 5.5 of the paper. Therefore, layer-wise and channel-wise dropout can reduce the sensitivity of the model to dropout and produce a tight generalization error bound. ■

A.2 Proof of Proposition 2

Proposition 2. Dropout in different layers perform as different data augmentations in the input space, thus the layer-wise dropout can generate more diverse augmented data than the single-layer dropout, i.e., increasing the dataset size N .

Proof. Consider the hypothesis h that is a network with n hidden layers, and the i -th hidden layer is denoted as $l_i(\cdot)$. Given the input image x , we denote the output of the i -th hidden layer as $f_i(x)$, and use $\widetilde{f}_i(x)$ to denote its perturbed version generated by dropout. Let $m^{(i)}$ be the dropout mask in the i -th hidden layer, we can formulate the relationship between $\widetilde{f}_i(x)$ and $f_i(x)$ as:

$$\widetilde{f}_i(x) = m_i \odot (f_i(x)). \quad (17)$$

Assuming that for a given $\widetilde{f}_i(x)$, we can find a x'_i in the input space that satisfies $\widetilde{f}_i(x) = f_i(x'_i)$. Then we will prove that for an input x , it is hard to find a single $x'_i = x'_j$ if $i \neq j$, where $\widetilde{f}_i(x) = f_i(x'_i)$

and $\widetilde{f}_j(x) = f_j(x'_j)$. For convenience, we consider the situation when $j < i$. From Eq. (17), we have:

$$\begin{aligned}\widetilde{f}_i(x) &= (m_i \odot l_i) \dots (m_{j+1} \odot l_{j+1})(\widetilde{f}_j(x)) \\ &= (m_i \odot l_i) \dots (m_{j+1} \odot l_{j+1})(f_j(x'_j)), \\ f_i(x'_i) &= l_i \circ l_{i-1} \dots \circ l_{j+1}(f_j(x'_j)).\end{aligned}\tag{18}$$

If there exist $x'_i = x'_j$, from Eq. (18) we can derive that $\{m_i, m_{i-1}, \dots, m_{j+1}\}$ does not apply any modification to $f_j(x'_j)$, which means $m = 1$ for $m \in \{m_i, m_{i-1}, \dots, m_{j+1}\}$ when $f_j(x'_j) > 0$. This situation is unlikely to happen unless the dropout rate $p = 0$ for all hidden layers between l_i and l_j . Therefore, inserting dropout into different layers performs as different data augmentations in the input space, which can generate diverse variants of data during training. Based on this conclusion, our layer-wise and channel-wise dropout can effectively increase the dataset size N , leading to a tighter generalization error bound than single-layer dropout. ■

B ANALYSIS FOR THE STRONG BASELINE DEEPALL⁺⁺

B.1 Effect of each components in DeepAll⁺⁺

We conduct ablation studies for each component in DeepAll⁺⁺ on PACS with ResNet-18 as the backbone. The experimental results are represented in Tab. 12. The DeepAll⁺⁺ is a strong baseline model for domain generalization, which consists of two augmentation methods, *i.e.*, the image-level augmentation method RandAug [7] and the feature-level augmentation method SwapStyle. Specifically, SwapStyle can augment the data at the feature level by exchanging the feature statistics of any two images as the statistics can indicate the representative style information [16, 68]. RandAug [7] is an effective automated data augmentation method to increase image-level diversity. As shown in Tab. 12, both Swapstyle and RandAug can effectively improve the model performance. SwapStyle improves the average accuracy performance of DeepAll by 1.55% (83.19% vs. 81.64%) while RandAug improves by 1.99% (83.63% vs. 81.64%). As we can observe, Swapstyle can effectively improve the model performance on Cartoon and Photo, while RandAug is more effective on Art and Sketch. By combining DeepAll with SwapStyle and RandAug, DeepAll⁺⁺ achieves the highest performance and outperforms DeepAll by 3.28% (84.92% vs. 81.64%), indicating that the two augmentation methods are complementary with each other for improving model generalization. We also test the performance of PLACE dropout on these different variants of DeepAll⁺⁺ and present the experimental results in Tab. 12. We observe that the performance of PLACE dropout on Photo

Table 12. Effect (%) of PLACE dropout on different variants of DeepAll⁺⁺ on the PACS dataset with ResNet-18 as the backbone network. We denote SwapStyle as S and RandAug as R for simplicity. The best is **bolded**.

Method	Art	Cartoon	Sketch	Photo	Avg.
DeepAll	80.19	77.19	73.48	95.71	81.64
+ PLACE dropout	82.60	78.33	81.47	95.65	84.51
DeepAll w. S	81.59	78.03	76.43	96.71	83.19
+ PLACE dropout	82.86	79.43	82.18	95.51	85.00
DeepAll w. R	82.90	76.89	78.55	96.17	83.63
+ PLACE dropout	84.14	79.89	79.95	96.11	85.02
DeepAll ⁺⁺	83.64	78.28	81.41	96.34	84.92
+ PLACE dropout	85.40	79.69	83.97	96.23	86.32

Table 13. Comparison of PLACE dropout with other start-or-the-art DG methods on DeepAll⁺⁺ using ResNet-18 as the backbone network. The experiments are conducted on the PACS dataset. The best is **bolded**.

Method	Art	Cartoon	Sketch	Photo	Avg.
MMLD [24]	81.28	77.16	72.29	96.09	81.83
+DeepAll ⁺⁺	83.40	79.56	83.15	94.19	85.08
SagNet [28]	81.01	77.47	77.32	95.99	82.95
+DeepAll ⁺⁺	83.89	79.67	80.63	96.41	85.15
PLACE (Ours)	82.60	78.33	81.47	95.65	84.51
+DeepAll ⁺⁺	85.40	79.69	83.97	96.23	86.32

drops a little, which may be due to the task being saturated in performance. Nevertheless, PLACE dropout can achieve significant improvements on all variants, which verifies its effectiveness and stability on different baselines.

B.2 Incorporating DeepAll⁺⁺ with other DG methods

Besides incorporating the strong baseline DeepAll⁺⁺ with our PLACE dropout, we also investigate its efficacy in improving the SOTA DG methods, including MMLD [24] and SagNet [28]. MMLD is an adversarial domain-invariant learning method that trains the domain-generalized model by extracting features to confuse the domain discriminator. SagNet is a regularization method that decouples content and style features from the representations based on statistics and reduces the style bias of the model to enhance its generalization. As presented in Tab. 13, both the SOTA methods gain significant improvement on DeepAll⁺⁺, *i.e.*, MMLD gets 3.25% (85.08% vs. 81.83%) and SagNet gets 2.20% (85.15% vs. 82.95%). The results verify the effectiveness of our strong baseline. Moreover, among these methods, PLACE dropout achieves the best performance, which proves its superiority and effectiveness to mitigate the overfitting issue of the model on source domains.