

Class-Aware Diversified Augmentation for Open-Set Single Domain Generalization

Jian Hu, Shaogang Gong, Weitong Cai, Junchi Yan

Abstract—In Open-Set Single Domain Generalization (OS-SDG), one only has access to a single labeled source domain for training. It assumes that the learned model generalizes well to target samples belonging to the source label space whilst classifies target samples outside the source label space into a single “unknown” class. The current method synthesizes new samples that are semantically unrelated to known classes to simulate target unknown classes. This ignores that unknown classes actually may semantically correlated to known classes, making it difficult to discriminate samples at the margins of class decision boundaries as “unknown”. In this work, we introduce a Class-Aware Diversified Augmentation (CADA) method to overcome this problem. Our key idea is to synthesize explicitly new multiple unknown target classes with diversified semantic and learn the inherent correlation among the known and unknown classes, so to both increase the coverage of multiple target unknown classes and to optimize class margin separation. CADA is optimized by enhanced diversity maximization and class-aware minimization. The former synthesizes more novel classes by considering both semantic relationships to known classes and domain shift between the source and target domains. The latter employs class-agnostic clustering with synthesized samples to simulate class correlations among target classes, maximizing class margin separation. Theoretical analysis and experiments on five benchmarks show the efficacy of our CADA.

Index Terms—Transfer Learning, Semi-supervised Learning, Domain Generalization.

I. INTRODUCTION

Deep learning in computer vision has shown exceptional performance with large amounts of labeled data [8, 9]. But it assumes that training data (source domain) and test data (target domain) are from the same distribution (i.i.d. assumption). Domain Adaptation (DA) [2, 6] relaxes this assumption by minimizing the distribution gap between domains. However, it needs access target domains during training. Alternatively, Domain Generalization (DG) [19] employs multiple sources for learning a domain-invariant representation in order to generalize the model to unseen target domains. Both DA and DG usually assume source and target domains share the same label space. This is not always true. Multi-Source Open Domain Generalization [32] trains a model across multiple source domains with unaligned label spaces, enabling it to generalize well to unseen domains with unknown classes. It has also been applied to large-scale language model training to improve generalization ability [1], but the performance of

Jian Hu, Shaogang Gong and Weitong Cai ({jian.hu, s.gong, weitong.cai}@qmul.ac.uk) are in Queen Mary University of London, UK. Junchi Yan (yanjunchi@sjtu.edu.cn) is in Shanghai Jiao Tong University, CN.

Manuscript received April 19, 2021; revised August 16, 2021.

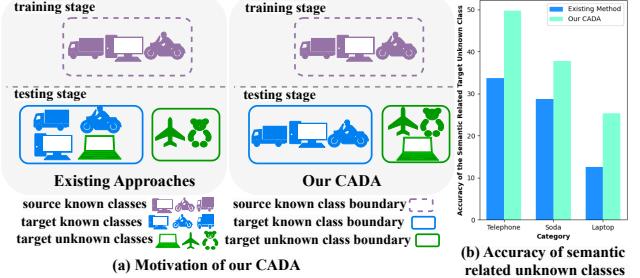


Fig. 1: (a) Motivation of our CADA. Open-Set Single Domain Generalization [51] synthesizes new unknown class examples that are semantically unrelated to the known classes, mimicing target unknown classes in model training. In practice, target unknown classes can either be close to or overlapping with the known classes in the feature space, therefore difficult to be discriminated, e.g. unknown class “bicycle” is similar to known class “motorcar”. CADA explicitly synthesizes unknown classes that are “close” to source known classes in model training in order to improve model generalization in a target domain with both existing known and new unknown classes. (b) Our synthesized unknown classes are semantically similar to known classes whilst CADA can discriminate better than the existing method [51].

such large models in specific scenarios is often unsatisfactory, i.e., medical image analysis [44]. Additionally, due to security and privacy restrictions, collecting labeled data from various domains in these scenarios is also a challenge. A unique approach to minimizing dependency on multiple datasets is Open-Set Single Domain Generalization (OS-SDG) called CrossMatch [51], where a model trained on a single labeled source domain not only generalizes to *unseen* target domains on known classes but also learns to recognize new unknown classes explicitly in the unseen target domain. Tab. I shows comparisons on different DA and DG scenarios.

However, OS-SDG is a harder problem to solve. It needs address both domain shift and label space expansion without accessing target domain training data. CrossMatch [51] synthesizes novel samples with as distinct semantic and domain shifts as possible from the source domain, to simulate unseen target domains during training. But CrossMatch ignores the fact that potential new unknown classes in an unseen target domain may not be well-separated in the feature space from known classes, resulting in target unknown classes being misclassified as target known ones (see Fig. 1(a)). This problem is further aggravated when there is a domain shift between source and

TABLE I: A comparison of different DA and DG settings.

Methods	unaligned label space	only one source domain	inaccessible target domain
Closed-Set Domain Adaptation [35]	x	✓	x
Domain Adaptation with category gap [11, 12, 45]	✓	✓	x
Multi-Source Domain Adaptation [26, 47]	x	x	x
Domain Generalization [19, 25]	x	x	✓
Single Source Domain Generalization [29]	x	✓	✓
Multi-Source Open Domain Generalization [32]	✓	x	✓
Open-Set Single Domain Generalization [51]	✓	✓	✓

target (see Fig. 1(b)). Moreover, since the number of target unknown classes is arbitrary and unknown, CrossMatch treats all simulated target unknown class samples as a single new “unknown” class to the known classes. It ignores any correlations among the unknown classes as well as their relationships to known classes, resulting in suboptimal model learning.

To address the above problems, we introduce Class-Aware Diversified Augmentation (CADA), consisting of two model optimization objectives: enhanced diversity maximization and class-aware minimization. During maximization, diversified unknown classes are first synthesized through unknown maximization. After append synthesized unknown classes into source data, it then synthesizes both known and unknown classes with domain shift to simulate unseen target domains through out-of-distribution maximization. Specifically, The synthesized unknown samples have unique semantics distinct from known classes while still maintaining correlations with the known classes used for their synthesis. After maximization, in order to improve model’s open-set generalization ability, synthesized samples are append to source data for class-aware minimization. Model is trained via class-agnostic clustering and modified supervised learning on expand source data, to maximize class margin whilst simultaneously minimize intra-class scattering in the presence of known classes nearby. **Our contributions are as follows:**

(1) We introduce a new Class-Aware Diversified Augmentation (CADA) method to explicitly synthesize potential target unknown class samples that are “close” to known classes when the target domain is unseen in model training. This is designed to address the limitations of the existing OS-SDG model CrossMatch which ignores any potential class similarities between known and unknown classes.

(2) To implement the CADA idea, we formulate a joint learning objective for both enhanced diversity maximization to synthesize diverse unknown class samples that are similar to those of known classes, and class-aware minimization to perform unsupervised clustering of the synthesized samples without knowing the number of target unknown classes.

(3) Theoretical analysis and comprehensive experiments show that CADA outperforms a wide range of existing generalization and adaptation methods, demonstrating its superiority to unseen target domains of different distribution shifts.

II. RELATED WORKS

Multi-Source Domain Generalization [25, 32, 33, 41] enhance target performance by learning cross-domain invariance from multiple source domains. Kernel-based methods [25], meta-learning [32], and data augmentation [38, 41] are proposed to address this, but fail to identify novel target classes.

[1] leverages meta-learning to extract invariance from source domains to classify known targets and identify unknowns. [33] uses CLIP’s generalization with prompt optimization for unknown class recognition. [3] distills a large model to help a small model generalize to unseen unknown classes. However, these methods rely on multiple source domains, which is impractical in real-world applications.

Single Source Domain Generalization [32] relaxes this assumption. It only needs a labeled source domain for training and can be generalized to multiple unseen target domains under the same label space. [49] presents adversarial gradient-based augmentation to address it and achieves good performance. CrossMatch [51] further introduces Open-Set Single Domain generalization. It assumes that the target domain includes novel classes that do not appear in the source. It simulates target unseen samples by synthesization, but synthesized unknown samples are semantically unrelated to known classes, ignores unknown samples closely related to the known classes.

Data Augmentation with Diversity. Data Augmentation [21, 38] is a max-min game that maximization synthesizes samples on the source distribution edge to simulate out-of-distribution target samples, while minimization learns from them to improve generalization ability. CrossMatch [51] synthesizes novel samples beyond source label space with domain shift. But these samples lack semantic diversity, making it difficult to classify boundary classes. Contrastly, CADA synthesizes diverse novel classes semantically related to the known classes to enhance its open-set generalization capability.

III. CLASS-AWARE DIVERSIFIED AUGMENTATION

Problem Setting. In OS-SDG, the model is trained on an annotated source domain D_s and tested on multiple unannotated target domains $D_t = \{D_{t_1}, D_{t_2}, \dots, D_{t_H}\}$, which are inaccessible during training. The source domain $D_s = \{x_s^i, y_s^i\}_{i=1}^{n_s}$ has n_s labeled samples, while each target domain $D_{t_h} = \{x_{t_h}^i\}_{i=1}^{n_{t_h}}$ has n_{t_h} unlabeled samples. Target domains include novel categories absent from the source domain. The source label space \mathcal{C}_s is a subset of the target \mathcal{C}_t , i.e., $\mathcal{C}_s \subset \mathcal{C}_t$, with novel target classes defined as $\mathcal{C}_t^u = \mathcal{C}_t \setminus \mathcal{C}_s$. During inference, all novel classes \mathcal{C}_t^u are grouped into a single “unknown” class. The label spaces of multiple target domains $\mathcal{C}_{t_1}, \dots, \mathcal{C}_{t_H}$ are not fixed but all include \mathcal{C}_s . Additionally, there is a distribution shift between source and target domains: $\mathbb{P}_s(x) \neq \mathbb{P}_t(x)$.

Overview. We introduce the Class-Aware Diversified Augmentation (CADA). Fig.2 shows that CADA consists of both enhanced diversity maximization and class-aware minimization. The maximization stage further comprises unknown maximization and out-of-distribution maximization. The former generates unknown samples without domain shift, preserving unique semantics distinct from known classes while maintaining correlations with the known samples used for synthesis. The latter extends this by synthesizing both known and unknown classes under domain shift, simulating unseen target domains. Minimization combines original and synthesized samples to improve open-set generalization ability.

Preliminaries In this section, we review the principle of worst-case problem [34] and adversarial domain generation [38]. In single domain generalization, worst-case problem is employed

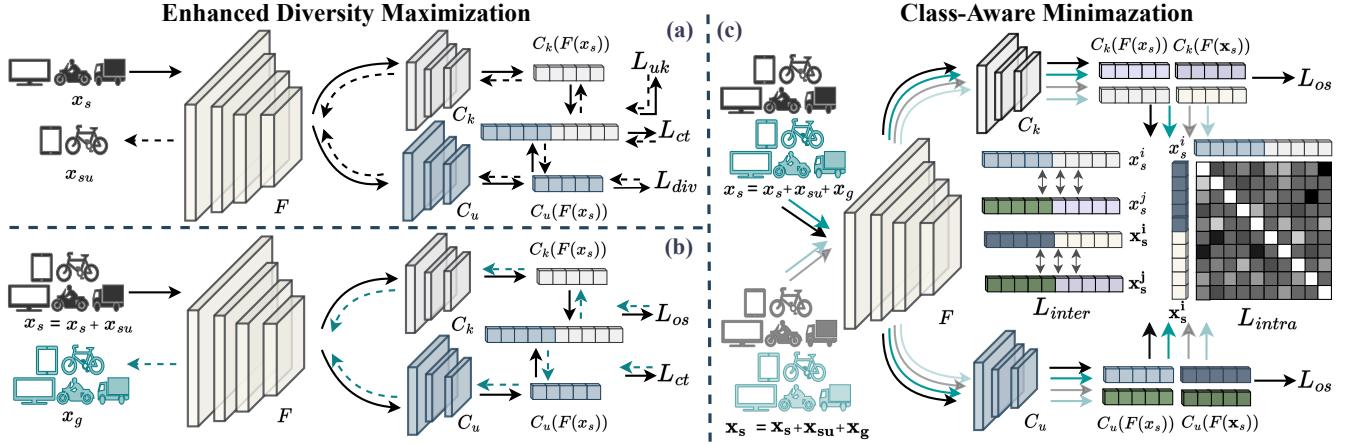


Fig. 2: CADA includes a feature extractor F , a known class classifier C_k , and an unknown class classifier C_u , with enhanced diversity maximization on the left and class-aware minimization on the right. Solid lines represent forward propagation, while dashed lines represent backward propagation. Maximization includes unknown and out-of-distribution maximization. (a) Unknown maximization, guided by unknown loss L_{uk} , semantic consistency loss L_{ct} , and diversified loss L_{div} , synthesizes unknown samples without domain shift x_{su} , which maintains correlation with the corresponding class of the source sample used for synthesizing them. (b) Out-of-distribution maximization appends x_{su} to x_s and synthesizes out-of-distribution samples x_g , to simulate target domain that has both known and unknown classes with open-set loss L_{os} and semantic consistency loss L_{ct} . (c) Minimization improves generalization using both source samples x_s and augmented samples x_s with open-set loss L_{os} , inter-sample loss L_{inter} and intra-sample loss L_{intra} . It explores the generalized class correlation and improves generalization performance. Maximization and minimization are performed alternatively.

173 to iteratively learn “hard” data points from fictitious target
174 distributions to learn generalization ability by addressing:

$$\min_{\theta} \sup_{D_t} \{ \mathbb{E}[L_{ce}(\theta; D_t)] : d(D_t, D_s) \leq \rho \}, \quad (1)$$

175 d is used to quantify the similarity between D_s and D_t . ρ
176 corresponds to the maximum allowable cost to transfer D_s to D_t . Meanwhile,
177 θ is the parameters of the model, which are refined through cross-entropy loss L_{ce} . Adversarial Data Augmentation
178 expands the source domain’s diversity by creating a domain D_g to approximate the unseen target domains D_t .
179 It redefines the worst-case problem (Eq.1) into a Lagrangian optimization problem with a fixed penalty parameter γ :

$$\min_{\theta} \sup_{D_g} \{ \mathbb{E}[L_{ce}(\theta; D_g)] - \gamma d(D_g, D_s) \}, \quad (2)$$

183 where $d(D_g, D_s) = L_{ct}(\theta_g; D_g, D_s) = \|F(x_g) - F(x_s)\|_2^2 + \infty \cdot \mathbf{1}_{\{y \neq y_s\}}$, F is the feature extractor. $\mathbf{1}\{\cdot\}$ is the 0-1
184 indicator function. It ensures consistent semantics between the
185 synthesized samples and their corresponding original samples
186 used for synthesis. The loss function is:
187

$$L_{ada}(\theta, \theta_g; D_s, D_g) = L_{ce}(\theta; D_g) - \gamma L_{ct}(\theta_g; D_g, D_s), \quad (3)$$

188 the training is processed iteratively between two phases:
189 the maximization and the minimization. During maximization,
190 a fictitious target domain D_g is synthesized from D_s by
191 optimizing L_{ada} with learning rate $\eta \geq 0$:

$$x_g \leftarrow x_s + \eta \nabla_{x_s} L_{ada}(\theta; x_s, x_g), \quad (4)$$

192 After the maximization phase, the synthesized domain D_g is
193 appended to D_s . In the minimization phase, θ is optimized by
194 minimizing L_{ce} with the updated D_s .

195 A. Enhanced Diversity Maximization

196 CrossMatch [51] employs an additional classifier to syn-
197 theses unknown samples during above training process to

198 address OS-SDG. In reality, the boundaries between known
199 and unknown classes are artificially defined, with many cat-
200 egories near the boundary having significant correlations.
201 But the unknown classes synthesized by CrossMatch exhibit
202 significant semantic differences from the known classes and
203 do not encompass these boundary unknown classes, leading to
204 poor performance for these classes (Fig. 1 (b)). To address it,
205 during the maximization, we synthesize unknown classes to
206 display correlations with known classes to cover the unknown
207 classes at the classification boundary for better open-set dis-
208 crimination ability. Samples synthesized by maximization need
209 to address both label space and domain shifts simultaneously
210 from the source domain. Enhanced Diversity Maximization
211 is divided into unknown maximization and out-of-distribution
212 maximization to address them respectively.

213 *1) Unknown Maximization:* Unknown maximization syn-
214 theses new unknown classes with different semantics while
215 avoiding domain shift. Traditional classification includes a
216 feature extractor F , and a known class classifier C_k . But
217 source supervised learning cannot distinguish unseen unknown
218 samples. Inspired by [50], we introduce a parallel unknown
219 class classifiers C_u to identify unknown classes. C_u is a linear
220 layer trained from sketch with K -head, and it shares the same
221 feature extractor F as the known classifier C_k . C_u and C_k are
222 used for two objectives simultaneously: classify known classes
223 and distinguish novel samples into a “unknown” classes. For
224 the first objective, augmented features $f(x_s)$ is defined as:

$$f(x_s) = [C_k(F(x_s)), \max(C_u(F(x_s)))], \quad (5)$$

225 where $[., .]$ is concatenation operation, $f(x_s)$ is applied
226 to classify known classes with cross-entropy loss L_{ce} . For
227 the second objective, we divide unknown samples into two
228 groups: unrelated to known classes and related to known
229 classes. Unknown samples unrelated to the known classes are
230 distinguished by Eq.11, details are in Sec.III-A2 and III-B. For

unknown samples related to the known classes, we synthesize unknown samples x_{su} without introducing domain shift. It simulates multiple unknown classes at the classification boundary between known and unknown classes using the unknown loss L_{uk} , which is defined as follows:

$$L_{uk}(\theta_{su}, \theta; D_{su}, D_s) = L_{ce}(\mathbf{f}(x_{su}), \mathbf{y}_{K+1}) + L_{ce}(\mathbf{f}^*(x_{su}), K+1), \quad (6)$$

where $\mathbf{y}_{K+1} = (1-\alpha)\mathbb{1}_{K+1} + \alpha/(K+1)$ is the smooth label, K is the number of known classes, α is set as 0.1. $\mathbf{f}^*(x_s)$ is the masked ground truth position features, and is defined as:

$$\mathbf{f}^*(x_s) = [C_k(F(x_s)) \circ (\mathbf{1}_K - \mathbf{1}_{y_s}), \max(C_u(F(x_s)))], \quad (7)$$

where $\mathbf{1}_K$ is the K -dimension all-one vector and $\mathbf{1}_{y_s}$ is K -dimension one-hot encoding with only the y_s -th element as 1. \circ is element-wise product. Eq.6 ensures that the synthesized classes belong to a “unknown” cluster. Meanwhile, since the real unknown samples usually have correlations with known classes instead of being completely unrelated (i.e., “computer” and “laptop” in Fig.1(a)). we further introduce diversified loss L_{div} , which ensures the correlation between the synthesized sample and the corresponding source sample used for synthesis. L_{div} is denoted as:

$$L_{div}(\theta_{su}; D_{su}, D_s) = L_{ce}(C_u(F(x_{su})), y_s). \quad (8)$$

To summary, the overall loss of unknown maximization is:

$$L_{um}(\theta, \theta_{su}; D_{su}, D_s) = L_{ct}(\theta_{su}; D_{su}, D_s) - \lambda L_{div}(\theta_{su}; D_{su}, D_s) - L_{uk}(\theta_{su}, \theta; D_{su}, D_s), \quad (9)$$

where λ is the hyper-parameter. By maximizing L_{um} , we synthesize multiple unknown samples with diverse semantics while maintaining some correlation with known classes. This ensures that the synthesized unknown samples belong to different classes, reflecting the actual diversity and interrelation of classes near the known and unknown classification boundary. Similar to Eq.3, We also transform Eq.9 into a Lagrangian relaxation to synthesize unknown classes. Source unknown samples x_{su} are synthesized through the following process:

$$x_{su} \leftarrow x_s + \eta \nabla_{x_s} \{L_{ct}(\theta_{su}; x_{su}, x_s) - \lambda L_{div}(\theta_{su}; x_{su}, y_s) - L_{uk}(\theta_{su}, \theta; x_{su}, x_s)\}, \quad (10)$$

After unknown maximization, synthesized unknown classes in D_{su} have distinct semantic yet without domain shift, and they are appended to D_s for the next stage.

2) *Out-of-distribution Maximization*: Unknown maximization synthesizes multiple unknown classes with novel semantic to simulate the label space shift. But the target domains are under both label space shift and distribution shift. Out-of-distribution maximization follows unknown maximization, making synthesized samples resemble out-of-distribution samples from unseen targets. We consider the synthesized unknown samples x_{su} as the $(K+1)$ -th class and transform the OS-SDG problem into a closed-set Single Domain Generalization problem that contains $K+1$ classes.

Additionally, since a sample is an open-set sample for all classes except its ground truth, we can distinguish unknown classes that are unrelated to both known and synthesized unknown classes by training with x_s and x_{su} using this relationship. For the remaining classes except for its ground truth, $\mathbf{f}^*(x_s)$ is treated as the features of an unknown sample. Using $\mathbf{f}^*(x_s)$ and $\mathbf{f}(x_s)$, we can distinguish semantically unrelated unknown classes by open-set loss L_{os} :

$$L_{os}(x_s, y_s) = L_{ce}(\mathbf{f}(x_s), \mathbf{y}_{y_s}) + L_{ce}(\mathbf{f}^*(x_s), K+1), \quad (11)$$

The first item differentiates known classes and semantically related unknown classes, while the second aligns the masked augmented prediction with the unknown class, enhancing discrimination against semantically unrelated unknown samples without accessing target domains. Since we already possess samples both within and outside the source label space, a natural idea is to synthesize samples that deviate from the source distribution across these spaces. By maximizing Eq. 11 we can generate synthetic classes that incorporate domain shifts at the classification boundary during training, thus enlarging the margins between known and unknown classes. Therefore, we aim to create fictitious target classes by maximizing out-of-distribution characteristics, with our overall objective being:

$$L_{om}(\theta, \theta_g; D_s, D_g) = -L_{ct}(\theta_g; x_g, x_s) + L_{os}(\theta_g; x_g, x_s), \quad (12)$$

where x_g are synthesized fictitious target samples, and $D_g = \{x_g^i\}_{i=1}^{n_g}$ contains n_g samples. The first term guarantees the synthesized samples in D_g maintain semantic connections with the samples in D_s used for their synthesis, the last two terms synthesize samples that aim to confuse the classifier to simulate the domain shift, thereby synthesizing known and unknown samples with domain shift. We apply lagrangian relaxation to synthesize samples as:

$$x_g \leftarrow x_s + \eta \nabla_{x_g} \{L_{os}(\theta; x_s, y_s) - L_{ct}(\theta_g; x_g, x_s)\}, \quad (13)$$

where D_g contains both known and unknown classes with distribution gap. D_g are appended to D_s after this stage.

B. Class-Aware Minimization

Minimization occurs both during and after maximization. During maximization, minimization focuses solely on supervised training, ensuring that C_k and C_k achieve sufficient discrimination to guide the maximization in generating desired samples. After maximization, the model needs adapt to synthesized known and unknown classes to enhance its cross-domain open-set discrimination ability. However, as the target unknown samples are inaccessible and their categories are agnostic during training, the inherent data structure among unknown classes cannot be fully explored. Furthermore, unknown classes near the classification boundary often correlate with nearby known classes. Learning these connections between known and unknown samples is essential.

1) *Learning Inter-Sample Correlation*: Although target unknown classes and their number are both agnostic during training, the synthesized known and unknown samples with diverse semantic and domain shift are sufficient for simulating the target domain. Given that all synthesized unknown samples exhibit some correlation with the source known classes used in their synthesis, we introduce a class-agnostic clustering method. This approach utilizes the similarity between synthesized unknown samples and K known classes to cluster the synthesized unknown samples into K clusters, effectively maximize the class margin between known and unknown classes along the classification boundary. we set the number of clusters for unknown samples to be K , and argue the k -th cluster corresponds to the k -th class in the known label space. Compared with other known classes, the k -th unknown cluster exhibits the highest similarity to the k -th known class. Meanwhile, unknown samples that are semantically unrelated to the known classes can be distinguished by minimizing Eq.11.

These two strategies simulates class correlations among target unknown classes without requiring knowing their quantity.

Our class-agnostic clustering includes inter-sample clustering and intra-sample clustering. Inspired by [43], inter-sample clustering first concatenate the outputs of C_k and C_u , and define $\mathbf{z}(x_s^i) = [C_k(F(x_s^i)), C_u(F(x_s^i))]$ for clustering with $2K$ heads. Two output vectors, $\mathbf{z}(x_s^i)$ and $\mathbf{z}(x_s^j)$, are chosen to compute cosine similarity matrix $S(x_s^i, x_s^j)$ for unlabeled input samples x_s^i and x_s^j . Then pseudo inter-sample label $h(x(i), x(j))$ can be obtained by setting a threshold μ on $S(x_s^i, x_s^j)$ ($\mu=0.9$ by default), where

$$h(x_s^i, x_s^j) = \begin{cases} 1, & \text{if } S(x_s^i, x_s^j) > \mu \\ 0, & \text{otherwise} \end{cases}, \quad (14)$$

Given pseudo inter-sample label $h(x_s^i, x_s^j)$, we employ binary cross entropy loss to train the model as follows:

$$\begin{aligned} L_{\text{inter}}(x_s^i, \mathbf{x}_s^j) = & \sum_{x_s^i, \mathbf{x}_s^j \in D_s} (L_{\text{bce}}(\mathbf{S}(x_s^i, x_s^j), h(x_s^i, x_s^j)) \\ & + L_{\text{bce}}(\mathbf{S}(\mathbf{x}_s^j, x_s^j), h(x_s^j, x_s^j))), \end{aligned} \quad (15)$$

here, we denote x_s^i and \mathbf{x}_s^i as the original and augmented input samples, respectively. Specifically, the augmentation of \mathbf{x}_s^i is achieved by applying randaugment, color jitter, and random horizontal flip to the samples, while x_s^i is obtained by only applying random horizontal flip to the samples, just like other samples during the training process. By learning inter-sample correlation across different augmentations, it can identify more robust and generalizable class-related characteristics.

2) *Learning Intra-Sample Invariance*: The inter-sample clustering explores generalizable category relationships, but it fails to fully exploit the invariance inherent within individual samples. Following [14], we address this by maximizing mutual information between augmented versions of the same sample for intra-sample clustering, denoted as:

$$L_{\text{intra}}(p(x_s^i), p(\mathbf{x}_s^i)) = \sum_{m=1}^{2K} \sum_{n=1}^{2K} P(x_s^i, \mathbf{x}_s^i) \log \left[\frac{P(x_s^i, \mathbf{x}_s^i)}{P(x_s^i)P(\mathbf{x}_s^i)} \right], \quad (16)$$

here, $p(x_s^i) = \text{Softmax}(z(x_s^i))$, and $P(x_s^i, \mathbf{x}_s^i)$ is the joint probability matrix: $P(x_s^i, \mathbf{x}_s^i) = \frac{1}{2K} \sum_{i=1}^{2K} p(x_i) p(\mathbf{x}_s^i)^T$, where m and n represent the m -th row and n -th column, respectively. The formula is optimized using both original samples x_s^i and their augmentations \mathbf{x}_s^i . This improves model generalization and produces more uniformly distributed clusters. Additionally, we minimize Eq.11 to enhance discrimination for unknown classes to the related known classes. The loss for class-aware minimization is:

$$\begin{aligned} L_{\text{min}}(x_s^i, \mathbf{x}_s^j, y_s^i) = & \sum_{x_s^i, \mathbf{x}_s^j, y_s^i \in D_s} (\sigma(L_{\text{inter}}(x_s^i, \mathbf{x}_s^j) \\ & - L_{\text{intra}}(x_s^i, \mathbf{x}_s^i)) + L_{\text{os}}(x_s^i, y_s^i)), \end{aligned} \quad (17)$$

σ is a trade-off and is set as 0.1. By minimizing Eq.17, CADA can classify both generalized seen and unseen samples well and explore the correlations among categories.

IV. THEORETICAL ANALYSIS

The source domain is defined as D_s , they are used to synthesize unknown samples D_{su} . Define $L_e: F \times F \rightarrow \mathbb{R}^+ \cup \{\infty\}$ as the cost of perturbing embedding x_s to x_{su} in embedding space. $L_i: X \times X \rightarrow \mathbb{R}^+ \cup \{\infty\}$ is the cost of perturbing x_s to x_{su} in input space. The distance between D_s and D_{su} in embedding space

is $d_{L_e}(\theta_{su}; D_s, D_{su}) := \inf_{M_F \in \Pi(D_s, D_{su})} \mathbb{E}_{M_F}[L_e(x_s, x_{su})]$, similarly, the distance between D_s and D_{su} in input space is $d_{L_i}(\theta_{su}; D_s, D_{su}) := \inf_{M_x \in \Pi(D_s, D_{su})} \mathbb{E}_{M_x}[L_i(x_s, x_{su})]$. M_F and M_x are measures in the embedding and input space respectively. $\Pi(D_s, D_{su})$ is joint distribution of D_s and D_{su} .

Analysis on unknown maximization. Unknown maximization synthesizes unknown samples D_{su} without introducing domain shift from D_s . D_{su} exhibits semantic differences from the D_s in embedding space, but possesses strong correlations with D_s in input space. Consequently, the goal of the worst-case problem for the unknown maximization stage becomes to synthesize the most challenge samples on boundary between known and unknown classes, which makes the model confused. The relaxed worst-case problem can be rewritten as:

$$\theta^* = \max_{\theta} \inf_{D_{su}} \mathbb{E}(L_{\text{task}}(\theta; D_{su})) = -\min_{\theta} \sup_{D_{su}} [-\mathbb{E}(L_{\text{task}}(\theta; D_{su}))], \quad (18)$$

where θ is parameters of the model. In unknown maximization, L_{task} is L_{uk} . D_{su} synthesizes unknown samples without domain shift. D_{su} exhibits both semantic differences from D_s in the semantic space, and correlations with the corresponding class D_s used to synthesize D_{su} in input space. Hence, $\{D_{su}: d_{L_i}(\theta_{su}; D_s, D_{su}) \leq \rho, d_{L_e}(\theta_{su}; D_s, D_{su}) \geq \eta\}$, we use Lagrangian relaxation under constraint with fixed parameters $\lambda \geq 0$ and $\beta \geq 0$ to solve Eq. 18 as follows:

$$\begin{aligned} \theta^* = & -\min_{\theta} \sup_{D_{su}} \{\mathbb{E}[-L_{uk}(\theta; D_{su})] - \lambda[W_{L_i}(\theta_{su}; D_{su}, D)] \\ & + \beta[W_{L_e}(\theta_{su}; D_{su}, D_s)]\} = -\min_{\theta} \{\mathbb{E}[\delta_{\lambda, \beta}(\theta_{su}, \theta; x_{su}, x_s)]\}, \end{aligned} \quad (19)$$

where $\delta_{\lambda, \beta}(\theta_{su}, \theta; x_{su}, x_s) = \sup_{x_{su}} \{-L_{uk}(\theta_{su}, \theta; x_{su}, x_s) - \lambda W_{L_i}(\theta_{su}; x_{su}, x_s) + \beta W_{L_e}(\theta_{su}; x_{su}, x_s)\}$. The problem reduces to minimizing $\delta_{\lambda, \beta}$. As shown in [34], $\delta_{\lambda, \beta}$ is smooth with respect to θ if λ and β are large enough and the Lipschitz smoothness assumption holds. Gradient can be computed as:

$$\nabla_{\theta} \delta_{\lambda, \beta}(\theta; x_s) = \nabla_{\theta} \{\mathbb{E}[-L_{uk}(\theta; x_s^*(x_s, \theta))]\}, \quad (20)$$

where $x_s^*(x_s, \theta) = \text{argmax}_{x_{su}} [\beta W_{L_e}(\theta_{su}; x_{su}, x_s) - \lambda W_{L_i}(\theta_{su}; x_{su}, x_s) - L_{uk}(\theta_{su}, \theta; x_{su}, x_s)] = \text{argmax}_{x_{su}} [L_{um}(\theta_{su}, \theta; x_{su}, x_s)]$, which is exactly unknown maximization in Eq. 9.

Analysis on out-of-distribution maximization. unknown samples within the source distribution belonging to D_{su} are seen as the $(K+1)$ -th class and are append to D_s . Open-Set Single Domain generalization problem becomes Single Domain generalization problem with $K+1$ classes.

Similar to unknown maximization, out-of-distribution maximization synthesizes new samples for better generalization ability. But out-of-distribution maximization synthesizes new samples D_g that out of the source distribution with D_s . D_g exhibits domain shift from D_s in the input space, but possesses strong semantic similarity with D_s in semantic space. Consequently, the synthesised unknown samples encourage classifier to distinguish samples from both D_g and D_s well even if the sample selected from . The relaxed worst-case problem can be rewritten as:

$$\theta^* = \min_{\theta} \sup_{D_g} \mathbb{E}(L_{\text{task}}(\theta; D_g)), \quad (21)$$

TABLE II: Classification accuracy (%) on *Digits* with LeNet-5, on *Office31* and *VisDA2017* with ResNet-18. Best are in **bold**.

Methods	Access to D_t	Type	MNIST \rightarrow Others			Amazon \rightarrow Others			Synthesis \rightarrow Real World			Real World \rightarrow Synthesis			Average		
			acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u
OSDAP[31] CombEmb[16]	Accessible	OSDA	57.3 \pm 0.3	54.2 \pm 0.4	49.6 \pm 0.2	76.1 \pm 0.3	79.6 \pm 0.2	84.7 \pm 0.3	54.7 \pm 0.3	59.2 \pm 0.5	69.1 \pm 0.4	68.6 \pm 0.3	50.9 \pm 0.2	38.6 \pm 0.3	61.7 \pm 0.3	55.1 \pm 0.4	53.9 \pm 0.4
	Accessible	OSDA	57.4 \pm 0.4	53.3 \pm 0.5	48.5 \pm 0.4	77.2 \pm 0.5	80.6 \pm 0.3	85.4 \pm 0.4	55.7 \pm 0.4	60.5 \pm 0.4	71.3 \pm 0.5	69.3 \pm 0.4	52.0 \pm 0.4	39.7 \pm 0.3	62.4 \pm 0.4	56.3 \pm 0.4	55.5 \pm 0.4
ERM[17] PROSER[50] ADA[38] MEADA[49] CrossMatch[51]	Inaccessible	SL	49.2 \pm 0.3	18.0 \pm 0.4	13.0 \pm 0.3	77.9 \pm 0.4	40.7 \pm 0.3	21.1 \pm 0.5	43.3 \pm 0.4	30.4 \pm 0.5	23.4 \pm 0.6	64.5 \pm 0.4	25.6 \pm 0.4	15.5 \pm 0.4	53.9 \pm 0.3	28.0 \pm 0.4	19.4 \pm 0.4
	Inaccessible	OSR	49.9 \pm 0.4	41.2 \pm 0.3	33.6 \pm 0.5	74.0 \pm 0.3	45.6 \pm 0.6	32.2 \pm 0.4	44.4 \pm 0.5	40.0 \pm 0.4	35.3 \pm 0.5	68.5 \pm 0.4	32.3 \pm 0.4	20.3 \pm 0.3	56.8 \pm 0.4	36.2 \pm 0.3	27.8 \pm 0.5
CADA	Inaccessible	SDG	50.2 \pm 0.4	20.1 \pm 0.3	15.1 \pm 0.4	77.0 \pm 0.4	37.2 \pm 0.3	24.0 \pm 0.4	44.7 \pm 0.3	31.7 \pm 0.5	24.5 \pm 0.2	67.9 \pm 0.5	29.5 \pm 0.3	18.2 \pm 0.4	56.3 \pm 0.5	30.7 \pm 0.3	21.4 \pm 0.4
	Inaccessible	SDG	52.9 \pm 0.4	30.4 \pm 0.5	29.8 \pm 0.3	76.8 \pm 0.4	35.0 \pm 0.5	22.2 \pm 0.3	45.1 \pm 0.4	33.7 \pm 0.5	25.7 \pm 0.3	68.2 \pm 0.4	28.7 \pm 0.5	17.6 \pm 0.3	56.6 \pm 0.4	31.2 \pm 0.5	21.6 \pm 0.3
CADA	OS-SDG	51.3 \pm 0.4	38.7 \pm 0.5	46.1 \pm 0.3	76.6 \pm 0.4	52.8 \pm 0.4	39.3 \pm 0.5	44.4 \pm 0.3	42.5 \pm 0.2	40.8 \pm 0.5	69.5 \pm 0.4	53.0 \pm 0.5	40.9 \pm 0.5	57.0 \pm 0.3	47.8 \pm 0.5	40.8 \pm 0.4	
	Inaccessible	Ours	51.0 \pm 0.4	49.5\pm0.3	47.2\pm0.4	76.0 \pm 0.5	57.7\pm0.5	45.4\pm0.4	44.6 \pm 0.3	45.5\pm0.5	47.0\pm0.3	69.4 \pm 0.5	56.2\pm0.3	45.1\pm0.4	57.0 \pm 0.5	50.8\pm0.3	46.1\pm0.4

TABLE III: Classification accuracy (%) on *Office-Home* with ResNet-18.

Methods	Access to D_t	Type	Art \rightarrow Others			Clipart \rightarrow Others			Product \rightarrow Others			Real World \rightarrow Others			Average		
			acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u
OSDAP[31] CombEmb[16]	Accessible	OSDA	45.6 \pm 0.3	43.4 \pm 0.4	49.8 \pm 0.5	52.8 \pm 0.4	52.3 \pm 0.3	51.7 \pm 0.4	41.4 \pm 0.5	46.6 \pm 0.3	54.9 \pm 0.4	53.5 \pm 0.3	47.1 \pm 0.4	41.6 \pm 0.4	48.3 \pm 0.5	48.4 \pm 0.3	49.5 \pm 0.4
	Accessible	OSDA	46.5 \pm 0.5	47.7 \pm 0.4	49.2 \pm 0.5	51.9 \pm 0.4	52.7 \pm 0.5	53.7 \pm 0.4	54.2 \pm 0.4	55.8 \pm 0.3	55.3 \pm 0.3	66.3 \pm 0.4	59.1 \pm 0.5	52.8 \pm 0.4	54.7 \pm 0.4	53.8 \pm 0.3	52.8 \pm 0.4
ERM[17] PROSER[50] ADA[38] MEADA[49] CrossMatch[51]	Inaccessible	SL	65.0 \pm 0.4	31.1 \pm 0.5	20.5 \pm 0.6	64.1 \pm 0.4	35.8 \pm 0.6	24.7 \pm 0.5	60.5 \pm 0.4	36.3 \pm 0.4	26.3 \pm 0.5	66.6 \pm 0.4	33.9 \pm 0.5	23.2 \pm 0.5	64.1 \pm 0.4	34.3 \pm 0.4	23.7 \pm 0.5
	Inaccessible	OSR	62.6 \pm 0.4	47.5 \pm 0.3	37.6 \pm 0.5	60.0 \pm 0.5	38.8 \pm 0.3	28.2 \pm 0.4	59.9 \pm 0.5	42.1 \pm 0.5	31.9 \pm 0.5	64.9 \pm 0.5	49.7 \pm 0.6	39.7 \pm 0.4	61.9 \pm 0.4	44.5 \pm 0.5	34.3 \pm 0.4
CADA	Inaccessible	SDG	68.3 \pm 0.4	32.9 \pm 0.5	22.1 \pm 0.4	65.1 \pm 0.4	42.1 \pm 0.4	31.2 \pm 0.3	60.5 \pm 0.5	34.7 \pm 0.4	24.6 \pm 0.5	67.1 \pm 0.3	34.9 \pm 0.3	22.9 \pm 0.5	65.2 \pm 0.4	36.2 \pm 0.3	25.4 \pm 0.4
	Inaccessible	SDG	68.3 \pm 0.6	33.3 \pm 0.5	22.4 \pm 0.4	65.3 \pm 0.5	42.1 \pm 0.4	31.3 \pm 0.5	60.4 \pm 0.4	35.7 \pm 0.4	25.6 \pm 0.5	67.0 \pm 0.4	34.7 \pm 0.4	23.7 \pm 0.5	65.0 \pm 0.3	36.4 \pm 0.4	25.7 \pm 0.5
CADA	OS-SDG	Inaccessible	65.9 \pm 0.4	53.2 \pm 0.3	45.3 \pm 0.5	62.9 \pm 0.4	48.9 \pm 0.5	37.8 \pm 0.4	58.4 \pm 0.3	45.3 \pm 0.5	37.7 \pm 0.6	67.1 \pm 0.4	50.8 \pm 0.4	41.3 \pm 0.3	63.6 \pm 0.5	49.6 \pm 0.5	40.5 \pm 0.3
	Ours	Inaccessible	65.8 \pm 0.4	54.9\pm0.3	46.4\pm0.5	59.0 \pm 0.3	53.4\pm0.3	48.3\pm0.4	60.1 \pm 0.3	52.4\pm0.4	45.8\pm0.4	65.0 \pm 0.3	56.3\pm0.4	49.1\pm0.5	62.5 \pm 0.4	54.2\pm0.3	47.4\pm0.4

TABLE IV: Classification accuracy (%) on *PACS* with ResNet-18.

Methods	Access to D_t	Type	Art Paint \rightarrow Others			Cartoon \rightarrow Others			Photo \rightarrow Others			Sketch \rightarrow Others			Average		
			acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u
OSDAP[31] CombEmb[16]	Accessible	OSDA	36.5 \pm 0.3	36.3 \pm 0.4	76.6 \pm 0.5	32.7 \pm 0.4	32.3 \pm 0.4	67.1 \pm 0.4	36.4 \pm 0.5	36.0 \pm 0.4	74.3 \pm 0.3	34.3 \pm 0.4	33.8 \pm 0.5	32.9 \pm 0.3	35.0 \pm 0.3	34.6 \pm 0.4	62.7 \pm 0.5
	Accessible	OSDA	48.4 \pm 0.4	53.5 \pm 0.5	72.2 \pm 0.4	52.3 \pm 0.3	56.4 \pm 0.5	67.4 \pm 0.5	42.6 \pm 0.3	49.1 \pm 0.3	75.1 \pm 0.4	46.4 \pm 0.3	50.3 \pm 0.4	61.3 \pm 0.4	47.4 \pm 0.3	52.3 \pm 0.4	69.0 \pm 0.4
ERM[17] PROSER[50] ADA[38] MEADA[49] CrossMatch[51]	Inaccessible	SL	46.2 \pm 0.4	32.7 \pm 0.3	23.5 \pm 0.5	49.0 \pm 0.4	23.5 \pm 0.4	15.0 \pm 0.4	34.7 \pm 0.4	26.7 \pm 0.5	20.2 \pm 0.3	41.3 \pm 0.4	30.7 \pm 0.4	22.7 \pm 0.5	42.8 \pm 0.4	28.4 \pm 0.3	20.3 \pm 0.3
	Inaccessible	OSR	48.6 \pm 0.4	39.4 \pm 0.5	30.8 \pm 0.5	49.3 \pm 0.3	35.0 \pm 0.3	25.2 \pm 0.4	36.5 \pm 0.5	33.8 \pm 0.3	30.3 \pm 0.5	42.7 \pm 0.6	41.6 \pm 0.4	39.8 \pm 0.3	44.3 \pm 0.5	37.4 \pm 0.4	31.4 \pm 0.5
CADA	Inaccessible	SDG	47.7 \pm 0.5	35.9 \pm 0.5	26.8 \pm 0.4	49.3 \pm 0.4	26.8 \pm 0.4	17.2 \pm 0.5	37.0 \pm 0.4	29.2 \pm 0.5	22.4 \pm 0.4	43.2 \pm 0.5	33.7 \pm 0.4	25.7 \pm 0.3	44.3 \pm 0.4	31.4 \pm 0.5	22.9 \pm 0.4
	Inaccessible	SDG	48.1 \pm 0.4	36.5 \pm 0.5	27.4 \pm 0.3	48.8 \pm 0.4	27.6 \pm 0.5	18.0 \pm 0.3	36.9 \pm 0.3	28.7 \pm 0.5	21.8 \pm 0.3	42.8 \pm 0.4	32.6 \pm 0.5	24.4 \pm 0.3	44.1 \pm 0.4	31.3 \pm 0.4	22.9 \pm 0.5
CADA	OS-SDG	Inaccessible	48.5 \pm 0.4	41.0 \pm 0.5	33.2 \pm 0.5	49.6 \pm 0.3	37.5 \pm 0.5	30.1 \pm 0.4	34.2 \pm 0.3	33.0 \pm 0.4	31.1 \pm 0.5	43.2 \pm 0.3	45.1 \pm 0.5	30.6 \pm 0.3	43.9 \pm 0.4	39.2 \pm 0.3	36.2 \pm 0.5
	Ours	Inaccessible	49.2 \pm 0.4	43.0\pm0.5	35.8\pm0.3	49.9 \pm 0.5	42.2\pm0.4	34.2\pm0.3	34.3 \pm 0.4	35.1\pm0.4	37.2\pm0.4	42.1 \pm 0.4	45.0 \pm 0.3	57.2\pm0.5	43.9 \pm 0.3	41.3\pm0.3	41.1\pm0.4

where θ is parameters of the model. In out-of-distribution maximization, L_{task} is L_{os} . D_g generates both unknown and known samples out of the source distribution. D_{su} not only exhibits semantic differences from D_s in the input space, but also maintains semantic correlations with D_s in the semantic space. Hence, $\{D_g : W_{L_i}(D_s, D_g) \geq \rho, W_{L_e}(D_s, D_g) \leq \eta\}$. It is hard to solve Eq. 21, so we use Lagrangian relaxation under constraint with fixed parameters $\lambda \geq 0$ and $\beta \geq 0$:

$$\theta^* = \min_{\theta} \sup_{D_g} \{ \mathbb{E}[L_{os}(\theta; D_g)] - \lambda [W_{L_i}(\theta_{s_g}; D_s, D_g) - \beta W_{L_e}(\theta_{s_g}; D_s, D_g)] \} \quad (22)$$

where $x_s^*(\theta; \theta_{s_g}) = \text{argmax}_{x_g} [L_{os}(\theta, \theta_{s_g}; x_s, x_g) - \lambda W_{L_i}(\theta_{s_g}; x_s, x_g) - \beta W_{L_e}(\theta_{s_g}; x_s, x_g)]$. The problem becomes minimizing $\delta_{\lambda, \beta}$. $\delta_{\lambda, \beta}$ is smooth w.r.t. θ if λ, β are large enough and the assumption of Lipschitzian smoothness holds. The gradient is computed as:

$$\nabla_{\theta} \delta_{\lambda, \beta}(\theta; x_s) = \nabla_{\theta} \{ \mathbb{E}[L_{os}(\theta; x_s^*(\theta; \theta_{s_g}), x_g)] \}, \quad (23)$$

where $x_s^*(\theta; \theta_{s_g}) = \text{argmax}_{x_g} [L_{os}(\theta, \theta_{s_g}; x_s, x_g) - \lambda W_{L_i}(\theta_{s_g}; x_s, x_g) - \beta W_{L_e}(\theta_{s_g}; x_s, x_g)]$. The target unknown label space \mathcal{C}_t^u consists of unknown categories ranging from 5 to 9.

TABLE V: Experiments on Face Anti-Spoofing

Method	C to H		I to H	
	HTER (%) ↓	AUC (%) ↑	HTER (%) ↓	AUC (%) ↑
PatchCNN [42]	39.54	64.54	35.03	73.24
PDEN [22]	35.76	69.10	31.03	74.35
LD [42]	38.33	65.12	32.43	73.57
PCGRL [15]	27.24	78.81	28.03	79.37
Ours	31.45	72.36	30.55	75.83

485 bottleneck layer and employ a weight normalization layer in
486 the last FC layer. Mini-batch SGD is adopted with momentum
487 0.9 and weight decay 1e-3, and set the learning rate as 1e-4.
488 Batch size is 64. We use cropping [40] and RandAugment
489 [5] to augment the training data x_s to \mathbf{x}_s . CADA compares
490 with other methods following their original procedures under
491 our settings except DA method accessing target data during
492 training. Time and space complexity are both O(N).

493 *4) Comparative Evaluations and Performance Metrics:*
494 CADA is compared against various baseline and SOTA
495 methods includes supervised learning (SL) method Empirical
496 Risk Minimization (ERM) [17], Open-Set DA method
497 Open-Set Domain Adaptation by Back propagation (OSDAP)
498 [31], Open-Set recognition (OSR) method PROSER [50] and
499 CombEmb [16], state-of-the-art SDG methods Adversarial
500 Data Augmentation (ADA) [38] and Maximum-Entropy Ad-
501 versarial Data Augmentation (MEADA) [49], and state-of-
502 the-art OS-SDG method CrossMatch (CrossMatch) [51] on
503 five datasets. In this paper, we use acc, hs and acc_u as our
504 three main metrics. Since in our setting, almost half data
505 are unknown samples, while per-class accuracy (acc) is the
506 mean accuracy averaged over all classes in all $K + 1$ classes,
507 unknown samples are seen as the $K + 1$ class. Therefore,
508 treating unknown samples, which account for over half of the
509 dataset, equally with the rest of the known classes is unfair.
510 While h-score (hs), $hs = \frac{2*acc_u*acc_k}{acc_u+acc_k}$, is the harmonic mean
511 of average per-class accuracy in known and unknown space,
512 which give equal importance weight to known and unknown
513 classes. This setting is more consistent with our practical
514 situation, as the hs score will only be high when both known
515 and unknown accuracies are high. Here, acc_k and acc_u are
516 per-class accuracy of known and unknown label space.

517 B. Results and Further Analysis

518 *1) Experiment Result:* Table II, III and IV reports the per-
519 formance on the Digits, Office31, VisDA-2017, Office-Home
520 and PACS datasets under the OS-SDG setting, respectively.
521 OSDA performs well on acc_u . However, it uses both source
522 and target data in training. Thanks to L_{os} , PROSER can distin-
523 guish some semantic unrelated unknown samples, but fails to
524 classify target known ones well due to distribution shift. SDG
525 methods are effective in improving accuracy on known classes.
526 But they struggle on unknown classes. This leads to low
527 accuracy on unknown classes and low hs . The SOTA OS-SDG
528 method, CrossMatch, improves the discrimination of unknown
529 samples by synthesizing them, but its assumption that these
530 samples are entirely unrelated to known classes limits its open-
531 set discriminative power. Our CADA outperforms CrossMatch
532 significantly in hs and acc_u through all the datasets. It shows
533 that CADA significantly improves generalization over existing
534 methods on unknown classes while maintaining accuracy on
535 known classes. More analysis is in supplemental materials.

TABLE VII: Comparison on Office-Home and VisDA2017

Methods	Venue	Office-Home			VisDA2017		
		acc	h_s	acc_u	acc	h_s	acc_u
NormAuG+PROSER [28, 50]	TIP24	60.5	51.2	43.8	45.3	40.3	40.8
MEDIC [1]	ICCV23	58.2	44.8	35.8	53.3	45.1	37.6
ODG-CLIP [33]	CVPR24	61.3	53.0	46.1	54.7	48.9	42.9
CADA	Ours	62.5	54.2	47.4	57.0	50.8	49.1

536 *2) Component Effectiveness Evaluation:* We conduct an
537 ablation study in Table VI to evaluate L_{div} , L_{uk} , L_{intra} ,
538 and L_{inter} . The first row is the baseline, a CADA variant
539 using only L_{os} without these four components. Although it
540 retains strong discriminative power for known classes, its
541 acc_u and h_s remain relatively low, indicating limited ability
542 to handle unknown classes. Comparing the first and last
543 rows shows that while L_{os} can only distinguish unrelated
544 unknown samples, the full CADA model can differentiate both
545 related and unrelated unknowns, validating our approach. The
546 difference between the first and second rows highlights that
547 L_{intra} and L_{inter} help capture class correlations, improving
548 performance. The comparison between the fourth and last rows
549 demonstrates that L_{div} enhances the diversity of synthesized
550 unknown samples, contributing to better results. Additionally,
551 from the first, second, and last rows, we observe that L_{uk}
552 strengthens discriminative power for unknown samples by gen-
553 erating diverse, known-related unknowns. The last four rows
554 highlight how L_{intra} and L_{inter} explore class relationships.

555 *3) Applicability and Comparison with Other Methods:*
556 We evaluate CADA’s applicability on anti-spoofing task in
557 Table V. Our method outperforms SDG methods [20, 39, 42],
558 which can be attributed to its strong open-set discrimination
559 capability. However, compared to PCGRL [15], an OS-SDG
560 method specifically designed for this scenario, our method
561 performs less effectively. This is because our CADA primarily
562 focuses on relative relationships between classes without cap-
563 turing fine-grained, attribute-level features. Such limitations
564 hinder performance on face anti-spoofing tasks, which heavily
565 rely on local detail features. Table VII compare CADA
566 with more SOTA approaches. It shows that even SOTA SDG
567 methods, when augmented with open-set recognition modules,
568 experience significant performance drops when target domain
569 data is inaccessible. This highlights both the challenge and
570 the practical significance of our OS-SDG setting. In real-
571 world scenarios, open-set domain generalization methods often
572 assume multiple source domains—an assumption that rarely
573 holds true. When these methods are applied in a single-source
574 setting, their performance deteriorates markedly. Notably, even
575 ODG-CLIP, which leverages the strong generalization ability
576 of CLIP, underperforms compared to our proposed CADA.
577 This further substantiates the superiority of our approach and
578 the challenges inherent in the OS-SDG setting.

579 *4) Visualization:* We conduct t-SNE [36] under the OS-
580 SDG setting with Real World domain → other task, Fig. 3a
581 show the comparison between CrossMatch and CADA. Red
582 points are source known classes, while blue and green are
583 target known and unknown classes. CrossMatch performs bet-
584 ter but sometimes misclassifies unknown samples into known
585 clusters (red circles) due to difficulty in identifying seman-
586 tically related unknowns. In contrast, CADA separates target
587 known samples more accurately and effectively distinguishes

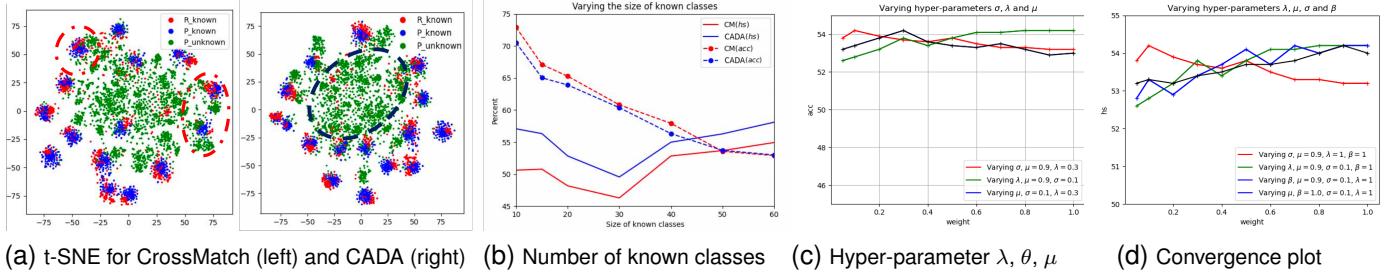


Fig. 3: Ablation study under OS-SDG setting. (a) Feature visualization on “Real World→Others” task. (b) Number of categories known on *Office-Home* dataset. (c) Varying hyper-parameters λ, θ and μ . (d) Convergence plot on *Office-Home* dataset.

TABLE VI: Ablation study of variants with our CADA on *Office-Home* dataset.

CADA's variant				settings on <i>Office-Home</i>														
<i>L</i> _{div}	<i>L</i> _{uk}	<i>L</i> _{intra}	<i>L</i> _{inter}	Art → Others			Clipart → Others			Product → Others			Real World → Others			Average		
				acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u	acc	hs	acc _u
				65.3	43.7	32.3	58.2	45.2	36.5	60.4	42.0	31.7	65.9	47.4	36.4	62.4	44.6	34.2
		✓		65.3	44.9	34.2	58.0	44.2	35.8	60.5	42.4	32.7	66.0	47.6	37.3	62.4	44.8	34.9
	✓			65.1	54.5	46.2	56.6	49.9	44.2	59.5	50.1	42.8	64.9	55.1	47.3	61.5	52.4	45.1
	✓	✓		64.8	53.2	44.5	58.3	52.0	46.5	57.9	50.7	44.6	64.1	56.4	49.7	61.3	53.1	46.3
	✓	✓	✓	64.0	53.1	45.3	58.83	51.84	46.3	58.5	48.8	41.8	63.6	55.8	49.7	61.2	52.4	45.8
	✓	✓	✓	64.8	54.4	46.2	58.3	51.9	46.2	59.6	50.3	42.9	64.7	56.4	49.4	61.9	53.2	46.2
	✓	✓	✓	65.8	54.9	46.4	59.0	53.4	48.3	60.1	52.4	45.8	65.0	56.3	49.1	62.5	54.2	47.4

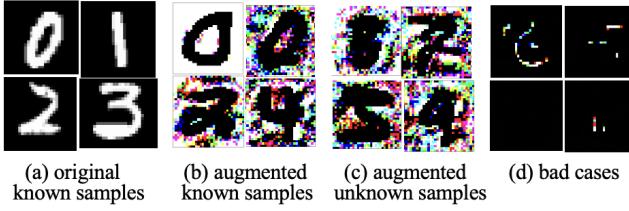


Fig. 4: Visualization of synthesized samples.

both related and unrelated unknown samples (black circle).

5) *Varying number of known classes*: Fig. 3b illustrates the adaptation task from the Real World to other domains, with the number of known classes ranging from 10 to 60. The red curves represent the performance of CrossMatch, while the blue curves correspond to CADA. Both models exhibit sensitivity to the number of known classes, with performance generally declining as the number increases. However, CADA consistently outperforms CrossMatch, demonstrating superior differentiation of unknown samples. This advantage leads to a higher harmonic mean score (hs) in terms of classifiability and separability while maintaining comparable accuracy levels.

6) *Hyper-parameter sensitivity*: In Fig. 3c, we conduct an ablation study on the Office-Home dataset under the OS-SDG setting, exploring the impact of varying hyper-parameters λ , θ , and μ on model performance. The evaluation metric used is hs . The figure shows three curves representing different hyper-parameter variations: the red curve varies σ with $\mu = 0.9$ and $\lambda = 0.3$, the green curve varies λ with $\mu = 0.9$ and $\sigma = 0.1$, and the blue curve varies μ with $\sigma = 0.1$ and $\lambda = 0.3$. The results show that CADA maintains stable performance across different settings, demonstrating its robustness to these hyper-parameters. Specifically, varying σ (red curve) causes mild fluctuations, with a slight peak around a weight of 0.3, suggesting an optimal region for σ . Adjusting λ (green curve) shows an increasing trend, indicating its importance for

TABLE VIII: Impact of different class ratios.

Source	Known	Source	Unknown	OOD	Known	OOD	Unknown	hs
1			1			1		54.9
1			0			1		25.6
1			1			1		36.8
1			0			1		50.8
1			0.5			1		52.4
1			0.5			0.5		51.3

TABLE IX: Accuracy vs. Complexity Comparison

Method	hs	Accuracy (Unknown)	Training Time (hours)	Inference Time (ms)	FLOPs (G)
CrossMatch	47.8%	40.5%		3.8	14.7
CADA (Ours)	50.8%	47.4%	5.2	14.9	2.1 2.2

regularization. Changes in μ (blue curve) result in a relatively stable trend, highlighting the model’s resilience to variations in this parameter. Overall, Fig. 3c shows that CADA is robust to hyper-parameter changes, validating the model’s stability.

7) *Impact of different class ratios*: As shown in Table VIII, the first row represents the full CADA model, achieving the best performance by generating sufficient known and unknown samples to enhance generalization and open-set recognition. The second row, which removes all unknown samples, sharply reduces open-set recognition ability, lowering hs despite strong known-class performance. In the third row, removing only out-of-distribution unknowns slightly improves performance over the second row but remains far below the full model. The fourth row removes source unknowns, causing a moderate performance drop, as Eq.(6) relies on class relationships, but the impact is limited. The last two rows, with adjusted generation ratios, show performance declines, highlighting the importance of each module in our approach.

8) *Accuracy vs complexity comparison*: We added an “Accuracy vs. Complexity” comparison with baseline CrossMatch, as shown in Table IX. Our method improves unknown accuracy by +6.9% and achieves a 3% higher hs compared to CrossMatch. Although the two-stage maximization increases

637 training time by approximately 1.5 hours, inference time and
 638 FLOPs remain similar to the baseline.

639 *9) Iteration of training:* The convergence plot is depicted
 640 in Fig. 3d. For the Office-Home dataset, the training includes
 641 2 stages with 200 epochs. The result exhibits convergence
 642 stability, our training includes two phases. Specifically, the
 643 first phase aims to synthesize the fictitious target domain
 644 to simulate the real target domain. During this phase, we
 645 first utilizes supervised learning with L_{os} to learn open-
 646 set recognition ability. The loss curve converges stably.
 647 Then we conduct unknown maximization and minimization to
 648 synthesize classes with both domain shift and label space shift
 649 to simulate the inaccessible target domain. The loss curve is
 650 dynamically changing. During the second phase, synthesized
 651 samples appends to the source domain, model is modified
 652 with both class-agnostic clustering and modified supervised
 653 learning with L_{os} to learn generalized open-set discrimination
 654 ability, the loss curve converges stably.

655 *10) Synthesized samples:* Fig. 4 displays the synthesized
 656 samples—both known and unknown—from the MNIST
 657 dataset. The synthesized known samples exhibit a distinct
 658 domain gap compared to the original images, yet they suc-
 659 cessfully preserve semantic integrity. In contrast, the syn-
 660 thethesized unknown samples, specifically from digits 0, 2, 3,
 661 and 4, demonstrate both domain and semantic variations
 662 when compared to their original counterparts. Despite these
 663 differences, there remains a significant correlation with the
 664 original samples, indicating that the essential characteristics
 665 of the digits are still recognizable. This balance highlights the
 666 effectiveness of the synthesis process in generating new, yet
 667 contextually related samples that expand the training dataset
 668 while maintaining a connection to the original data.

VI. CONCLUSION

669 This work presents Class-Aware Diversified Augmentation,
 670 which synthesizes more realistic unknown samples semanti-
 671 cally correlated with the source known classes to simulate
 672 unseen target domains for generalization. Theoretical analysis
 673 and experiments on benchmarks show its superiority.

674 **Future Works and Limitations** CADA has some limita-
 675 tions for future research. First, our method requires a two-
 676 stage sample generation. Although it produces more diverse
 677 samples than previous methods, reducing generation steps
 678 and resource consumption remains a challenge. Second, our
 679 method currently constrains semantics through simple inter-
 680 class relations; in the future, inspired by [23, 24], we aim to
 681 incorporate part-whole relational properties to construct finer-
 682 grained inter-class relationships for better synthesis.

REFERENCES

- 685 [1] Shirsha Bose, Ankit Jha, Hitesh Kandala, and Biplob Banerjee.
 686 Beyond boundaries: A novel data-augmentation discourse for
 687 open domain generalization. *Transactions on Machine Learning
 688 Research*, 2023.
- 689 [2] Xinyang Chen, Sinan Wang, Mingsheng Long, and Jianmin
 690 Wang. Transferability vs. discriminability: Batch spectral pe-
 691 nalization for adversarial domain adaptation. In *International
 692 conference on machine learning*, pages 1081–1090. PMLR,
 693 2019.
- 694 [3] Zining Chen, Weiqiu Wang, Zhicheng Zhao, Fei Su, Aidong
 695 Men, and Hongying Meng. Practicalldg: Perturbation distillation
 696 on vision-language models for hybrid domain generalization. In
 697 *Proceedings of the IEEE/CVF Conference on Computer Vision
 698 and Pattern Recognition*, pages 23501–23511, 2024.
- [4] Ivana Chingovska, André Anjos, and Sébastien Marcel. On
 699 the effectiveness of local binary patterns in face anti-spoofing.
 700 In *2012 BIOSIG-proceedings of the international conference of
 701 biometrics special interest group (BIOSIG)*, pages 1–7. IEEE,
 702 2012.
- [5] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le.
 703 Randaugment: Practical automated data augmentation with a
 704 reduced search space. In *Proceedings of the IEEE/CVF con-
 705 ference on computer vision and pattern recognition workshops*,
 706 pages 702–703, 2020.
- [6] Abolfazl Farahani, Sahar Voghieri, Khaled Rasheed, and
 707 Hamid R Arabnia. A brief review of domain adaptation.
 708 *Advances in Data Science and Information Engineering: Pro-
 709 ceedings from IC DATA 2020 and IKE 2020*, pages 877–894,
 710 2021.
- [7] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain
 711 adaptation by backpropagation. In *International conference on
 712 machine learning*, pages 1180–1189. PMLR, 2015.
- [8] Kai Han, An Xiao, Enhua Wu, Jianyuan Guo, Chunjing Xu, and
 713 Yunhe Wang. Transformer in transformer. *Advances in Neural
 714 Information Processing Systems*, 34:15908–15919, 2021.
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
 715 Deep residual learning. *Image Recognition*, 7, 2015.
- [10] Guillaume Heusch, Anjith George, David Geissbühler, Zohreh
 716 Mostaani, and Sébastien Marcel. Deep models and shortwave
 717 infrared information to detect face presentation attacks. *IEEE
 718 Transactions on Biometrics, Behavior, and Identity Science*,
 719 2(4):399–409, 2020.
- [11] Jian Hu, Hongya Tuo, Chao Wang, Lingfeng Qiao, Haowen
 720 Zhong, Junchi Yan, Zhongliang Jing, and Henry Leung. Dis-
 721 criminative partial domain adversarial network. In *Computer
 722 Vision-ECCV 2020: 16th European Conference, Glasgow, UK,
 723 August 23–28, 2020, Proceedings, Part XXVII 16*, pages 632–
 724 648. Springer, 2020.
- [12] Jian Hu, Haowen Zhong, Fei Yang, Shaogang Gong, Guile Wu,
 725 and Junchi Yan. Learning unbiased transferability for domain
 726 adaptation by uncertainty modeling. In *Computer Vision-ECCV
 727 2022: 17th European Conference, Tel Aviv, Israel, October 23–
 728 27, 2022, Proceedings, Part XXXI*, pages 223–241. Springer,
 729 2022.
- [13] Jonathan J. Hull. A database for handwritten text recognition
 730 research. *IEEE Transactions on pattern analysis and machine
 731 intelligence*, 16(5):550–554, 1994.
- [14] Xu Ji, Joao F Henriques, and Andrea Vedaldi. Invariant
 732 information clustering for unsupervised image classification and
 733 segmentation. In *Proceedings of the IEEE/CVF International
 734 Conference on Computer Vision*, pages 9865–9874, 2019.
- [15] Fangling Jiang, Qi Li, Weining Wang, Min Ren, Wei Shen, Bing
 735 Liu, and Zhenan Sun. Open-set single-domain generalization for
 736 robust face anti-spoofing. *International Journal of Computer
 737 Vision*, pages 1–22, 2024.
- [16] Geeho Kim, Junoh Kang, and Bohyun Han. Open-set rep-
 738 resentation learning through combinatorial embedding. In
 739 *Proceedings of the IEEE/CVF Conference on Computer Vision
 740 and Pattern Recognition*, pages 19744–19753, 2023.
- [17] Vladimir Koltchinskii. *Oracle inequalities in empirical risk
 741 minimization and sparse recovery problems: École D’Été
 742 de Probabilités de Saint-Flour XXXVIII-2008*, volume 2033.
 743 Springer Science & Business Media, 2011.
- [18] Yann LeCun, Bernhard Boser, John S Denker, Donnie Hen-
 744 derson, Richard E Howard, Wayne Hubbard, and Lawrence D
 745 Jackel. Backpropagation applied to handwritten zip code recog-
 746 nition. *Neural computation*, 1(4):541–551, 1989.
- [19] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M
 747 Hospedales. Deeper, broader and artier domain generalization.
 748 In *Proceedings of the IEEE international conference on com-*

- puter vision, pages 5542–5550, 2017.
- [20] Lei Li, Ke Gao, Juan Cao, Ziyao Huang, Yepeng Weng, Xiaoyue Mi, Zhengze Yu, Xiaoya Li, and Boyang Xia. Progressive domain expansion network for single domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 224–233, 2021.
- [21] Shuang Li, Kaixiong Gong, Binhui Xie, Chi Harold Liu, Weipeng Cao, and Song Tian. Critical classes and samples discovering for partial domain adaptation. *IEEE Transactions on Cybernetics*, 53(9):5641–5654, 2022.
- [22] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *International Conference on Machine Learning*, pages 6028–6039. PMLR, 2020.
- [23] Yi Liu, De Cheng, Dingwen Zhang, Shoukun Xu, and Jungong Han. Capsule networks with residual pose routing. *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [24] Yi Liu, Dingwen Zhang, Qiang Zhang, and Jungong Han. Part-object relational visual saliency. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3688–3704, 2021.
- [25] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *International conference on machine learning*, pages 10–18. PMLR, 2013.
- [26] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1406–1415, 2019.
- [27] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*, 2017.
- [28] Lei Qi, Hongpeng Yang, Yinghuan Shi, and Xin Geng. Normaug: Normalization-guided augmentation for domain generalization. *IEEE Transactions on Image Processing*, 33:1419–1431, 2024.
- [29] Fengchun Qiao, Long Zhao, and Xi Peng. Learning to learn single domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12556–12565, 2020.
- [30] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *Computer Vision–ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010, Proceedings, Part IV 11*, pages 213–226. Springer, 2010.
- [31] Kuniaki Saito, Shohei Yamamoto, Yoshitaka Ushiku, and Tatsuya Harada. Open set domain adaptation by backpropagation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 153–168, 2018.
- [32] Yang Shu, Zhangjie Cao, Chenyu Wang, Jianmin Wang, and Mingsheng Long. Open domain generalization with domain-augmented meta-learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9624–9633, 2021.
- [33] Mainak Singha, Ankit Jha, Shirsha Bose, Ashwin Nair, Moloud Abdar, and Biplob Banerjee. Unknown prompt the only lacuna: Unveiling clip’s potential for open domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13309–13319, 2024.
- [34] Aman Sinha, Hongseok Namkoong, Riccardo Volpi, and John Duchi. Certifying some distributional robustness with principled adversarial training. *arXiv preprint arXiv:1710.10571*, 2017.
- [35] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017.
- [36] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [37] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5018–5027, 2017.
- [38] Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John C Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. *Advances in neural information processing systems*, 31, 2018.
- [39] Fei Wang, Xing Zhang, Yong Jiang, Li Kong, and Xiaotong Wei. Patchcnn: An explicit convolution operator for point clouds perception. *IEEE Geoscience and Remote Sensing Letters*, 18(4):726–730, 2020.
- [40] Meng Wang, Changzhi Luo, Richang Hong, Jinhui Tang, and Jiashi Feng. Beyond object proposals: Random crop pooling for multi-label image recognition. *IEEE Transactions on Image Processing*, 25(12):5678–5688, 2016.
- [41] Xiran Wang, Jian Zhang, Lei Qi, and Yinghuan Shi. Generalizable decision boundaries: Dualistic meta-learning for open set domain generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11564–11573, 2023.
- [42] Zijian Wang, Yadan Luo, Ruihong Qiu, Zi Huang, and Mahsa Baktashmotagh. Learning to diversify for single domain generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 834–843, 2021.
- [43] Jianlong Wu, Keyu Long, Fei Wang, Chen Qian, Cheng Li, Zhouchen Lin, and Hongbin Zha. Deep comprehensive correlation mining for image clustering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8150–8159, 2019.
- [44] Zhiling Yan, Kai Zhang, Rong Zhou, Lifang He, Xiang Li, and Lichao Sun. Multimodal chatgpt for medical applications: an experimental study of gpt-4v. *arXiv preprint arXiv:2310.19061*, 2023.
- [45] Kaichao You, Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Universal domain adaptation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2720–2729, 2019.
- [46] Netzer Yuval. Reading digits in natural images with unsupervised feature learning. In *Proceedings of the NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [47] Kun Zhang, Mingming Gong, and Bernhard Schölkopf. Multi-source domain adaptation: A causal view. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- [48] Zhiwei Zhang, Junjie Yan, Sifei Liu, Zhen Lei, Dong Yi, and Stan Z Li. A face antispoofting database with diverse attacks. In *2012 5th IAPR international conference on Biometrics (ICB)*, pages 26–31. IEEE, 2012.
- [49] Long Zhao, Ting Liu, Xi Peng, and Dimitris Metaxas. Maximum-entropy adversarial data augmentation for improved generalization and robustness. *Advances in Neural Information Processing Systems*, 33:14435–14447, 2020.
- [50] Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Learning placeholders for open-set recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2021.
- [51] Ronghang Zhu and Sheng Li. Crossmatch: Cross-classifier consistency regularization for open-set single domain generalization. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25–29, 2022*. OpenReview.net, 2022.