

# Mixup-Induced Domain Extrapolation for Domain Generalization

Meng Cao, Songcan Chen\*

College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics  
MIIT Key Laboratory of Pattern Analysis and Machine Intelligence

## Abstract

Domain generalization aims to learn a well-performed classifier on multiple source domains for unseen target domains under domain shift. Domain-invariant representation (DIR) is an intuitive approach and has been of great concern. In practice, since the targets are variant and agnostic, only a few sources are not sufficient to reflect the entire domain population, leading to biased DIR. Derived from PAC-Bayes framework, we provide a novel generalization bound involving the number of domains sampled from the environment ( $N$ ) and the radius of the Wasserstein ball centred on the target ( $r$ ), which have rarely been considered before. Herein, we can obtain two natural and significant findings: when  $N$  increases, 1) the gap between the source and target sampling environments can be gradually mitigated; 2) the target can be better approximated within the Wasserstein ball. These findings prompt us to collect adequate domains against domain shift. For seeking convenience, we design a novel yet simple *Extrapolation Domain* strategy induced by the *Mixup* scheme, namely EDM. Through a reverse Mixup scheme to generate the extrapolated domains, combined with the interpolated domains, we expand the interpolation space spanned by the sources, providing more abundant domains to increase sampling intersections to shorten  $r$ . Moreover, EDM is easy to implement and be plugged-and-played. In experiments, EDM has been plugged into several methods in both closed and open set settings, achieving up to 5.73% improvement.

## Introduction

In conventional classification, the training set and test set generally follow independent identical distribution (i.i.d.) assumption. However, it is impractical in real-world applications due to domain shift (Li et al. 2022), including changing background, style, color, etc. To alleviate this issue, a learning paradigm, namely Domain Generalization (DG), has been presented and received increasing attention. DG aims to induce a well-performed (meta-)classifier from a set of given source domains so that it can generalize to unseen but related target domains.

Up to now, abundant methods have been proposed for DG (Wang et al. 2022). Domain-invariant representation (DIR) (Lu et al. 2022), as one of the dominant approaches, has been

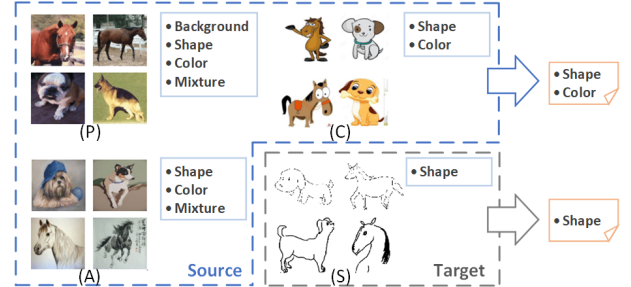


Figure 1: The illustration of biased DIR on PACS dataset, where (P) denotes Photo, (C) denotes Cartoon, (A) denotes Art, and (S) denotes Sketch. The representations provided from each domain are listed roughly in corresponding box.

widely studied, which can be divided into the following categories: causal inference (Arjovsky et al. 2019), information bottleneck (Li et al. 2022), adversarial learning (Ganin et al. 2016), and others (Ding et al. 2022). These approaches aim to remove the impurity representations, that is domain dependent representations, and to find common representations depending on the downstream task as much as possible. It seems reliable. In practice, since the target domains are variant and agnostic, it is difficult to extract unbiased DIR among domains from limited sample domains that are insufficient to reflect the entire domain population. Taking PACS dataset as a sample, shown in Fig. 1, we can observe that the required representations are inconsistent between the source and target domains. If the sketch domain is selected as the target domain, shape and color will be provided from the source domains, but only the shape is required. In this case, the color is redundant on the entire domain population, which may lead to poor generalization ability. What's more, (Zhu et al. 2022) has discovered that the domains cannot be mixed and can be obviously observed within each class. This phenomenon indicates that the impurity representations remain to some extent and the discrepancy cannot be completely eliminated, which further implies biased DIR with limited domain sampling. Herein, a question arises spontaneously: will the generalization ability be improved as domain sampling increases? Going one step further, what are the factors that influence the generalization ability?

\*Corresponding author: s.chen@nuaa.edu.cn

To answer this question, we first have to shift our perspective in DG. Almost previous works in DG, to the best of our knowledge, can boil down to a task-oriented framework, which pays attention to the divergence between pairwise domains (Lu et al. 2023), leading to inflexibility in theory. In essence, DG is an inductive learning paradigm on multiple related tasks, which follows a two-step sampling process, namely an environment-task framework (Baxter 2000). In analogy to the standard single-task learning where data is sampled from an unknown distribution, tasks in DG are sampled from an unknown task distribution, i.e., the environment. Herein, these tasks are more commonly referred to as domains in DG. In this way, the gap between environments reflects the similarity between the domain population learned from source domains and that of target domains. And, the gap between tasks reflects the relationship between observed tasks, such as the relationship between animal categorization in two environments. Therefore, compared to the former, the environment-task perspective can be more flexible and not limited to closed-set scenario, such as open-set scenario (Shu et al. 2021).

Following this perspective, we provide a novel generalization bound for DG, derived from PAC-Bayes framework (McAllester 1998), whose key is change of measure inequality. This bound involves the number of domains sampled from the environment  $N$  and the radius of the Wasserstein ball centred on the target, which have received little attention. Herein, we can obtain two natural and significant findings: when  $N$  increases, 1) the gap between the source and target sampling environments can be gradually reduced; 2) the target can be better approximated within the Wasserstein ball. These findings prompt us to collect adequate domains against domain shift. Indeed, previous works have made efforts to generate varied domain samples through data augmentation. (Shankar et al. 2018) develop an adversarial strategy by reversing the gradient of the domain classifier. (Xiao et al. 2021) generate a new domain through an extra network module. Obviously, these methods are inflexible and computationally expensive. Mixup (Zhang et al. 2017) is another popular and widely used technique. However, the generated samples are usually mixed from pairwise instances and lie in the interpolation space spanned by those instances.

In this manuscript, in search of modeling convenience, we design a novel yet simple *Extrapolation Domain* strategy induced by the *Mixup* scheme, namely EDM. It is a two-stage Dir-Mixup (Shu et al. 2021) strategy, i.e., extrapolation followed by interpolation. In extrapolation stage, the extrapolated domains are generated through a reverse Mixup scheme. And then, in interpolation stage, the new domains are generated by mixing the generated extrapolated domains. Along this line, the interpolation space spanned by the sources can be expanded, so that the domains can be obtained not only inside but also outside this space, called interpolation domains and extrapolation domains, respectively. In this way, more abundant domains can be provided with unrestricted of the interpolation space, and then the intersections of sampled-domain sets are increased to provide a better target approximation. Moreover, EDM inherits the lightweight and flexible characteristics of Mixup, so that

it can be easy to implement and be directly plugged-and-played. To sum up, our contributions are listed as follows:

1. A novel generalization bound for DG is provided, which is flexible for task settings and guides us to pay attention to the number of domains from the environment sampling perspective and the radius of the Wasserstein ball.
2. A two-stage Dir-Mixup strategy, namely EDM, is initially designed to provide more abundant domains, where the extrapolation domains outside the interpolation space can increase the sampling intersections.
3. In experiments, EDM has been plugged into several methods in both closed and open set settings, achieving up to 5.73% improvement.

## Related Works

### Mainstream Methods for Domain Generalization

DIR aims to find task-dependent but domain-independent representations. Classic moment-based methods align the statistics in representation space, such as MMD (Grubinger et al. 2015), CORAL (Sun and Saenko 2016). The adversarial-based methods, e.g., DANN (Ganin et al. 2016), make an attempt to confuse the domain classifier to remove the domain-related representations, so that the task-related representations can be retained. Unlike DANN, which adds an extra network module, information bottleneck (Li et al. 2022) leverages the information entropy between input and hidden representations and the one between hidden representations and output. To avoid retaining spurious correlations, causal inference has been introduced, where IRM (Albuquerque et al. 2019), a strategy with gradient penalty, is one of the well-known methods. And, VRex (Krueger et al. 2021) provides a variance penalty regularization in loss to obtain invariant representations. Meanwhile, (Sagawa et al. 2020) argue that due to minimizing average loss via empirical risk minimization, spurious correlations arise from typical examples, so that they regroup domains with underlying correlation representations, e.g., background, to avoid learning models that rely on spurious correlations. RSC (Huang et al. 2020) iteratively forces a CNN to activate features that are less dominant in the training domain, but still correlated with labels. In essence, almost of them, especially (Ding et al. 2022), make the effort to remove impurities across domains, which are domain-specific features. In recent years, some researches imply that only a few source domains are insufficient to reflect the entire domain population. For example, (Zhu et al. 2022) discover that at local regions, the domains cannot be mixed and are clustered.

Data augmentation, which is another hot topic in DG, is one of the cheap and simple ways to increase the quality and diversity of the training data. Mixup (Zhang et al. 2017) is a popular technique to achieve this goal. In (Lu et al. 2023), the samples mixed by the same class but different domains are generated to enlarge the diversity, and the samples mixed by the same domain but different classes are generated to reduce the influence of redundant domain information. Meanwhile, (Zhou et al. 2021) generate new samples by mixing the statistical information of paired samples in multiple hid-

den layers. (Mancini et al. 2020) advocate that mixing samples of both different domains and classes allows to obtain samples that cannot be categorized in a single class and domain of the one available during training, so that they construct some novel semantic-visual samples from triple tuple samples to recognize unseen categories in unseen domains. To further enlarge the diversity of the generated sample, (Shu et al. 2021) proposes a multi-sample mixing strategy, namely Dir-Mixup, whose mixing coefficients sample from Dirichlet distribution. Although the Mixup scheme is flexible, all of them are limited by interpolation space, and the domain distribution cannot be reflected due to the mixing statistics from a single sample. In contrast, a min-max game has been designed to generate new examples, such as Cross-Grad (Shankar et al. 2018), which utilizes gradient ascent to expand both class and domain space. (Xiao et al. 2021) add a network module to sample a new domain with a meta-learning learner. Nevertheless, high computational complexity cannot be ignored and its flexibility is not well.

### Generalization Bounds for Domain Generalization

In recent years, some generalization bounds for DG have been emerged to demonstrate the effectiveness of corresponding methods. (Albuquerque et al. 2019) provide a generalization bound w.r.t. the specific linear combination of empirical errors from the source domains and the divergence between the real target distribution and the approximate fake target distribution within the space spanned by the source domains. In analogy, (Dai et al. 2023) provide a similar formulation, replacing  $\mathcal{H}$ -divergence with Wasserstein distribution. And then, (Lu et al. 2023) explicitly provide an additional term, maximizing the divergence across the sources domains, to further reduce loss caused by alignment. These bounds focus on the distribution divergence and motivate DG methods based on DIR. Besides, based on kernel mean embedding, the bound w.r.t. the marginal distribution is given in (Blanchard, Lee, and Scott 2011), which implies that the generalization ability is related to the number of sampled domains. Obviously, this bound is valid for the domain shift depending on the marginal distributions and is difficult to explain methods in deep network.

### Methodology

In this section, we introduce our motivation, theoretical framework, and proposed method EDM in detail.

#### Preliminaries

In DG, a common setting is to provide  $N$  domains under domain shift. Let  $D_n$  denotes  $n$ -th observed domain, which is a set of  $M_n$  independent samples from a space of examples  $\mathcal{Z}$ , i.e.,  $D_n = \{z_n^m\}_{m=1}^{M_n}$ . Each sample is drawn from an unknown distribution  $\mathcal{D}_n$ , namely  $z_n^m \sim \mathcal{D}_n$ , and  $z_n^m = (x_n^m, y_n^m)$ , where  $x_n^m$  denotes an input instance and  $y_n^m$  denotes the corresponding label. Due to domain shift,  $\mathcal{D}_{ni} \neq \mathcal{D}_{nj}, \forall ni \neq nj$ . According to the environment-task perspective discussed in Introduction, we argue that these domains are generated i.i.d. from an unknown hyper-distribution  $\tau$ , i.e., distribution over distribution or  $\mathcal{D}_n \sim \tau$ .

Let  $h \in \mathcal{H}$  denotes a hypothesis  $h$  belongs to a hypothesis space  $\mathcal{H}$ . In analogy to the standard single-task learning where a single hypothesis  $h$  is learned based on an observed sampling set  $D$ , the selected hypothesis  $h$  is induced from observed sampling sets (domains)  $\{D_n\}_{n=1}^N$ . To select an appropriate hypothesis, the PAC-Bayes framework, whose starting point is model average, construct a probability distribution set over  $\mathcal{H}$ , namely  $\mathcal{M}(\mathcal{H})$ . That is,  $h \sim P$ , where  $P \in \mathcal{M}(\mathcal{H})$  denotes a probability measure in  $\mathcal{M}(\mathcal{H})$ , and is described as the prior, which is data-dependent. Based on the observed sampling set  $D$  and the prior  $P$ , the learner output a posterior distribution  $Q$  over  $\mathcal{H}$  when learning a new task. Herein, the prior and posterior notations are utilized to describe the relationship between the observed task and the new task, without the need for a likelihood function to connect them. Following the environment-task framework, the above standard PAC-Bayes framework should be extended to adapt to changes in the environment. To this end, we assume a hyper-distribution over the distribution measure space  $\mathcal{M}(\mathcal{H})$ , e.g.,  $\mathcal{P} \in \mathcal{M}(\mathcal{P})$ , where  $\mathcal{P}$  denotes a hyper-prior distribution.  $\mathcal{Q}$  is similar and denotes a hyper-posterior distribution. The expected error is denoted as  $er(Q, \tau) \triangleq \mathbb{E}_{Q \sim \mathcal{Q}} er(Q, \tau)$ . Since  $er(Q, \tau)$  is not computable, we can evaluate its corresponding empirical error  $\hat{er}(Q, \tau) \triangleq \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{P \sim \mathcal{Q}} \hat{er}(P, D_n)$ .

### Generalization Bound

**Theorem 1** (Domain Generalization Generalization bound). *Giving a hypothesis space  $\mathcal{H}$ , and  $N$  domains  $\{D_n\}_{n=1}^N$  sampled from  $\tau$ , where each domain  $D_n$  consists of  $M_n$  samples. Let  $\mathcal{P}$  denotes a hyper-prior distribution  $\mathcal{P} \in \mathcal{M}(\mathcal{P})$ , where  $P \in \mathcal{M}(\mathcal{H})$  and  $\mathcal{M}(\mathcal{S})$  denotes the set of all probability over  $\mathcal{S}$ . Then, for any  $\delta \in (0, 1]$ , the following inequality holds uniformly for all hyper-posterior distributions  $\mathcal{Q}$  with probability at least  $1 - \delta$ :*

$$\begin{aligned} \mathbb{E}_{Q \sim \mathcal{Q}} er(Q, \tau) &\leq \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{P \sim \mathcal{Q}} \hat{er}(P, D_n) \\ &\quad + \sqrt{\frac{L_0 \cdot W_1(\mathcal{Q}, \mathcal{P}) + \ln(N/\delta)}{2(N-1)}} \\ &\quad + \frac{1}{N} \sum_{n=1}^N \sqrt{\frac{L_n \cdot W_1(Q, P_n) + \ln(NM_n/\delta)}{2(M_n-1)}} \end{aligned} \quad (1)$$

where  $W_1(\cdot, \cdot)$  is the 1st order Wasserstein Distance, and  $L_0$  and  $L_n$  are the Lipschitz constants. Its corresponding proof is provided in Appendix in detail.

From Theorem 1, we can find that the generalization error of DG is bounded by the empirical error from the source domains plus two complexity terms. The first complexity term is a so-called environment-complexity term, which measures the gap between environments  $W_1(Q, \mathcal{P})$ , where environments sample the source and target domains, respectively. This gap is caused by observing only a finite number of tasks. Meanwhile, the second complexity term is an average task complexity term, which measures the divergence of tasks between target domain and each source domain  $W_1(Q, P_n)$ . For a clearer explanation, due to the domains following a hyper-distribution, we can further assume

that the divergence between environments  $W_1(Q, P)$  can be relaxed to the radius of the Wasserstein ball centred on the target  $r$ , that is  $W_1(Q, P) \leq r$ . To this end, upper bound of generalization error is with respect to  $N$ ,  $M_n$ ,  $r$  and  $W_1(Q, P_n)$ . When each  $M_n$  approaches to infinity, the task complexity term will converge to zero, and when  $N$  approaches to infinity, both the task and environment complexity term will converge to zero. These findings are very natural and intuitive, prompting us to collect adequate observed samples for each domain while collecting adequate augmented domains. In this way, the domain population learned from the observed source domains can be a better approximation of the entire population, ensuring good performance on a novel task, i.e., the task of target domain. Moreover, according to Theorem 3.4 in (Mohajerin Esfahani and Kuhn 2018), the radius  $r$  will be inversely proportional to the number of observed domains  $N$ , if the unknown distribution is light-tailed in the sense. Therefore, through increasing sampling domains, the target can be better represented on domain population, analogous to the situation with test examples on the training distribution. In summary, increasing the sampling domains is another way to improve the generalization ability of DG, similar to increasing the sampling samples from each domain. By the way, due to loose assumptions for tasks, our generalization bound is very flexible and can explain multiple settings through the gap of tasks  $W_1(Q, P_n)$ , such as open-set setting.

### Extrapolation Domain induced by Mixup

We know that it's unrealistic to collect an infinite number of domains. As an alternative, we design a novel yet simple *Extrapolation Domain* strategy induced by the *Mixup* scheme, namely EDM. Compared to previous methods (Zhou et al. 2021; Shu et al. 2021; Lu et al. 2023), most of which mix paired samples, EDM has two significant characteristics: 1) mixing the statistics from multiple source domains; 2) constructing an extrapolation space surrounding the interpolation space spanned by source domains.

The reason behind the former is intuitive, that is, our aim is to augment domains rather than samples. Moreover, the mixing strategy in classic methods boils down to a linear interpolation strategy, which only generates new domains between two domains (the lines between vertices as shown in Fig. 2). This pairwise mixing strategy is obviously limited by the lack of domains mixed from multiple domains, i.e., the whole blue area.

The reason behind the latter is that due to the finite observed domains and the variant and agnostic target domain, as shown in Fig. 2, the interpolation space or the environment obtained from the source domains i.e., the blue area, may be biased. To mitigate this issue, inspired by  $W_1(Q, P) \leq r$ , where  $r \propto \exp(1/N)$  if  $P_n$  follows a light-tailed sampling, we would like to expand the interpolation space to increase the intersections with the domains sampling from the target environment, i.e. the green area, in order to satisfy the sampling assumption as much as possible. In this way, not only more abundant domains can be generated further, but also theoretical generalization ability can be guaranteed to some extent.

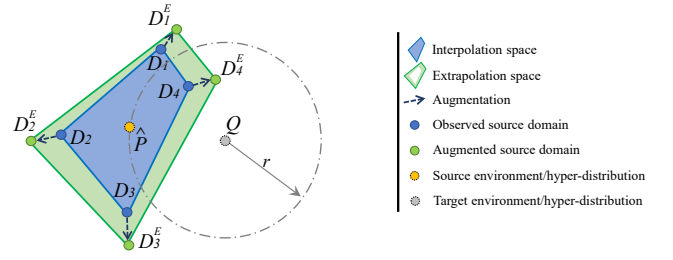


Figure 2: The illustration of EDM. Through Mixup scheme, each domain  $D_n$  is pushed outward to generate the extrapolation space, represented by the light green region, based on the corresponding new domain  $D_n^E$ .

To realize EDM, similar to Dir-Mixup (Shu et al. 2021), we formulate the foundation of EDM, i.e., the multiple domain mixing scheme, as follows:

$$\mathcal{D}_\lambda = \sum_{n=1}^N \lambda_n \mathcal{D}_n, \forall \lambda_n \geq 0 \text{ and } \sum_{n=1}^N \lambda_n = 1 \quad (2)$$

where  $\lambda_n$  denotes the mixing coefficient of  $n$ -th domain. Unlike classic Mixup scheme, where the mixing coefficient is sampled from Beta Distribution,  $\lambda$  are sampled from Dirichlet Distribution parameterized by a parameter  $\alpha$ , i.e.,  $\lambda \sim \text{Dirichlet}(\alpha)$ .

For a general case, we assume that  $\mathcal{D}_n \triangleq \mathcal{N}(\mu_n, \sigma_n^2)$ , where  $\mathcal{N}(\mu, \sigma^2)$  denotes a Gaussian Distribution with the mean  $\mu$  and the standard deviation  $\sigma$ . Therefore, Eq. (2) can be reformulated as:

$$\mathcal{N}(\mu_\lambda, \sigma_\lambda^2) = \mathcal{N}\left(\sum_{n=1}^N \lambda_n \mu_n, \sum_{n=1}^N \lambda_n^2 \sigma_n^2\right) \quad (3)$$

And, each pair of parameters  $(\mu_n, \sigma_n^2)$  can be calculated based on the corresponding training data  $x_n^m \in \mathbb{R}^{C \times H \times W}$  in each batch, formulated as:

$$\mu_n = \frac{1}{B_n H W} \sum_{m=1}^{B_n} \sum_{h=1}^H \sum_{w=1}^W (x_n^m)_{h,w} \\ \sigma_n^2 = \frac{1}{B_n H W} \sum_{m=1}^{B_n} \sum_{h=1}^H \sum_{w=1}^W \left( (x_n^m)_{h,w} - \mu_n \right)^2 \quad (4)$$

where  $B_n$  is the number of training data on  $n$ -th domain.

To avoid wasting information, a momentum strategy is adopted, which utilizes historical information through a moving average weight  $\rho$ . Then, we have

$$\mu_n^t = \rho \mu_n^{t-1} + (1 - \rho) \mu_n, \sigma_n^t = \rho \sigma_n^{t-1} + (1 - \rho) \sigma_n \quad (5)$$

where  $\mu_n^{t-1}$  denotes the mean of  $n$ -th domain at  $(t-1)$ -th iteration, and  $\sigma_n^{t-1}$  is similar.

Combining Eq. (3) and Eq. (5), we can obtain the statistics of a new domain. It is noted that these new domains only lie in the interpolation space spanned by observed source domains, which is not our intention.

To generate extrapolated domains, we reverse the above multiple domain mixing scheme to obtain the corresponding supported domains for the extrapolation space. For example,

$D_1^E$  is a supported domain corresponding to  $D_1$  in Fig. 2. Specifically, each domain  $D_n$  can be regarded as an interpolated domain through mixing the corresponding extrapolated domain  $D_n^E$  and the other source domains, respectively. In other words, similar to Eq. (2), we have:

$$\mathcal{D}_n = \lambda_n \mathcal{D}_n^E + \sum_{i=1, i \neq n}^N \lambda_i \mathcal{D}_i, \forall \lambda_i \geq 0 \text{ and } \sum_{i=1}^N \lambda_i = 1 \quad (6)$$

In this way,  $\mathcal{D}_n^E$  can be reformulated as:

$$\mathcal{D}_n^E = \frac{1}{\lambda_n} \mathcal{D}_n - \sum_{i=1, i \neq n}^N \frac{\lambda_i}{\lambda_n} \mathcal{D}_i \quad (7)$$

where  $\lambda_n$  has the same constraint as Eq. (6). And, its Gaussian Distribution can be referred to as

$$\begin{aligned} & \mathcal{N}\left((\mu_{\lambda}^E)_n, ((\sigma_{\lambda}^E)_n)^2\right) \\ &= \mathcal{N}\left(\frac{1}{\lambda_n} \mu_n - \sum_{i=1, i \neq n}^N \frac{\lambda_i}{\lambda_n} \mu_i, \frac{1}{\lambda_n^2} \sigma_n^2 + \sum_{i=1, i \neq n}^N \left(\frac{\lambda_i}{\lambda_n}\right)^2 \sigma_i^2\right) \end{aligned} \quad (8)$$

Next, based on these supported domains outside the interpolation space, we once again employ the multiple domain mixing scheme, similar to Eq. (3), to generate the new domains, which will be located inside and outside the interpolation space spanned by the observed source domains  $\{D_n\}_{n=1}^N$ , as the green and blue areas in Fig. 2.

Note that this twice Mixup scheme, i.e., Eq. (8) followed by Eq. (3) employed with the corresponding supported domains, is one of the augmentation schemes to indirectly obtain the extrapolation domains. We can also select more domains with the positive coefficient in Eq. (7) to directly obtain the extrapolation domains. Obviously, it is too complex and difficult to control.

Finally, through AdaIN scheme (Zhou et al. 2021), the samples sampling from new domains can be represented as:

$$x_a^m = \frac{x_n^m - \mu_n}{\sigma_n} \sigma_a + \mu_a \quad (9)$$

where  $\mu_a$  and  $\sigma_a$  denote the mean and standard deviation of augmented domains through twice Mixup scheme, respectively. These generated samples are directly fed to the training model along with the original samples. And, their class labels are the same as the corresponding original samples, and a new domain label will be assigned. Detailed Algorithm is provided in Appendix.

## Experiments

In this section, extensive experiments are constructed to comprehensively evaluate the effectiveness of EDM on two datasets both in closed and open set settings.

### Datasets and Settings

For the architecture, we use ResNet-18 as backbone on three datasets, i.e., PACS, Office-Home, and DomainNet datasets. For both settings, we follow corresponding settings from the previous methods, i.e., the same closed-set setting as (Lu

et al. 2022), and the same open-set setting as (Shu et al. 2021). In closed-set setting, we compare with twelve recent strong comparison methods and two other representative methods. Except ERM, they can be divided into four categories: 1) domain-invariant representation based methods DANN (Ganin et al. 2016), MMD (Grubinger et al. 2015), CORAL (Sun and Saenko 2016), VREx (Krueger et al. 2021), DIFEX (Lu et al. 2022); 2) data augmentation based method Mixup (Zhang et al. 2017), CrossGrad (Shankar et al. 2018), MixStyle (Zhou et al. 2021); 3) learning robust features based methods: GroupDRO (Sagawa et al. 2020), RSC (Huang et al. 2020); 4) model optimization based method ANDMask (Parascandolo et al. 2020), SAGM (Wang et al. 2023). And in open-set setting, following (Shu et al. 2021), we compare with seven other popular methods, which are divided into the following categories: 1) data augmentation based method CuMix (Mancini et al. 2020); 2) learning robust features based methods PAR (Wang et al. 2019), RSC (Huang et al. 2020); 3) heterogeneous method FC (Li et al. 2019b); 4) meta-learning based methods MLDG (Li et al. 2018), Epi-FCR (Li et al. 2019a), DAML (Shu et al. 2021). For more details on datasets, comparison methods, and settings, please refer to Appendix. The code is available at <https://github.com/Alrash/EDM>.

## Results and Analysis

Tab. 1 reports accuracy results in closed-set setting, and Tab. 2 reports accuracy and H-score (Fu et al. 2020) results in open-set setting.

From **Tab. 1**, we can observe two key findings as follows: 1) compared with previous methods, the learner + EDM can give the best results, achieving up to 5.73% improvement in the sketch on PACS; 2) cross-compared with similar methods with high complexity, EDM has shown its lightweight and flexible characteristics.

Specifically, the former key finding can be reflected in the following aspects. First, on both datasets, the learner attaching EDM can achieve the best results on average. Second, in each domain on both datasets, the best and second results can almost be obtained through attaching EDM. These phenomena can testify to the effectiveness of EDM. Third, the learner + Inter, i.e., attaching augmented interpolation domains, is usually slightly weaker than the learner + EDM, i.e., attaching both augmented interpolation and extrapolation domains, but is better than the corresponding basic learner. This fact indicates that domain augmentation is beneficial and can improve the generalization ability for DG. And, additional extrapolation space, which reflects the more complete domain population combined with interpolation space, can further improve the performance. Fourth, different basic learners have received different gains in each domain. The results in the sketch on PACS and in the clipart on Office-Home have received significant improvement, especially for DANN + EDM, which achieves up to 5.73% and 2.09%, respectively. And, the second improvements are in cartoon and product, respectively. These phenomena further indicate that EDM can simulate more severe drift and make the learner perform well in scenarios with significant domain shift. Fifth, compared with SAGM, there is no sig-

	PACS					Office-Home					DomainNet
	Art-Painting	Cartoon	Photo	Sketch	Avg	Art	Clipart	Product	Real-World	Avg	Avg
ERM	81.10	77.94	95.03	76.94	82.75	57.77	50.63	71.30	74.45	63.54	40.10
DANN	82.86	78.33	96.11	76.99	83.57	57.60	48.52	71.16	72.99	62.57	40.20
Mixup	81.84	75.43	95.27	76.51	82.26	58.71	51.00	72.20	75.42	64.33	39.24
RSC	82.13	77.99	94.43	79.87	83.60	57.07	50.77	71.93	73.63	63.35	37.13
MMD	80.32	76.45	92.46	<b>83.63</b>	83.21	59.29	50.52	72.34	74.43	64.15	39.14
CORAL	79.39	77.90	91.98	82.03	82.83	59.29	50.15	72.25	74.20	63.97	39.16
GroupDRO	79.15	76.75	91.32	81.52	82.19	59.09	50.22	71.91	74.48	63.92	31.78
CrossGrad $\ddagger$	80.37	74.87	<u>96.59</u>	74.98	81.70	58.67	51.18	71.66	74.80	64.08	39.49
Mixstyle $\ddagger$	82.51	79.09	95.65	79.23	84.12	55.50	51.00	70.62	73.19	62.57	39.98
ANDMask	80.81	73.29	95.81	71.95	80.47	53.61	47.54	69.36	72.23	60.69	23.92
VREx	81.54	78.11	95.39	80.35	83.85	59.09	49.81	71.64	74.82	63.84	37.96
DIFEX-ori $\ddagger$	82.86	78.46	94.97	79.41	83.93	57.89	50.82	71.61	73.40	63.43	38.33
DIFEX-norm $\ddagger$	<u>83.40</u>	<u>79.74</u>	95.03	79.10	84.32	58.09	51.50	72.08	73.62	63.82	38.53
SAGM $\ddagger$	82.62	78.50	96.05	79.64	84.20	59.13	51.23	72.67	75.90	64.73	38.86
<hr/>											
ERM + Inter	83.25	77.60	95.99	81.09	84.48	58.01	50.65	72.02	74.62	63.83	-
DANN + Inter	81.93	77.82	95.99	81.57	84.33	57.27	50.13	71.53	74.04	63.24	-
Mixup + Inter	83.01	76.32	<b>96.65</b>	78.04	83.50	<b>59.46</b>	52.30	72.88	75.60	<u>65.06</u>	-
SAGM + Inter	82.28	79.01	96.29	80.22	84.45	58.55	<b>52.71</b>	72.85	75.51	64.91	-
<hr/>											
ERM + EDM	82.32	79.27	96.53	81.24	84.84	58.67	51.84	72.38	75.35	64.56	<u>40.34</u>
DANN + EDM	82.96	78.07	96.47	<u>82.72</u>	85.06	58.51	50.61	72.22	74.59	63.98	<b>40.40</b>
Mixup + EDM	<b>83.50</b>	79.14	<u>96.59</u>	81.04	<u>85.07</u>	<u>59.33</u>	51.94	<b>73.15</b>	<u>75.97</u>	<b>65.10</b>	40.04
SAGM + EDM	82.47	<b>80.38</b>	<u>96.59</u>	80.86	<b>85.08</b>	58.84	<u>52.33</u>	<u>72.94</u>	<b>76.06</b>	65.04	39.23

Table 1: Accuracy results on PACS, Office-Home and DomainNet datasets in closed-set settings. + Inter denotes that the method attaches augmented interpolation domains. + EDM denotes that the method attaches augmented interpolation and extrapolation domains. The **bold** and underline items are the best and the second-best results, respectively.  $\ddagger$  denotes our reproduced results on PACS and Office-Home, and  $\ddagger$  denotes our reproduced results on Office-Home. All results on DomainNet are reproduced.

	PACS		Office-Home	
	Acc	H-score	Acc	H-score
ERM	55.17	44.78	50.43	47.41
MLDG	57.43	45.00	51.07	47.58
FC	58.13	46.69	51.03	48.02
Epi-FCR	60.64	48.47	50.25	48.48
PAR	56.56	44.95	51.26	49.03
RSC	58.92	45.05	49.56	47.89
CuMix	57.85	41.05	51.67	49.40
DAML	65.49	<u>51.88</u>	56.45	53.34
<hr/>				
DAML + Inter	69.22	51.83	<u>59.15</u>	<u>53.64</u>
DAML + EDM	<b>70.78</b>	<b>54.12</b>	<b>59.58</b>	<b>54.19</b>

Table 2: Accuracy and H-score results both on PACS and Office-Home datasets in open-set settings.

nificant improvement with SAGM + Inter or SAGM + EDM. We think that the model perturbation mechanism can indirectly simulate the drift between domains so that the effect of EDM has been counteracted to some extent.

And, the latter key finding can be reflected in three aspects. First, Mixup + EDM can achieve better performance than Mixup, the results of ERM + Inter and ERM + EDM are almost better than Mixup while the gains on ERM and Mixup are approximately similar. These facts imply that domain augmentation and data augmentation can be paral-

lel to each other to improve the generalization ability, and theoretical analysis can be empirically testified simultaneously. Second, compared with DIFEX-ori and DIFEX-norm, as a representative of domain-invariant representation based methods, DANN + EDM can be slightly better. Third, compared with CrossGrad, which contains data augmentation and domain augmentation, the results of ERM + EDM are almost better, not to mention those of Mixup + EDM. These phenomena demonstrate the convenience of EDM, which does not require extra networks or tasks (regularizers).

From **Tab. 2**, we can observe the following findings. First, although DAML is the SOTA in open-set setting, in which augmented samples are utilized, its performance can also be improved when augmented domains are attached. Second, DAML + EDM achieves the best results. This fact indicates that extrapolation domains should be considered and can further improve performance. Third, significant improvement in accuracy results both of DAML + Inter and DAML + EDM can be observed, but their improvement in H-score results is a little bit inferior. The reason is that our proposed domain augmentation strategy does not directly solve the issue of unknown class detection in open-set setting. In fact, new domains containing random classes are further generated through Inter or EDM. Therefore, the learner can capture more complete known category information to boost performance through CE loss of mixup samples in DAML. For more detailed experimental results and corresponding analysis of Tab. 2, please refer to Appendix.



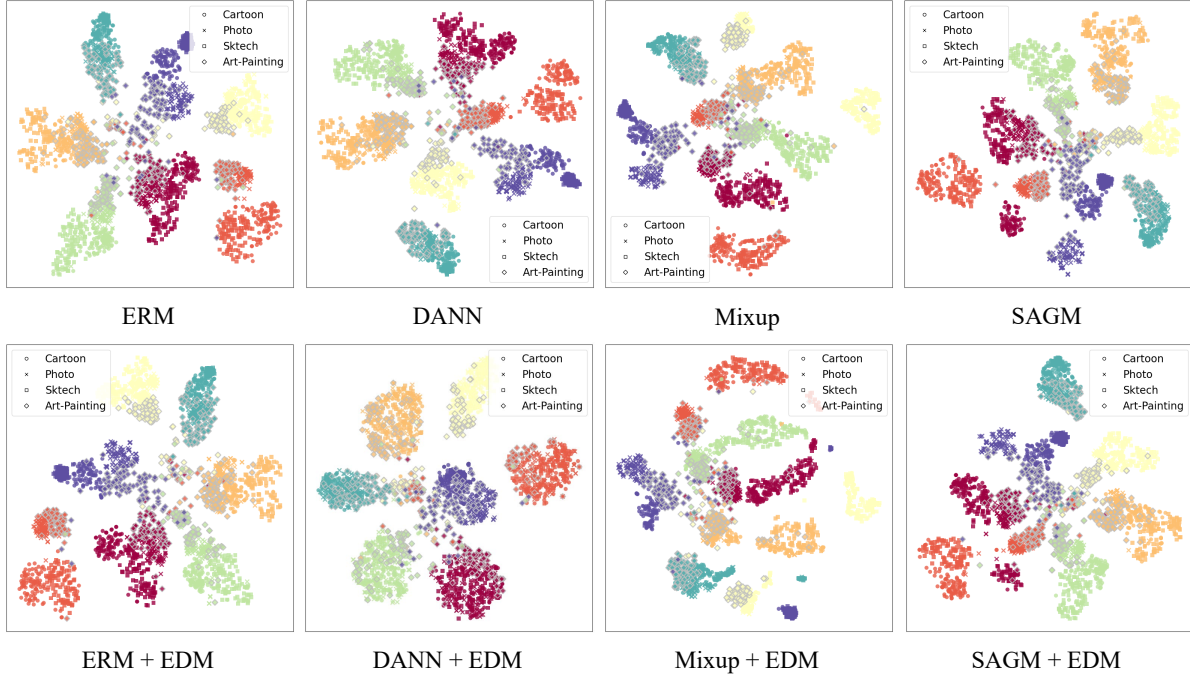


Figure 3: Visualization of the t-SNE embeddings of learned representation spaces for PACS with different methods. Different colors correspond to different classes and different shapes correspond to different domains. Note that the diamond with the grey edge denote the target domain.

In our experiments, **ablation study** is equivalent to whether adding new domains can improve the performance of the corresponding method. These results can be found both in Tab. 1 and Tab. 2, so we will no longer report them. Overall, adding new domains can improve performance, and combining interpolation and extrapolation domains can achieve better performance. More details can be referred to in the aforementioned analysis.

From **Fig. 3**, we can discover that the methods belonging to different categories exhibit different phenomena. DANN, as a domain-invariant representation based method, aims to obtain a representation space in which the domains can be confused in each class. However, the domains can be obviously observed and each class cluster seems not tight. In contrast, in DANN + EDM, the domains can be more scattered within each class and each class cluster can be tighter. These phenomena reveal that a few domains only receive biased domain-invariant representations, which contain undesired domain-dependent yet task-dependent representations. Domain augmentation, which mixes domain information rather than class information, can alleviate this issue to some extent. ERM and Mixup, as representative of aggregation methods, have not paid too much attention to domain information. Therefore, the domains exhibit a certain degree of clustering within each class. With EDM, the discriminative ability of tasks has not been harmed, and the classes in target domain prefer to classify the classes on a more similar source domain, such as the photo. For example, compared to ERM, the orange and the blue in ERM + EDM are closed to the corresponding classes in the photo,

and the domains are more scattered within each class. These phenomena also imply that domain augmentation can rich domain information, and obtain more essential representations even on non-domain-invariant representation based methods. Finally, compared to SAGM, a model optimization based method, EDM as a data perturbation strategy does not compromise the model perturbation mechanism. The tightness of each class has ups and downs on both sides.

## Conclusion

Domain generalization (DG) is regarded as an inductive learning paradigm on multiple related tasks, which belongs to an environment-task framework. Following this perspective, we give a novel generalization bound for DG, derived from PAC-Bayes framework. In light of this bound, we argue that the factors that influence the generalization ability involve four aspects. In this manuscript, we focus on two factors: the number of observed domains and the gap between sampling environments, which have received little attention in previous methods. After relaxing this gap to the radius of a Wasserstein ball centred on the target, we discover once again that increasing the sampling domains can improve the generalization ability. To this end, we design a novel yet simple Extrapolation Domain strategy induced by the Mixup scheme, namely EDM, which indirectly constructs an extrapolation space surrounding the interpolation space spanned by source domains to provide more abundant domains. In addition, EDM is easy to implement and can be plugged and played. Finally, extensive experiments are conducted to testify to the effectiveness of EDM.

## Acknowledgments

This work was supported by the Natural Science Foundation of China (NSFC) (Grant No.62076124).

## References

- Albuquerque, I.; Monteiro, J.; Falk, T. H.; and Mitliagkas, I. 2019. Generalizing to unseen domains via distribution matching. *arXiv preprint arXiv:1911.00804*, 8.
- Amit, R.; Epstein, B.; Moran, S.; and Meir, R. 2022. Integral Probability Metrics PAC-Bayes Bounds. In *Advances in Neural Information Processing Systems*, volume 35, 3123–3136.
- Arjovsky, M.; Bottou, L.; Gulrajani, I.; and Lopez-Paz, D. 2019. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*.
- Baxter, J. 2000. A model of inductive bias learning. *Journal of artificial intelligence research*, 12: 149–198.
- Blanchard, G.; Lee, G.; and Scott, C. 2011. Generalizing from several related classification tasks to a new unlabeled sample. In *Advances in neural information processing systems*, volume 24.
- Dai, R.; Zhang, Y.; Fang, Z.; Han, B.; and Tian, X. 2023. Moderately Distributional Exploration for Domain Generalization. In *International Conference on Learning Representations*.
- Ding, Y.; Wang, L.; Liang, B.; Liang, S.; Wang, Y.; and Chen, F. 2022. Domain generalization by learning and removing domain-specific features. In *Advances in Neural Information Processing Systems*, volume 35, 24226–24239.
- Fu, B.; Cao, Z.; Long, M.; and Wang, J. 2020. Learning to detect open classes for universal domain adaptation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16*, 567–583. Springer.
- Ganin, Y.; Ustinova, E.; Ajakan, H.; Germain, P.; Larochelle, H.; Laviolette, F.; Marchand, M.; and Lempitsky, V. 2016. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1): 2096–2030.
- Grubinger, T.; Birlutiu, A.; Schöner, H.; Natschläger, T.; and Heskes, T. 2015. Domain generalization based on transfer component analysis. In *Advances in Computational Intelligence: 13th International Work-Conference on Artificial Neural Networks, IWANN 2015, Palma de Mallorca, Spain, June 10–12, 2015. Proceedings, Part I 13*, 325–334. Springer.
- Huang, Z.; Wang, H.; Xing, E. P.; and Huang, D. 2020. Self-challenging improves cross-domain generalization. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, 124–140. Springer.
- Krueger, D.; Caballero, E.; Jacobsen, J.-H.; Zhang, A.; Binias, J.; Zhang, D.; Le Priol, R.; and Courville, A. 2021. Out-of-distribution generalization via risk extrapolation (rex). In *International Conference on Machine Learning*, 5815–5826. PMLR.
- Li, B.; Shen, Y.; Wang, Y.; Zhu, W.; Li, D.; Keutzer, K.; and Zhao, H. 2022. Invariant information bottleneck for domain generalization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 7399–7407.
- Li, D.; Yang, Y.; Song, Y.-Z.; and Hospedales, T. 2018. Learning to generalize: Meta-learning for domain generalization. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Li, D.; Yang, Y.; Song, Y.-Z.; and Hospedales, T. M. 2017. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, 5542–5550.
- Li, D.; Zhang, J.; Yang, Y.; Liu, C.; Song, Y.-Z.; and Hospedales, T. M. 2019a. Episodic training for domain generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 1446–1455.
- Li, Y.; Yang, Y.; Zhou, W.; and Hospedales, T. 2019b. Feature-critic networks for heterogeneous domain generalization. In *International Conference on Machine Learning*, 3915–3924. PMLR.
- Lu, W.; Wang, J.; Li, H.; Chen, Y.; and Xie, X. 2022. Domain-invariant Feature Exploration for Domain Generalization. *Transactions on Machine Learning Research*.
- Lu, W.; Wang, J.; Wang, Y.; and Xie, X. 2023. Towards Optimization and Model Selection for Domain Generalization: A Mixup-guided Solution. In *The KDD’23 Workshop on Causal Discovery, Prediction and Decision*, 75–97. PMLR.
- Mancini, M.; Akata, Z.; Ricci, E.; and Caputo, B. 2020. Towards recognizing unseen categories in unseen domains. In *European Conference on Computer Vision*, 466–483. Springer.
- McAllester, D. 2003. Simplified PAC-Bayesian margin bounds. In *Learning Theory and Kernel Machines: 16th Annual Conference on Learning Theory and 7th Kernel Workshop, COLT/Kernel 2003, Washington, DC, USA, August 24–27, 2003. Proceedings*, 203–215. Springer.
- McAllester, D. A. 1998. Some pac-bayesian theorems. In *Proceedings of the eleventh annual conference on Computational learning theory*, 230–234.
- Mohajerin Esfahani, P.; and Kuhn, D. 2018. Data-driven distributionally robust optimization using the Wasserstein metric: performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1–2): 115–166.
- Parascandolo, G.; Neitz, A.; ORVIETO, A.; Gresele, L.; and Schölkopf, B. 2020. Learning explanations that are hard to vary. In *International Conference on Learning Representations*.
- Peng, X.; Bai, Q.; Xia, X.; Huang, Z.; Saenko, K.; and Wang, B. 2019. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, 1406–1415.
- Sagawa, S.; Koh, P. W.; Hashimoto, T. B.; and Liang, P. 2020. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations*.



Shankar, S.; Piratla, V.; Chakrabarti, S.; Chaudhuri, S.; Jyothi, P.; and Sarawagi, S. 2018. Generalizing Across Domains via Cross-Gradient Training. In *International Conference on Learning Representations*.

Shu, Y.; Cao, Z.; Wang, C.; Wang, J.; and Long, M. 2021. Open domain generalization with domain-augmented meta-learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9624–9633.

Sun, B.; and Saenko, K. 2016. Deep coral: Correlation alignment for deep domain adaptation. In *Computer Vision—ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8–10 and 15–16, 2016, Proceedings, Part III 14*, 443–450. Springer.

Vaze, S.; Han, K.; Vedaldi, A.; and Zisserman, A. 2021. Open-set recognition: A good closed-set classifier is all you need? *arXiv preprint arXiv:2110.06207*.

Venkateswara, H.; Eusebio, J.; Chakraborty, S.; and Panchanathan, S. 2017. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5018–5027.

Villani, C.; et al. 2009. *Optimal transport: old and new*, volume 338. Springer.

Wang, H.; Ge, S.; Lipton, Z.; and Xing, E. P. 2019. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, volume 32.

Wang, J.; Lan, C.; Liu, C.; Ouyang, Y.; Qin, T.; Lu, W.; Chen, Y.; Zeng, W.; and Yu, P. 2022. Generalizing to unseen domains: A survey on domain generalization. *IEEE Transactions on Knowledge and Data Engineering*.

Wang, P.; Zhang, Z.; Lei, Z.; and Zhang, L. 2023. Sharpness-aware gradient matching for domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3769–3778.

Xiao, Z.; Shen, J.; Zhen, X.; Shao, L.; and Snoek, C. 2021. A bit more bayesian: Domain-invariant learning with uncertainty. In *International Conference on Machine Learning*, 11351–11361. PMLR.

Zhang, H.; Cisse, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*.

Zhou, K.; Yang, Y.; Qiao, Y.; and Xiang, T. 2021. Mixstyle neural networks for domain generalization and adaptation. *arXiv preprint arXiv:2107.02053*.

Zhu, W.; Lu, L.; Xiao, J.; Han, M.; Luo, J.; and Harrison, A. P. 2022. Localized adversarial domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7108–7118.

## Algorithm for EDM

---

### Algorithm 1: A learner attaching EDM

---

**Input:** training set, batch size  $B$ , max epoch  $T$ , number of source domains  $N$ , Dirichlet Distribution parameter  $\alpha_i$  for interpolation, Dirichlet Distribution parameter  $\alpha_e$  for extrapolation, momentum weight  $\rho$ , a threshold  $\eta$ , hyperparameters of this learner  $\beta$ , learning rate  $\text{lr}$

**Output:** The model of this learner  $f_\theta$

```

1: Initialize  $f_\theta$  with pretrained model
2: Give the corresponding loss function of this learner  $\ell$ 
3: for  $t \leftarrow 1$  to  $T$  do
4:   Randomly sample a batch from training set:  $S = \{(x^b, y^b, d^b)\}_{b=1}^B$ 
5:   # Calculate statistics for each domain
6:   for  $n \leftarrow 1$  to  $N$  do
7:      $\mu_n, \sigma_n \leftarrow \text{Eq. (4) in current batch}$ 
8:      $\mu_n^t, \sigma_n^t \leftarrow \text{Eq. (5) with } \rho$ 
9:   end for
10:  # Extrapolation stage: Calculate statistics for the supported domain according to each domain
11:  for  $n \leftarrow 1$  to  $N$  do
12:    repeat
13:       $\lambda \sim \text{Dirichlet}(\alpha_e)$ 
14:    until  $\lambda_n > \eta$ 
15:     $(\mu_\lambda^E)_n, (\sigma_\lambda^E)_n \leftarrow \text{Eq. (8) with } \{\mu_n^t\}_{n=1}^N, \{\sigma_n^t\}_{n=1}^N$ 
16:  end for
17:  # Interpolation stage: Calculate statistics for the augmented domain
18:   $\lambda \sim \text{Dirichlet}(\alpha_i)$ 
19:   $\mu_a, \sigma_a \leftarrow \text{Eq. (3) with } \{(\mu_\lambda^E)_n\}_{n=1}^N, \{(\sigma_\lambda^E)_n\}_{n=1}^N$ 
20:  # Convert each sample  $x^b$  into the augmented domain
21:  for  $b \leftarrow 1$  to  $B$  do
22:     $x_a^b \leftarrow \text{Eq. (9) with } x^b, \mu_{d^b}^t, \sigma_{d^b}^t \text{ and } \mu_a, \sigma_a$ 
23:     $y_a^b \leftarrow y^b$ 
24:     $d_a^b \leftarrow N + 1$ 
25:  end for
26:   $S_A = \{(x_a^b, y_a^b, d_a^b)\}_{b=1}^B$ 
27:  # Training model
28:   $S_{New} = S \cup S_A$ 
29:   $R(\theta) \leftarrow \frac{1}{2B} \sum_{b=1}^{2B} \ell(f_\theta, \beta; x^b, y^b, d^b)$ 
30:   $\theta \leftarrow \theta - \text{lr} \nabla R(\theta)$ 
31: end for

```

---

We should re-emphasize that our EDM only generates a set of additional input data with a different domain distribution. Therefore, it can be plugged into almost previous methods. In other words,  $\ell$  is not limited to Cross-Entropy (CE) loss. For example,  $\ell$  can be a CE loss for classification plus a domain discrimination loss in DANN. Moreover, Line #12-14 control the similarity between  $D_n$  and  $D_n^E$  indirectly.

### Proof of Theorem 1

In this section, we prove Theorem 1 in detail, which constitutes two components: *task generalization error* and *environment generalization error*.

The notations involved in this section are the same as those in Preliminaries subsection. In addition, let  $\ell(\cdot, \cdot) : \mathcal{H} \times \mathcal{Z} \rightarrow [0, 1]$  denotes a convex loss function, where  $\mathcal{Z}$  consists of the input space  $\mathcal{X}$  and the output space  $\mathcal{Y}$ . The expected error can be defined as  $er(h, \mathcal{D}) \triangleq \mathbb{E}_{z \sim \mathcal{D}} \ell(h, z)$ . In practice, since the real distribution  $\mathcal{D}$  is unknown,  $er(h, \mathcal{D})$  can not be computed directly. Therefore, as a substitute, its corresponding empirical error  $\hat{er}(h, \mathcal{D}) \triangleq \frac{1}{M} \sum_{m=1}^M \ell(h, z_m)$ .

Due to domain shift, the choice of  $h$  from different domains has changed. Therefore, similar to PAC-Bayes framework, whose key is change of measure inequality, the task generalization bound can be given as follows:

**Theorem 2** (Task Generalization Bound). *Let  $\mathcal{M}$  denote the set of all probability measures over  $\mathcal{H}$ .  $P \in \mathcal{M}$  is the prior distribution over  $\mathcal{H}$ . Then, for any  $\delta \in (0, 1]$ , the following inequality holds uniformly for all posteriors distribution  $Q \in \mathcal{M}$  with probability at least  $1 - \delta$ ,*

$$\mathbb{E}_{h \sim Q} er(h, \mathcal{D}) \leq \mathbb{E}_{h \sim Q} \hat{er}(h, \mathcal{D}) + \sqrt{\frac{L \cdot W_1(Q, P) + \ln(M/\delta)}{2(M-1)}}$$

where  $W_1(Q, P)$  is the 1st order Wasserstein Distance and  $L$  is the Lipschitz constant.

*Proof.* The proof follows the similar structure as the classical derivation (McAllester 2003) and the IPM-PB (Amit et al. 2022).

**Proposition 3** (Kantorovich-Rubinstein Duality (Villani et al. 2009)). *For any  $L \geq 0$ , and any two probability measure  $P, Q \in \mathcal{M}$ ,*

$$L \cdot W_1(Q, P) = \sup_{f \in \mathcal{F}_L^{Lip}} |\mathbb{E}_{h \sim P} f(h) - \mathbb{E}_{h \sim Q} f(h)|$$

where  $\mathcal{F}_L^{Lip}$  is the set of  $L$ -Lipschitz functions w.r.t.  $\ell(\cdot, \cdot)$ .

Along this line, we can observe that for any pairs  $(P, Q)$ ,

$$\mathbb{E}_{h \sim Q} f(h) - \mathbb{E}_{h \sim P} f(h) \leq L \cdot W_1(Q, P)$$

Therefore, we have

$$\begin{aligned} \exp(\mathbb{E}_{h \sim Q} f(h) - L \cdot W_1(Q, P)) &\leq \exp(\mathbb{E}_{h \sim P} f(h)) \\ &\leq \mathbb{E}_{h \sim P} [\exp(f(h))] \end{aligned}$$

where the last inequality is by Jensen's inequality.

And then, giving an expectation over samples  $z \sim \mathcal{D}$  in the supremum over  $Q \in \mathcal{M}$ , we have that for any  $P \in \mathcal{M}$ ,

$$\begin{aligned} \mathbb{E}_{z \sim \mathcal{D}} \sup_Q \{ \exp(\mathbb{E}_{h \sim Q} f(h) - L \cdot W_1(Q, P)) \} \\ \leq \mathbb{E}_{z \sim \mathcal{D}} \mathbb{E}_{h \sim P} [\exp(f(h))] \\ = \mathbb{E}_{h \sim P} \mathbb{E}_{z \sim \mathcal{D}} [\exp(f(h))] \end{aligned} \quad (10)$$

Considering Lemma 5 in (McAllester 2003), by setting  $f(h) = 2(M-1)\Delta^2(h)$ , we have

$$\mathbb{E}_{z \sim \mathcal{D}} [\exp(f(h))] = \mathbb{E}_{z \sim \mathcal{D}} [\exp(2(M-1)\Delta^2(h))] \leq M \quad (11)$$

Combining inequalities (10) and (11), by Markov's inequality, for any  $t > 0$ , we have

$$\mathbb{P}_{z \sim \mathcal{D}} \left( \ln \left( \sup_Q \{ \exp(\mathbb{E}_{h \sim Q} f(h) - L \cdot W_1(Q, P)) \} \right) \geq \ln(t) \right) \leq \frac{M}{t}$$

Since the  $\ln(\cdot)$  and  $\sup(\cdot)$  operations are interchangeable, we have

$$\mathbb{P}_{z \sim \mathcal{D}} \left( \sup_Q \{ (\mathbb{E}_{h \sim Q} f(h) - L \cdot W_1(Q, P)) \} \geq \ln(t) \right) \leq \frac{M}{t}$$

Let  $\delta \in (0, 1)$ , by setting  $t = \frac{M}{\delta}$  and replacing  $f(h)$  with  $2(M-1)\Delta^2(h)$ , we get

$$\mathbb{P}_{z \sim \mathcal{D}} \left( \sup_Q \{ (\mathbb{E}_{h \sim Q} (2(M-1)\Delta^2(h)) - L \cdot W_1(Q, P)) \} < \ln\left(\frac{M}{\delta}\right) \right) \leq 1 - \delta$$

Therefore, for any  $P \in \mathcal{M}$ , with a probability of at least  $1 - \delta$  over samples  $z \sim \mathcal{D}$ , the following inequality holds for all  $Q \in \mathcal{M}$

$$\mathbb{E}_{h \sim Q} (\Delta^2(h)) < \frac{L \cdot W_1(Q, P) + \ln(M/\delta)}{2(M-1)}$$

Recall Jensen's inequality, we have

$$(\mathbb{E}_{h \sim Q} \Delta(h))^2 < \frac{L \cdot W_1(Q, P) + \ln(M/\delta)}{2(M-1)}$$

So that, we get

$$\mathbb{E}_{h \sim Q} \text{er}(h, \mathcal{D}) \leq \mathbb{E}_{h \sim Q} \hat{\text{er}}(h, D) + \sqrt{\frac{L \cdot W_1(Q, P) + \ln(M/\delta)}{2(M-1)}}$$

□

In analogy to the standard single-task learning where data is sampled from an unknown distribution, tasks are sampled from an unknown task distribution. Therefore, a task can be regarded as a point sampling from task distribution. In this way, with an abuse of notations, the environment generalization bound similar to the task generalization bound aforementioned can be given.

**Theorem 4** (Environment Generalization Bound). *Let  $\mathcal{P}$  denotes a prior distribution over the priors, i.e. a hyper-prior distribution, while  $\mathcal{Q}$  is a hyper-posterior distribution analogous to  $\mathcal{P}$ . And,  $\tau$  is a hyper-distribution over domains. Then, for any  $\delta \in (0, 1]$ , the following inequality holds uniformly for all hyper-posteriors distribution  $\mathcal{Q}$  with probability at least  $1 - \delta$ ,*

$$\mathbb{E}_{Q \sim \mathcal{Q}} \text{er}(Q, \tau) \leq \mathbb{E}_{P \sim \mathcal{Q}} \hat{\text{er}}(P, T) + \sqrt{\frac{L_0 \cdot W_1(\mathcal{Q}, \mathcal{P}) + \ln(N/\delta)}{2(N-1)}}$$

where  $W_1(\mathcal{Q}, \mathcal{P})$  is the 1st order Wasserstein Distance,  $L_0$  is the Lipschitz constant, and  $T = \{\mathcal{D}_n\}_{n=1}^N$ .

*Proof.* The proof is similar to task generalization bound in Theorem 2. □

Finally, combining Theorem 2 and Theorem 4, for any  $\delta > 0$ , set  $\delta_0 \triangleq \frac{\delta}{2}$  and  $\delta_i \triangleq \frac{\delta}{2N}$  for  $i = 1, \dots, N$ .

$$\begin{aligned} \mathbb{E}_{Q \sim \mathcal{Q}} \text{er}(Q, \tau) &\leq \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{P \sim \mathcal{Q}} \hat{\text{er}}(P, D_n) \\ &\quad + \sqrt{\frac{L_0 \cdot W_1(\mathcal{Q}, \mathcal{P}) + \ln(N/\delta)}{2(N-1)}} \\ &\quad + \frac{1}{N} \sum_{n=1}^N \sqrt{\frac{L_n \cdot W_1(Q, P_n) + \ln(NM_n/\delta)}{2(M_n-1)}} \end{aligned}$$

□

## Review of PAC-Bayes Framework

### EDM Meets Open-set Scenario

EDM is a domain augmentation strategy, which is unrelated to special task. Meanwhile, we can find that each augmented domain contains all the task information in a batch through Algo. 1. Therefore, in DAML + EDM, through its Mixup scheme, more complete task information can be fused to promote better detection of known classes. In this way, EDM can indirectly improve the performance of unknown classes detection. This intuition is consistent with the viewpoint in (Vaze et al. 2021), and from Tab. 8 to Tab. 11, it can be empirically testified.

### Is There a Need for the Threshold $\eta$ ?

Intuitively, each  $D_n^E$  is usually a worse case than the corresponding  $D_n$  in training. Therefore, according to the motivation in distributionally robust optimization (Mohajerin Esfahani and Kuhn 2018), that is a zero-sum game, the supported domains with proper discrepancy are required. With the use of threshold  $\eta$ , the similarity between  $D_n$  and  $D_n^E$  can be controlled indirectly. If  $\eta$  is large,  $D_n^E$  will approximate  $D_n$ . Conversely, if  $\eta$  is too small, it is more likely that  $D_n^E$  will diverge from  $D_n$ . This intuition has been empirically presented in Fig. 5.

## Additional Experiments

### Datasets

We adopt three conventional DG benchmarks: 1) PACS (Li et al. 2017), 2) Office-Home (Venkateswara et al. 2017), 3) DomainNet (Peng et al. 2019).

PACS is an object classification benchmark with four domains (art-painting, cartoon, phtot, sketch). There exist large discrepancies in image styles among different domains. Each domain contains seven classes and there are 9,991 images in total.

Office-Home consists of images from 4 different domains: art, clipart, product and real-world. It contains around 15,500 images from 65 different categories.

DomainNet contains 6 domains, i.e., clipart, infograph, painting, quickdraw, real and sketch, and about 0.6 million images distributed among 345 categories.

### Settings

In this subsection, the settings in both closed and open set are introduced in detail.

In common setting, we adopt ResNet-18 architecture as the backbone, and each input image is resized and cropped to  $224 \times 224$ . The optimizer is Stochastic gradient descent (SGD), where the momentum and the weight decay are set to 0.9 and  $5e-4$ , respectively. And, in EDM, all Dirich Distribution parameters  $\alpha_i$  and  $\alpha_e$  are set to 0.5 while the threshold  $\eta$  is set to 0.9, and the momentum weight  $\rho$  is set to 0.99.

For closed-set scenario, we follow the settings in (Lu et al. 2022). Specifically, the training data are randomly split into two parts: 80% for training and 20% for validation. The best model on the validation split is selected to evaluate the target domain. Batch size is set to 32 for each

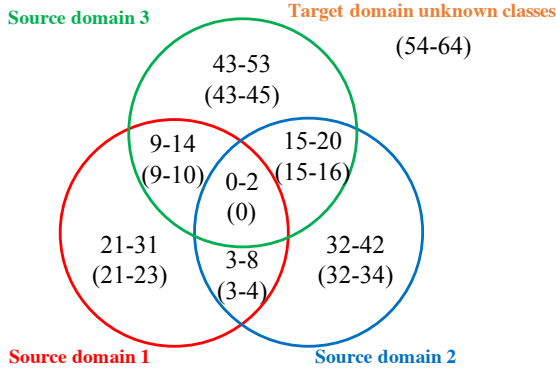


Figure 4: Illustration of the open-domain split of Office-Home dataset. Indices without brackets show the distribution of categories among source domains, while indices in brackets indicate the categories of the target domain.

Domain	Classes
Source 1	3, 0, 1
source 2	4, 0, 2
source 3	5, 1, 2
Target	0, 1, 2, 3, 4, 5, 6

Table 3: Open-domain split of PACS dataset.

source domain, and the number of epochs is set to 120. The learning rate is chosen from  $\{5e-4, 1e-3, 5e-3\}$ , and is decayed by 0.1 twice at the 70% and 90% of the max epoch, respectively. The hyper-parameter  $\alpha$  in SAGM is chosen from  $\{1e-5, 5e-5, 1e-4, 5e-4, 1e-3, 5e-3, 1e-2\}$ . The hyper-parameter  $\alpha$ ,  $\beta$ , and  $\lambda$  in DIFEX are chosen from  $\{1e-3, 1e-2, 1e-1\}$ ,  $\{0, 1e-2, 1e-1, 5e-1, 1e0, 1e1\}$ , and  $\{0, 1e-2, 1e-1, 1e0\}$ , respectively. Thanks to Wang and Lu for providing the implementations of previous methods in DeepDG (Deep Domain Generalization Toolkit) published on GitHub. All codes are run on Python 3.9 and Torch 1.12 on Arch Linux with an NVIDIA GeForce RTX 2080Ti GPU.

For open-set scenario, we follow the settings in (Shu et al. 2021). Specifically, default training and validation splits are adopted, and the best model on the validation split with known classes is selected to evaluate the target domain. The specific categories contained in each domain are shown in Tab. 3 and Fig. 4, which are the same as (Shu et al. 2021). Batch size is set to 24 for each source domain, and the number of epochs is set to 30. The learning rate is set to  $1e-3$ , and is decayed after 24 epochs by a factor of 10. The four hyper-parameters are chosen from  $\{0.1, 1, 5, 10\}$ ,  $\{0.1, 1, 5\}$ ,  $\{0.1, 1, 5\}$ , and  $\{0.1, 1, 5, 10\}$ , respectively. And, the other, such as temperature and meta step, are set to the corresponding default values. All codes are run on Python 3.8 and Torch 1.11 on Ubuntu with an NVIDIA GeForce RTX 3090 GPU.

## Learners attaching EDM

In this subsection, the implementations of the learners attaching EDM are described in detail. EDM, as shown in Algo. 1, generates the new samples from an augmented domain, whose domain label is  $N + 1$ . Therefore, these generated samples are fed into the network as the samples sampling from  $(N + 1)$ -th domain.

In particular, in DANN + EDM, the output dimension of the domain discriminator is  $N + 1$ , and the coding dimension of the one-hot coding is increased by 1 accordingly. Although the new domain will be generated for each iteration, all of them are only regarded as  $(N + 1)$ -th domain.

For the sample mixing scheme, the augmented domain together with the observed source domains is deployed to generate new samples. For example, in Mixup + EDM, the augmented samples are mixed through the samples sampling from two domains, where the domains are randomly selected from the augmented domains and the observed source domains. The same goes for DAML + EDM, where the augmented domain is split into  $N$  parts and each part is regarded as a new domain to be fed into a sub-network.

## Appending Results

In this subsection, we report more detailed results for EDM. Accuracy results are reported for both settings, and H-score (Fu et al. 2020) results are additionally reported for open-set settings, which measures the quality of unknown class detection.

**Tab. 5** and **Tab. 6**, which report more detailed accuracy results in closed-set setting both on PACS and Office-Home datasets, reveal that the basic learner attaching the augmented domains can improve performance. We also can find that the target domain with significant discrepancies compared to the source domains can achieve better performance improvements.

**Tab. 7** report more detailed accuracy results in closed-set setting on DomainNet dataset. We can observe: 1) Deploying EDM can indeed improve all learners' performance as highlighted, and can achieve the best of most domains. 2) Mixstyle is worse than our proposed methods. To sum up, EDM can be effective on large-scale dataset as well and has generality.

**Tab. 8** and **Tab. 9** report more detailed accuracy results in open-set setting while **Tab. 10** and **Tab. 11** report more detailed H-score results. In general, DAML attaching the augmented domains not only can improve accuracy but also can improve H-score for each domain. Moreover, the improvement in accuracy is significantly greater than that in H-score. These phenomena reveal the effectiveness of EDM in open-set setting. And, these indicate that EDM is not limited to a specific setting and has well flexibility, simultaneously. Furthermore, similar to the closed-set setting, the target domain with significant discrepancies can better performance improvements.

From **Fig. 5**, we can find that the selection of  $\eta$  can affect the performance, and the best results can usually be obtained when  $\eta = 0.9$ . This phenomenon testifies that the supported domains with proper discrepancy outside the interpolation space are beneficial for improving model performance.

The computational cost for each batch of EDM is approximately 68.57 ms. In addition, we report the average training computational cost here. From **Tab. 4**, we can observe: 1) The computation cost naturally increases due to data augmentation, but this increase is less than double. 2) ERM + EDM can achieve good performance with less computational cost. 3) Thanks to Mixup scheme, the computational cost of learner + EDM is much lower than that of CrossGrad, which is an optimization-based data augmentation.

	A	C	P	S	Avg
ERM	239.86	215.87	217.22	218.09	222.76
ERM + EDM	363.04	349.61	349.61	352.72	353.74
DANN	217.47	219.30	219.09	219.78	218.91
DANN + EDM	371.23	407.17	405.85	371.44	388.92
Mixup	379.85	313.86	366.87	376.33	359.23
Mixup + EDM	504.18	524.75	554.53	552.25	533.93
SAGM	366.64	369.04	369.54	379.00	371.06
SAGM + EDM	667.01	664.53	669.53	684.50	671.39
CrossGrad	972.79	987.57	986.95	1021.38	992.17

Table 4: Time cost (ms) of each batch on PACS dataset.

## References

- Amit, R.; Epstein, B.; Moran, S.; and Meir, R. 2022. Integral Probability Metrics PAC-Bayes Bounds. In *Advances in Neural Information Processing Systems*, volume 35, 3123–3136.
- Fu, B.; Cao, Z.; Long, M.; and Wang, J. 2020. Learning to detect open classes for universal domain adaptation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16*, 567–583. Springer.
- Li, D.; Yang, Y.; Song, Y.-Z.; and Hospedales, T. M. 2017. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, 5542–5550.
- McAllester, D. 2003. Simplified PAC-Bayesian margin bounds. In *Learning Theory and Kernel Machines: 16th Annual Conference on Learning Theory and 7th Kernel Workshop, COLT/Kernel 2003, Washington, DC, USA, August 24–27, 2003. Proceedings*, 203–215. Springer.
- Mohajerin Esfahani, P.; and Kuhn, D. 2018. Data-driven distributionally robust optimization using the Wasserstein metric: performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1-2): 115–166.
- Peng, X.; Bai, Q.; Xia, X.; Huang, Z.; Saenko, K.; and Wang, B. 2019. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, 1406–1415.
- Vaze, S.; Han, K.; Vedaldi, A.; and Zisserman, A. 2021. Open-set recognition: A good closed-set classifier is all you need? *arXiv preprint arXiv:2110.06207*.
- Venkateswara, H.; Eusebio, J.; Chakraborty, S.; and Panchanathan, S. 2017. Deep hashing network for unsupervised

domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5018–5027.

Villani, C.; et al. 2009. *Optimal transport: old and new*, volume 338. Springer.

	Art-Painting	Cartoon	Photo	Sketch	Avg
ERM	81.10	77.94	95.03	76.94	82.75
DANN	82.86	78.33	96.11	76.99	83.57
Mixup	81.84	75.43	95.27	76.51	82.26
RSC	82.13	77.99	94.43	79.87	83.60
MMD	80.32	76.45	92.46	<b>83.63</b>	83.21
CORAL	79.39	77.9	91.98	82.03	82.83
GroupDRO	79.15	76.75	91.32	81.52	82.19
CrossGrad $\ddagger$	80.37	74.87	<u>96.59</u>	74.98	81.70
Mixstyle $\ddagger$	82.51	79.09	95.65	79.23	84.12
ANDMask	80.81	73.29	95.81	71.95	80.47
Vrex	81.54	78.11	95.39	80.35	83.85
DIFEX-ori	82.86	78.46	94.97	79.41	83.93
DIFEX-norm	<u>83.40</u>	<u>79.74</u>	95.03	79.10	84.32
SAGM $\ddagger$	82.62	78.50	96.05	79.64	84.20
ERM + Inter	83.25 (2.15 $\uparrow$ )	77.60 (0.34 $\downarrow$ )	95.99 (0.96 $\uparrow$ )	81.09 (4.15 $\uparrow$ )	84.48 (1.73 $\uparrow$ )
DANN + Inter	81.93 (0.93 $\downarrow$ )	77.82 (0.51 $\downarrow$ )	95.99 (0.12 $\downarrow$ )	81.57 (4.58 $\uparrow$ )	84.33 (0.76 $\uparrow$ )
Mixup + Inter	83.01 (1.17 $\uparrow$ )	76.32 (0.89 $\uparrow$ )	<b>96.65</b> (1.38 $\uparrow$ )	78.04 (1.53 $\uparrow$ )	83.50 (1.24 $\uparrow$ )
SAGM + Inter	82.28 (0.34 $\downarrow$ )	79.01 (0.51 $\uparrow$ )	96.29 (0.24 $\uparrow$ )	80.22 (0.58 $\uparrow$ )	84.45 (0.25 $\uparrow$ )
ERM + EDM	82.32 (1.22 $\uparrow$ )	79.27 (1.33 $\uparrow$ )	96.53 (1.50 $\uparrow$ )	81.24 (4.30 $\uparrow$ )	84.84 (2.09 $\uparrow$ )
DANN + EDM	82.96 (0.10 $\uparrow$ )	78.07 (0.26 $\downarrow$ )	96.47 (0.36 $\uparrow$ )	<u>82.72</u> (5.73 $\uparrow$ )	85.06 (1.49 $\uparrow$ )
Mixup + EDM	<b>83.50</b> (1.66 $\uparrow$ )	79.14 (3.71 $\uparrow$ )	<u>96.59</u> (1.32 $\uparrow$ )	81.04 (4.53 $\uparrow$ )	<u>85.07</u> (2.81 $\uparrow$ )
SAGM + EDM	82.47 (0.15 $\downarrow$ )	<b>80.38</b> (1.88 $\uparrow$ )	<u>96.59</u> (0.54 $\uparrow$ )	80.86 (1.22 $\uparrow$ )	<b>85.08</b> (0.88 $\uparrow$ )

Table 5: More detailed accuracy results on PACS in closed-set settings.  $\uparrow$  denotes improved performance while  $\downarrow$  denotes the opposite.  $\ddagger$  denotes our reproduced results.

	Art	Clipart	Product	Real-World	Avg
ERM	57.77	50.63	71.3	74.45	63.54
DANN	57.6	48.52	71.16	72.99	62.57
Mixup	58.71	51	72.2	75.42	64.33
RSC	57.07	50.77	71.93	73.63	63.35
MMD	59.29	50.52	72.34	74.43	64.15
CORAL	59.29	50.15	72.25	74.2	63.97
GroupDRO	59.09	50.22	71.91	74.48	63.92
CrossGrad $\ddagger$	58.67	51.18	71.66	74.80	64.08
Mixstyle $\ddagger$	55.50	51.00	70.62	73.19	62.57
ANDMask	53.61	47.54	69.36	72.23	60.69
Vrex	59.09	49.81	71.64	74.82	63.84
DIFEX-ori $\ddagger$	57.89	50.82	71.61	73.40	63.43
DIFEX-norm $\ddagger$	58.09	51.50	72.08	73.62	63.82
SAGM $\ddagger$	59.13	51.23	72.67	75.90	64.73
ERM + Inter	58.01 (0.24 $\uparrow$ )	50.65 (0.02 $\uparrow$ )	72.02 (0.72 $\uparrow$ )	74.62 (0.17 $\uparrow$ )	63.83 (0.29 $\uparrow$ )
DANN + Inter	57.27 (0.33 $\downarrow$ )	50.13 (1.61 $\uparrow$ )	71.53 (0.37 $\uparrow$ )	74.04 (1.05 $\uparrow$ )	63.24 (0.67 $\uparrow$ )
Mixup + Inter	<b>59.46</b> (0.75 $\uparrow$ )	52.30 (1.30 $\uparrow$ )	72.88 (0.68 $\uparrow$ )	75.60 (0.18 $\uparrow$ )	<u>65.06</u> (0.73 $\uparrow$ )
SAGM + Inter	58.55 (0.58 $\downarrow$ )	<b>52.71</b> (1.48 $\uparrow$ )	72.85 (0.18 $\uparrow$ )	75.51 (0.39 $\downarrow$ )	64.91 (0.18 $\uparrow$ )
ERM + EDM	58.67 (0.90 $\uparrow$ )	51.84 (1.21 $\uparrow$ )	72.38 (1.08 $\uparrow$ )	75.35 (0.90 $\uparrow$ )	64.56 (1.02 $\uparrow$ )
DANN + EDM	58.51 (0.91 $\uparrow$ )	50.61 (2.09 $\uparrow$ )	72.22 (1.06 $\uparrow$ )	74.59 (1.60 $\uparrow$ )	63.98 (1.41 $\uparrow$ )
Mixup + EDM	<u>59.33</u> (0.62 $\uparrow$ )	51.94 (0.94 $\uparrow$ )	<b>73.15</b> (0.95 $\uparrow$ )	<u>75.97</u> (0.55 $\uparrow$ )	<b>65.10</b> (0.77 $\uparrow$ )
SAGM + EDM	58.84 (0.29 $\downarrow$ )	<u>52.33</u> (1.10 $\uparrow$ )	<u>72.94</u> (0.27 $\uparrow$ )	<b>76.06</b> (0.16 $\uparrow$ )	65.04 (0.31 $\uparrow$ )

Table 6: More detailed accuracy results on Office-Home in closed-set settings.  $\uparrow$  denotes improved performance while  $\downarrow$  denotes the opposite.  $\ddagger$  and  $\ddagger$  denote our reproduced results.



	Clipart	Infograph	Painting	Quickdraw	Real	Sketch	Avg
ERM	58.12	17.50	45.21	13.60	57.87	48.28	40.10
DANN	58.04	<u>17.52</u>	<u>45.71</u>	13.77	<b>58.21</b>	47.95	40.20
Mixup	55.63	16.84	44.83	<b>15.15</b>	55.01	47.98	39.24
RSC	52.96	16.99	42.44	13.50	51.25	45.62	37.13
MMD	57.98	16.49	42.81	12.93	56.15	48.46	39.14
CORAL	57.96	16.79	43.69	13.25	56.32	<b>48.73</b>	39.46
GroupDRO	46.93	14.07	32.31	10.77	46.78	39.79	31.78
CrossGrad	56.96	16.87	45.07	13.77	57.21	47.04	39.49
Mixstyle	58.17	17.20	45.23	13.52	57.75	48.02	39.98
ANDMask	33.90	8.82	27.96	11.72	30.29	30.81	23.92
Vrex	54.14	16.72	42.20	13.71	55.24	45.72	37.96
DIFEX-ori	55.81	15.01	42.29	13.30	55.93	47.65	38.33
DIFEX-norm	56.11	15.33	43.24	13.51	55.73	47.28	38.53
SAGM	56.47	16.79	43.35	13.90	55.51	47.12	38.86
ERM + EDM	<b>58.86</b> (0.74 ↑)	<b>17.53</b> (0.03 ↑)	45.65 (0.44 ↑)	14.20 (0.60 ↑)	57.52 (0.35 ↓)	48.29 (0.01 ↑)	<u>40.34</u> (0.24 ↑)
DANN + EDM	<u>58.72</u> (0.68 ↑)	17.36 (0.16 ↓)	<b>45.84</b> (0.13 ↑)	14.18 (0.41 ↑)	<u>58.07</u> (0.14 ↓)	48.20 (0.25 ↑)	<b>40.40</b> (0.20 ↑)
Mixup + EDM	58.17 (2.54 ↑)	17.08 (0.24 ↑)	44.93 (0.10 ↑)	14.66 (0.49 ↓)	57.33 (2.32 ↑)	48.07 (0.09 ↑)	40.04 (0.80 ↑)
SAGM + EDM	57.09 (0.62 ↑)	17.21 (0.42 ↑)	43.94 (0.59 ↑)	<u>14.21</u> (0.31 ↑)	55.61 (0.10 ↑)	47.32 (0.20 ↑)	39.23 (0.37 ↑)

Table 7: More detailed accuracy results on DomainNet in closed-set settings. ↑ denotes improved performance while ↓ denotes the opposite. All results have been reproduced.

	Art-Painting	Cartoon	Photo	Sketch	Avg
ERM	51.35	66.43	53.15	49.75	55.17
MLDG	44.59	71.64	62.20	51.29	57.43
FC	51.12	69.32	60.94	51.15	58.13
Epi-FCR	54.16	72.00	70.03	46.35	60.64
PAR	52.97	67.77	51.86	53.62	56.56
RSC	50.47	67.51	67.53	50.17	58.92
CuMix	53.85	<b>74.16</b>	65.67	37.70	57.85
DAML	54.10	<u>73.65</u>	75.69	58.50	65.49
DAML + Inter	<u>59.70</u> (5.60 ↑)	69.70 (3.95 ↓)	<u>80.40</u> (4.71 ↑)	<u>67.10</u> (8.60 ↑)	<u>69.22</u> (3.73 ↑)
DAML + EDM	<b>59.80</b> (5.70 ↑)	71.90 (1.75 ↓)	<b>81.90</b> (6.21 ↑)	<b>69.50</b> (11.00 ↑)	<b>70.78</b> (5.29 ↑)

Table 8: More detailed accuracy results on PACS in open-set settings.

	Art	Clipart	Product	Real-World	Avg
ERM	42.22	42.83	54.27	62.40	50.43
MLDG	42.58	41.82	56.89	62.98	51.07
FC	44.13	41.80	54.41	63.79	51.03
Epi-FCR	46.33	37.13	54.95	62.60	50.25
PAR	42.40	41.27	55.37	65.98	51.26
RSC	44.19	38.60	54.61	60.85	49.56
CuMix	42.76	41.54	57.74	64.63	51.67
DAML	53.13	45.13	61.54	65.99	56.45
DAML + Inter	<b>56.40</b> (3.27 ↑)	<u>48.60</u> (3.47 ↑)	<u>62.30</u> (0.76 ↑)	<u>69.30</u> (3.21 ↑)	<u>59.15</u> (2.70 ↑)
DAML + EDM	<b>56.40</b> (3.27 ↑)	<b>49.80</b> (4.67 ↑)	<b>62.61</b> (1.07 ↑)	<b>69.50</b> (3.51 ↑)	<b>59.58</b> (3.13 ↑)

Table 9: More detailed accuracy results on Office-Home in open-set settings.

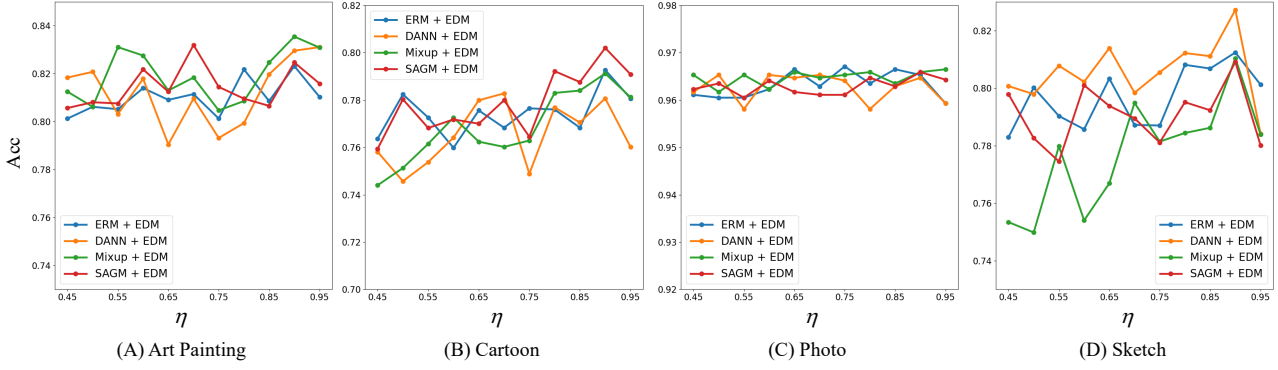


Figure 5: The illustration of threshold  $\eta$  sensitivity analysis on PACS.

	Art-Painting	Cartoon	Photo	Sketch	Avg
ERM	38.87	48.98	44.19	47.09	44.78
MLDG	31.54	<u>55.20</u>	43.35	49.91	45.00
FC	39.01	52.67	45.79	49.28	46.69
Epi-FCR	41.16	<b>58.19</b>	48.38	46.14	48.47
PAR	39.21	52.05	36.53	52.00	44.95
RSC	38.43	47.35	49.82	44.59	45.05
CuMix	38.67	47.53	49.28	28.71	41.05
DAML	43.02	54.47	53.29	<u>56.73</u>	<u>51.88</u>
DAML + Inter	<u>47.07</u> (4.05 $\uparrow$ )	53.30 (1.17 $\downarrow$ )	<b>55.87</b> (2.58 $\uparrow$ )	51.09 (5.64 $\downarrow$ )	51.83 (0.05 $\downarrow$ )
DAML + EDM	<b>48.13</b> (5.11 $\uparrow$ )	54.52 (0.05 $\uparrow$ )	<u>54.52</u> (1.23 $\uparrow$ )	<b>59.32</b> (2.59 $\uparrow$ )	<b>54.12</b> (2.24 $\uparrow$ )

Table 10: More detailed H-score results on PACS in open-set settings.

	Art	Clipart	Product	Real-World	Avg
ERM	40.87	44.98	50.11	53.67	47.41
MLDG	40.97	41.26	52.25	55.84	47.58
FC	43.25	41.65	52.02	55.16	48.02
Epi-FCR	44.46	42.05	52.68	54.73	48.48
PAR	42.62	41.77	54.13	57.60	49.03
RSC	44.77	38.39	54.66	53.73	47.89
CuMix	40.72	43.07	55.79	58.02	49.40
DAML	<b>51.11</b>	43.12	59.00	60.13	53.34
DAML + Inter	<u>49.76</u> (1.35 $\downarrow$ )	<u>43.42</u> (0.30 $\uparrow$ )	<u>59.48</u> (0.48 $\uparrow$ )	<b>61.91</b> (1.78 $\uparrow$ )	<u>53.64</u> (0.30 $\uparrow$ )
DAML + EDM	<u>51.07</u> (0.04 $\downarrow$ )	<b>45.03</b> (1.91 $\uparrow$ )	<b>59.58</b> (0.58 $\uparrow$ )	<u>61.06</u> (0.93 $\uparrow$ )	<b>54.19</b> (0.85 $\uparrow$ )

Table 11: More detailed H-score results on Office-Home in open-set settings.