

# MEMREIN: REIN THE DOMAIN SHIFT FOR CROSS-DOMAIN FEW-SHOT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Few-shot learning aims to enable models generalize to new categories (query instances) with only limited labeled samples (support instances) from each category. Metric-based mechanism is a promising direction which compares feature embeddings via different metrics. However, it always fail to generalize to unseen domains due to the considerable domain gap challenge. In this paper, we propose a novel framework, MemREIN, which considers Memorized, Restitution, and Instance Normalization for cross-domain few-shot learning. Specifically, an instance normalization algorithm is explored to alleviate feature dissimilarity, which provides the initial model generalization ability. However, naively normalizing the feature would lose fine-grained discriminative knowledge between different classes. To this end, a memorized module is further proposed to separate the most refined knowledge and remember it. Then, a restitution module is utilized to reconstitute the discrimination ability from the learned knowledge. A novel reverse contrastive learning strategy is proposed to stabilize the distillation process. Extensive experiments on five popular benchmark datasets demonstrate that MemREIN well addresses the domain shift challenge, and significantly improves the performance up to 16.37% compared with state-of-the-art baselines.

## 1 INTRODUCTION

In recent years, machine learning especially deep learning methods have made amazing achievements in the field of computer vision, image classification (He et al., 2016), semantic segmentation (Ren et al., 2016; He et al., 2017), etc. However, the high performance heavily relies on the large amount of well-labeled training data, which provides comprehensive and diverse samples to cover all corner cases. Such a huge scale makes it difficult in real practice, thus leads to a new topic of few-shot learning (Lake et al., 2015). Few-shot learning aims to enable models generalize to new categories (query instances) with only limited labeled samples (support instances) from each category.

Among existing few-shot learning methods, metric-based methods have attracted more attention because of their effectiveness and intelligibility. In general, the core idea of this kind of methods is to make classification based on the similarity between the query images and the support images via proposed similarity measurements. It usually consists of two main components: (1) feature encoder and (2) metric function. Given a task with few labeled images (support set) and unlabeled images (query set), the visual features are firstly extracted via the feature encoder and then passed through the defined metric function to determine the categories of the query images. The underlying assumption is that both training and testing are from the same dataset, namely the same domain. While, when it comes to different domains, the generalization ability of the metric-based methods greatly decreases (Chen et al., 2019; Tseng et al., 2020). However, such ability to generalize to unseen domains is of great importance in real practice, e.g., expensive human annotation or time-consuming data collection. As a result, considering the domain shift scenario within the few-shot learning has become an important yet challenging task.

Various unsupervised domain adaptation methods have been proposed (Yang et al., 2018; Ding et al., 2018). These methods aim to minimize the domain gap either by learning domain-invariant representations via representation learning, projection learning, or adversarial strategies (Long et al., 2015; Kumar et al., 2018; Ganin et al., 2016; Tzeng et al., 2015; 2017; Kundu et al., 2019). However, these methods assume that the complete unlabeled samples from the target domain are accessible

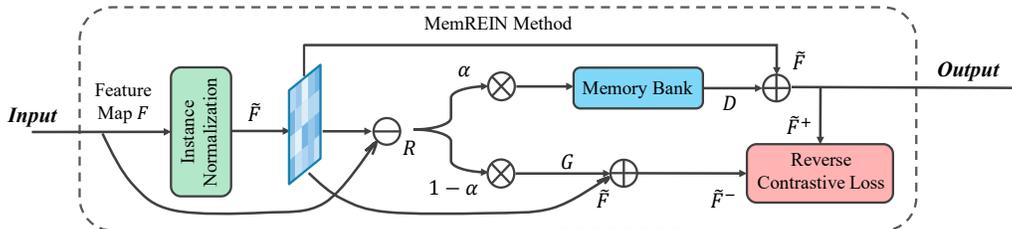


Figure 1: Framework of our MemREIN method. With instance normalization approach, the sample-specific features  $F$  can be reduced, and then with memorized and restitution approach, the long-term discriminative information can be distilled and restituted to refined features.

while training. We argue that this assumption may not hold in real situations, and it could lead to high computational cost in testing phase. Domain shift problem could be addressed by various domain generalization methods (Blanchard et al., 2011; Muandet et al., 2013; Motiian et al., 2017). However, these methods assume that the source and target domains share the same categories. In contrast, our goal is to recognize novel categories from the target domain with only a few (e.g., 1 or 5) of samples selected from novel categories.

As argued above, there are two main challenges in cross-domain few-shot learning task. First, how to minimize the discrepancy between the source domain and the target domain. Second, how to recognize novel/unseen classes with only limited samples.

To this end, we propose a novel MemREIN approach, which includes Memorized, Restitution, and Instance Normalization as crucial modules, to “rein” the domain shift level in few-shot scenario. The core idea of MemREIN is to enhance the generalization ability while still be able to balance the discrimination ability for subsequent classification. In specific, on the training stage, we first present an instance normalization layer operating on features with respect to samples at the channel level. This operation aims to reserve spatial feature dependency and meanwhile remove the image-specific features, i.e., alleviate the discrepancy of these training samples. By this way, the generalization ability across different samples is enhanced. Then, the filtered out features are extracted from a residual structure. Normally, the filtered out features are considered as useless feature which could be discarded. However, we consider it still contains fine-grained distinctive knowledge which could be “remembered” and “restituted”. To this end, we manage to adaptively distill the long-term discriminative information from them via our proposed novel memorized approach. Then, such discriminative information is restituted to the above refined features to maintain the discrimination ability for subsequent classification. A novel reverse contrastive loss constraint in the restitution step to encourage the better separation of discriminative features and general features, which ensures the distillation process. Contributions of our work are summarized as follows:

- A novel memorized and restitution strategy is proposed for discriminative information distillation. It is able to distill the long-term discriminative information from filtered out features to maintain the discrimination ability of original features for better classification.
- An instance normalization strategy is adopted to alleviate the discrepancy across training samples, which reduces the sample-specific features and greatly enhances the overall generalization ability across sample features.
- A novel reverse contrastive loss is proposed to encourage the better separation of discriminative features and general features, which is able to ensure the distillation process.

Our MemREIN method is simple yet effective. It is a universal and method-agnostic method that can be applied to various existing metric-based methods for enhancing their generalization ability to unseen domains. Extensive experiments demonstrate the effectiveness of MemREIN, which achieves consistent superior performance than existing state-of-the-art methods under the cross-domain setting.

## 2 RELATED WORK

### 2.1 FEW-SHOT CLASSIFICATION

Few-shot classification aims to recognize novel classes with a limited amount of labeled samples. Among these existing methods, metric-based methods have attracted considerable attention and

achieved promising performance. This kind of methods usually consists of two components: (1) feature encoder and (2) metric function. The feature encoder is used to extract features from both query and support samples. The metric function is used to calculate the similarity for classification. For instance, MatchingNet (Vinyals et al., 2016) utilizes cosine similarity with an attention Bi-LSTM for classification and ProtoNet (Snell et al., 2017) applies euclidean distance for classification. RelationNet (Sung et al., 2018) uses convolutional neural networks and GNN (Satorras & Estrach, 2018) uses the graph convolutional framework as the metric function. Although these methods have achieved promising performance, they always fail to generalize to unseen domains (datasets) since the distributions among different domains have huge shifts. Recent work (Chen et al., 2019) reveals that the performance of existing few-shot learning methods degrades significantly under the cross-domain setting. The motivation of our work aims to enhance the generalization ability of metric-based few-shot learning methods so that these methods can better generalize to unseen domains.

## 2.2 DOMAIN ADAPTATION

The domain adaptation methods can be divided into four categories depends on the relation between the source classes  $\mathcal{C}_s$  and target domain classes  $\mathcal{C}_t$ . **Closed-set domain adaptation (CDA)** assumes the source and target domain share the same categories, existing methods focus on how to alleviate the feature distribution gap (Borgwardt et al., 2006; Long et al., 2013; Sun & Saenko, 2016; Gopalan et al., 2011; Ganin & Lempitsky, 2015; Tzeng et al., 2017; Hoffman et al., 2018). **Partial domain adaptation (PDA)** assumes that the source domain is large enough to contain all the classes of target domain. Many works (Cao et al., 2018a; Zhang et al., 2018; Cao et al., 2018b) have been proposed to solve this problem. Studies on partial domain adaptation promote the setting to a more common and practical level. **Open-set domain adaptation (OSDA)** was first proposed by Busto *et al.* (Panareda Busto & Gall, 2017), which made a new definition of “unknown” category if this class is private to source or target domain. The setting of open-set domain adaptation is more practical compared with partial domain adaptation. **Universal Domain Adaptation (UniDA)** was first proposed in the work (You et al., 2019), in which the relation between source categories and target categories is inaccessible. Following this setting, recent methods Fu et al. (2020); Kundu et al. (2020); Saito et al. (2020); Saito & Saenko (2021) have been proposed to address this challenge.

## 2.3 CROSS-DOMAIN FEW-SHOT LEARNING

Recently, promoted by the pioneer work (Chen et al., 2019), cross-domain few-shot learning problem has attracted many attentions. As an emerging task, work () carried out a broader study and introduced a new benchmark. Some methods (Tseng et al., 2020; Sun et al., 2021; Phoo & Hariharan, 2020; Zou et al., 2021; Islam et al., 2021) have been proposed and achieve promising performance under this benchmark. Work (Cai et al., 2020) relaxes this setting where a large number of unlabeled target samples are accessible in the training phase. Most recently, method ATA (Wang & Deng, 2021) introduced an adversarial task augmentation method to improve the robustness of the inductive bias under the cross-domain few-shot learning setting. In addition, a noise-enhanced supervised auto-encoder method was proposed in (Liang et al., 2021) to obtain the broader variations of the feature distributions to greatly boost the generalization capability of the model. Paper (Fu et al., 2021) proposed an effective mix-up module into the meta-learning mechanism and a novel disentangle module with a domain classifier to obtain domain-irrelevant and domain-specific features, which achieves promising performance. In our work, we propose a simple yet effective method from the perspective of feature level, which is a universal and method-agnostic method.

# 3 METHOD

## 3.1 PRELIMINARIES

In the few-shot classification problem, define a task  $T$  characterized as  $N_w$  way and  $N_s$  shot, which represents the number of categories and the number of labeled samples in each category, respectively. At each iteration, the metric-based few-shot learning method first randomly samples  $N_w$  categories as a task  $T$ , and then constructs a support set  $S = \{(\mathcal{X}_s, \mathcal{Y}_s)\}$  and a query set  $Q = \{(\mathcal{X}_q, \mathcal{Y}_q)\}$ . These two sets are constructed by randomly selecting  $N_s$  and  $N_q$  samples for each of the  $N_w$  categories.

Once the data is prepared, the feature encoder  $E$  first extracts features of the samples from both support set  $S$  and the query set  $Q$ . Then, the defined metric function  $M$  makes predictions of the

query samples  $\mathcal{X}_q$  based on three parts: the label of support samples  $\mathcal{Y}_s$ , encoded query image  $E(\mathcal{X}_q)$ , and the encoded support images  $E(\mathcal{X}_s)$ , which is formulated as follows:

$$\hat{\mathcal{Y}}_q = M(\mathcal{Y}_s, E(\mathcal{X}_q), E(\mathcal{X}_s)). \quad (1)$$

After all, the objective of the metric-based few-shot learning method is the classification loss of the samples in the query set, which is formulated as follows:

$$\mathcal{L} = \mathcal{L}_{cls}(\mathcal{Y}_q, \hat{\mathcal{Y}}_q). \quad (2)$$

The main difference among existing metric-based few-shot learning methods lies in the different metric functions. Differently, we propose a universal method that can be applied in all the metric-based few-shot learning methods to achieve better performance under the cross-domain setting.

In this paper, we tackle the cross-domain few-shot classification problem. Given a set of few-shot classification tasks  $\mathcal{T} = \{T_1, T_2, \dots, T_n\}$  as a domain (dataset). At the training stage, given  $N$  accessible domains  $\{\mathcal{T}_1^{seen}, \mathcal{T}_2^{seen}, \dots, \mathcal{T}_N^{seen}\}$ , we aim to learn a metric-based few-shot learning model with these seen domains, then the model can generalize to an unseen domain  $\mathcal{T}^{unseen}$ .

## 3.2 MEMREIN METHOD

The core idea of our MemREIN method is to enhance the generalization ability, contain the ability to balance the discrimination of metric-based few-shot learning methods, and achieve promising performance on arbitrary unseen domains. The overall framework of our MemREIN method is illustrated in Figure 1. MemREIN is method-agnostic that can be applied to existing metric-based few-shot learning methods to improve the performance to unseen domains. In addition, it is a universal framework that can be applied by various neural networks for different applications, e.g. classification, segmentation, detection. In this paper, we delve into the cross-domain few-shot learning problem and propose our MemREIN method to “rein” the domain shift level in few-shot classification.

### 3.2.1 INSTANCE NORMALIZATION

As argued above, images with the same category from different domains have large discrepancies in many aspects e.g. , image style, color, quality. Generally speaking, the discrepancy between the source domain and the target domain hinders the generalization ability of the model to some extent.

To this end, we reduce the discrepancy cross samples by instance normalization in our proposed MemREIN method as follows. Denote the input feature map by  $F \in \mathbb{R}^{c \times h \times w}$  and the output feature map by  $\tilde{F} \in \mathbb{R}^{c \times h \times w}$ , where  $c, h, w$  denote the number of channel, height, width, respectively.

$$\tilde{F} = \text{IN}(F) = \gamma \left( \frac{F - \mu(F)}{\sigma(F)} \right) + \beta, \quad (3)$$

where  $\mu(\cdot)$  and  $\sigma(\cdot)$  denote the mean and standard deviation calculated at the channel level for each sample,  $\gamma \in \mathbb{R}^c$  and  $\beta \in \mathbb{R}^c$  are two trainable parameters. Instance normalization was originally used in style transfer (Dumoulin et al., 2016; Ulyanov et al., 2016; Huang & Belongie, 2017), which is helpful to enhance the generalization ability by reducing the feature dissimilarity. It can remove instance/sample specific features out of the input, which makes more general features remained (Shankar et al., 2018; Volpi et al., 2018).

However, instance normalization inevitably removes some discriminative information from the original feature maps (Jin et al., 2020; 2021), which weakens the discrimination ability of the extracted features. To address this emerging problem, we propose a memorized restitution approach to distill the discriminative information from the filtered out features and then reconstitute it as the final output feature maps.

### 3.2.2 MEMORIZED RESTITUTION

As discussed above, in order to maintain the discrimination ability of the refined features, we propose a following memorized restitution approach to distill discriminative information. We first obtain the filtered out feature  $R$  via a residual structure, which is defined as follows:

$$R = F - \tilde{F}, \quad (4)$$

where  $R \in \mathbb{R}^{c \times h \times w}$ , denoting the features that we have filtered out via the instance normalization operation. Since instance normalization operation will inevitably remove discriminative information from the original features. Hence, there exist discriminative features that we need to distill and purify from the residual feature  $R$ , in order to maintain the discrimination ability of extracted features.

At the training stage, given the feature map  $R$  at each iteration (we omit the subscript of feature map  $R$  for brevity), we assume  $R$  consists of two parts:  $D \in \mathbb{R}^{c \times h \times w}$  with relative more discriminative information, and  $G \in \mathbb{R}^{c \times h \times w}$  with relatively more general information, which is defined as follows:

$$\begin{cases} D(k, :, :) = \alpha_k R(k, :, :), \\ G(k, :, :) = (1 - \alpha_k) R(k, :, :), \end{cases} \quad (5)$$

where  $k$  denotes the  $k^{th}$  channel of the feature map,  $\alpha_k$  denotes the learnable attention parameters to split the residual feature map  $R$ . Note that we split the residual feature map  $R$  at the channel level.

Then, the attention vector  $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_c]$  is derived by SE-like channel attention (Hu et al., 2018) as follows:

$$\alpha = \delta(W_2 \eta(W_1 \text{avepooling}(R))), \quad (6)$$

where *avepooling* represents the average pooling layer,  $W_1$  and  $W_2$  are parameters to be learned,  $\delta$  and  $\eta$  represent the ReLU activation function and sigmoid activation function, respectively.

Since there are limited labeled samples under the few-shot learning framework, it is highly possible that the model would overfit. Thus we further propose a memorized mechanism with a memory vector  $M^{(l)} \in \mathbb{R}^c$  to store the long-term feature maps  $D$ , which is defined as follows:

$$\begin{aligned} M^{(l)} &= [M_1^{(l)}, \dots, M_k^{(l)}, \dots, M_c^{(l)}], \\ M_k^{(l+1)} &= D^{(l)}(k, :, :), \end{aligned} \quad (7)$$

where  $M_k^{(l)} \in \mathbb{R}^{h \times w}$ ,  $(l)$  represents the  $l^{th}$  iteration,  $k$  denotes the  $k^{th}$  channel. At the  $l^{th}$  iteration, we concatenate the feature map  $D$  to the memory bank at the channel level, and update  $D$  as follows:

$$D(k, :, :) = \text{maxpooling}(\text{concat}(M_k^{(l)}, D(k, :, :))), \quad (8)$$

where *concat* represents the concatenation operation, *maxpooling* represents the max pooling layer.

Once we obtain the updated feature map  $D$ , we restitute it to refined feature  $\tilde{F}$  as the final output  $\tilde{F}^+$  of our proposed MemREIN method, and we also restitute the relatively unimportant feature map  $G$  with feature  $\tilde{F}$  as the ‘‘contaminated’’ feature  $\tilde{F}^-$  for following loss optimization as follows:

$$\tilde{F}^+ = \tilde{F} + D, \quad \tilde{F}^- = \tilde{F} + G. \quad (9)$$

### 3.2.3 REVERSE CONTRASTIVE LOSS

Apart from the conventional cross-entropy loss defined in Equation 2, we also propose a novel reverse contrastive loss  $\mathcal{L}_{rcl}$  to promote the disentanglement of feature  $D$  and feature  $G$ . It consists of two parts:  $\mathcal{L}_{rcl}^+$  and  $\mathcal{L}_{rcl}^-$ , e.g.,  $\mathcal{L}_{rcl} = \mathcal{L}_{rcl}^+ + \mathcal{L}_{rcl}^-$ . Given a mini-batch  $\mathcal{X}_b = \{\mathcal{X}_1, \dots, \mathcal{X}_N\}$  contains  $N$  samples at the training phase, we first randomly select one anchor sample referred as  $\mathcal{X}_a$ , and then we denote samples with the same category as the positive samples  $\mathcal{X}_{pos}$ , samples with different categories as the negative samples  $\mathcal{X}_{neg}$ . Note that the corresponding features of these samples are denoted with their subscripts such as  $\tilde{F}_a$ ,  $\tilde{F}_{pos}$ , and  $\tilde{F}_{neg}$  in the following paragraphs.

We first reshape features  $\tilde{F}^+$  and  $\tilde{F}^-$  to the size of  $\mathbb{R}^{chw \times 1}$  and then pass them through one fully-connected layer following the *softmax* function to obtain the feature vectors  $\tilde{f}^+$  and  $\tilde{f}^-$ , which is defined as follows. Note that these two vectors have the same size of  $\in \mathbb{R}^{K \times 1}$ .

$$\tilde{f}^+ = \text{softmax}\left(W^+ \text{reshape}(\tilde{F}^+)\right), \quad \tilde{f}^- = \text{softmax}\left(W^- \text{reshape}(\tilde{F}^-)\right), \quad (10)$$

Table 1: Classification accuracy (%) of 5-way 1/5-shot tasks trained on the dataset mini-ImageNet.

5-way 1-shot	Classification Accuracy (%)			
	CUB	Cars	Places	Plantae
MNet (Vinyals et al., 2016)	35.89 ± 0.51%	30.77 ± 0.47%	49.86 ± 0.79%	32.70 ± 0.60%
MNet+FT (Tseng et al., 2020)	36.61 ± 0.53%	29.82 ± 0.44%	51.07 ± 0.68%	34.48 ± 0.50%
MNet+ATA (Wang & Deng, 2021)	41.59 ± 0.40%	35.14 ± 0.30%	51.86 ± 0.50%	37.02 ± 0.30%
<b>MNet+MemREIN (Ours)</b>	<b>43.72 ± 0.45%</b>	<b>37.64 ± 0.43%</b>	<b>53.44 ± 0.62%</b>	<b>39.83 ± 0.50%</b>
RNet (Sung et al., 2018)	42.44 ± 0.77%	29.11 ± 0.60%	48.64 ± 0.85%	33.17 ± 0.64%
RNet+FT (Tseng et al., 2020)	44.07 ± 0.77%	28.63 ± 0.59%	50.68 ± 0.87%	33.14 ± 0.62%
RNet+LRP (Sun et al., 2021)	41.57 ± 0.40%	30.48 ± 0.30%	48.47 ± 0.50%	32.11 ± 0.30%
RNet+ATA (Wang & Deng, 2021)	43.02 ± 0.40%	31.79 ± 0.30%	51.16 ± 0.50%	33.72 ± 0.30%
<b>RNet+MemREIN (Ours)</b>	<b>47.33 ± 0.49%</b>	<b>34.92 ± 0.33%</b>	<b>55.75 ± 0.52%</b>	<b>36.27 ± 0.32%</b>
GNN (Satorras & Estrach, 2018)	45.69 ± 0.68%	31.79 ± 0.51%	53.10 ± 0.80%	35.60 ± 0.56%
GNN+FT (Tseng et al., 2020)	47.47 ± 0.75%	31.61 ± 0.53%	55.77 ± 0.79%	35.95 ± 0.58%
GNN+LRP (Sun et al., 2021)	48.29 ± 0.51%	32.78 ± 0.39%	54.83 ± 0.56%	37.49 ± 0.43%
GNN+ATA (Wang & Deng, 2021)	45.00 ± 0.50%	33.61 ± 0.40%	53.57 ± 0.50%	34.42 ± 0.40%
<b>GNN+MemREIN (Ours)</b>	<b>49.94 ± 0.50%</b>	<b>36.04 ± 0.44%</b>	<b>57.35 ± 0.52%</b>	<b>39.09 ± 0.46%</b>
5-way 5-shot	Classification Accuracy (%)			
	CUB	Cars	Places	Plantae
MNet (Vinyals et al., 2016)	51.37 ± 0.77%	38.99 ± 0.64%	63.16 ± 0.77%	46.53 ± 0.68%
MNet+FT (Tseng et al., 2020)	55.23 ± 0.83%	41.24 ± 0.65%	64.55 ± 0.75%	41.69 ± 0.63%
MNet+ATA (Wang & Deng, 2021)	59.33 ± 0.40%	48.78 ± 0.40%	66.31 ± 0.40%	51.56 ± 0.30%
<b>MNet+MemREIN (Ours)</b>	<b>63.87 ± 0.66%</b>	<b>49.47 ± 0.60%</b>	<b>69.08 ± 0.65%</b>	<b>52.98 ± 0.34%</b>
RNet (Sung et al., 2018)	57.77 ± 0.69%	37.33 ± 0.68%	63.32 ± 0.76%	44.00 ± 0.60%
RNet+FT (Tseng et al., 2020)	59.46 ± 0.71%	39.91 ± 0.69%	66.28 ± 0.72%	45.08 ± 0.59%
RNet+LRP (Sun et al., 2021)	57.70 ± 0.40%	41.21 ± 0.40%	65.35 ± 0.40%	43.70 ± 0.30%
RNet+ATA (Wang & Deng, 2021)	59.36 ± 0.40%	42.95 ± 0.40%	66.90 ± 0.40%	45.32 ± 0.30%
<b>RNet+MemREIN (Ours)</b>	<b>63.31 ± 0.42%</b>	<b>46.75 ± 0.44%</b>	<b>70.84 ± 0.52%</b>	<b>49.52 ± 0.40%</b>
GNN (Satorras & Estrach, 2018)	62.25 ± 0.65%	44.28 ± 0.63%	70.84 ± 0.65%	52.53 ± 0.59%
GNN+FT (Tseng et al., 2020)	66.98 ± 0.68%	44.90 ± 0.64%	73.94 ± 0.67%	53.85 ± 0.62%
GNN+LRP (Sun et al., 2021)	64.44 ± 0.48%	46.20 ± 0.46%	74.45 ± 0.47%	54.46 ± 0.46%
GNN+ATA (Wang & Deng, 2021)	66.22 ± 0.50%	49.14 ± 0.40%	75.48 ± 0.40%	52.69 ± 0.40%
<b>GNN+MemREIN (Ours)</b>	<b>72.41 ± 0.56%</b>	<b>49.98 ± 0.43%</b>	<b>77.71 ± 0.54%</b>	<b>56.64 ± 0.45%</b>

where  $W^+$  and  $W^-$  are trainable parameters with the same size of  $\mathbb{R}^{K \times chw}$ ,  $K$  is the number of classes in the few-shot classification task. Then, the reverse contrastive loss is defined as follows:

$$\mathcal{L}_{rcl}^+ = -\mathbb{E} \left[ \log \frac{\exp(\tilde{f}_a^+ \top \tilde{f}_{pos}^+)}{\sum_{\mathcal{X}_{pos} \in \mathcal{X}} \exp(\tilde{f}_a^+ \top \tilde{f}_{neg}^+)} \right], \quad (11)$$

$$\mathcal{L}_{rcl}^- = -\mathbb{E} \left[ \log \frac{\sum_{\mathcal{X}_{neg} \in \mathcal{X}} \exp(\tilde{f}_a^- \top \tilde{f}_{neg}^-)}{\exp(\tilde{f}_a^- \top \tilde{f}_{pos}^-)} \right]. \quad (12)$$

The goal of our proposed reverse contrastive loss is to promote the disentanglement of feature  $D$  and feature  $G$ , where feature  $D$  contains more discriminative information and  $G$  contains more general information. Combining feature  $D$  with the refined feature  $\tilde{F}$ , defined in Equation 9, results in better discrimination capability of feature  $\tilde{F}^+$ , in other words, the sample features with same category are closer and those with different identities are farther apart. Therefore, we propose  $\mathcal{L}_{rcl}^+$  to promote the features of positive samples  $\tilde{f}_{pos}^+$  gather closer and separate the features of negative samples  $\tilde{f}_{neg}^+$  from the anchor feature as well. On the other hand, combining feature  $G$  with the refined feature  $\tilde{F}$  results in decreasing the discrimination capability, which means the feature  $\tilde{F}^-$  is more general that not capable of distinguishing samples with the same category correctly. Therefore, we propose  $\mathcal{L}_{rcl}^-$  to separate the the features of positive samples  $\tilde{f}_{pos}^-$  from both features with negative samples  $\tilde{f}_{neg}^-$  and the anchor feature  $\tilde{f}_a^-$ . The whole objective loss is defined as follows:

Table 2: Classification accuracy (%) of 5-way 1/5-shot tasks under the leave-one-out setting.

5-way 1-shot	Classification Accuracy (%)			
	CUB	Cars	Places	Plantae
MNet (Vinyals et al., 2016)	37.90 ± 0.55%	28.96 ± 0.45%	49.01 ± 0.65%	33.21 ± 0.51%
MNet+LFT (Tseng et al., 2020)	43.29 ± 0.59%	30.62 ± 0.48%	52.51 ± 0.67%	35.12 ± 0.54%
<b>MNet+MemREIN (Ours)</b>	<b>46.37 ± 0.50%</b>	<b>35.65 ± 0.45%</b>	<b>54.92 ± 0.64%</b>	<b>38.82 ± 0.48%</b>
RNet (Sung et al., 2018)	44.33 ± 0.59%	29.53 ± 0.45%	47.76 ± 0.63%	33.76 ± 0.52%
RNet+LFT (Tseng et al., 2020)	48.38 ± 0.63%	32.21 ± 0.51%	50.74 ± 0.66%	35.00 ± 0.52%
<b>RNet+MemREIN (Ours)</b>	<b>52.02 ± 0.52%</b>	<b>36.38 ± 0.38%</b>	<b>54.82 ± 0.57%</b>	<b>36.74 ± 0.45%</b>
GNN (Satorras & Estrach, 2018)	49.46 ± 0.73%	32.95 ± 0.56%	51.39 ± 0.80%	37.15 ± 0.60%
GNN+LFT (Tseng et al., 2020)	51.51 ± 0.80%	34.12 ± 0.63%	56.31 ± 0.80%	42.09 ± 0.68%
<b>GNN+MemREIN (Ours)</b>	<b>54.26 ± 0.62%</b>	<b>37.55 ± 0.50%</b>	<b>59.98 ± 0.64%</b>	<b>45.69 ± 0.64%</b>
5-way 5-shot	Classification Accuracy (%)			
	CUB	Cars	Places	Plantae
MNet (Vinyals et al., 2016)	51.92 ± 0.80%	39.87 ± 0.51%	61.82 ± 0.57%	47.29 ± 0.51%
MNet+LFT (Tseng et al., 2020)	61.41 ± 0.57%	43.08 ± 0.55%	64.99 ± 0.59%	48.32 ± 0.57%
<b>MNet+MemREIN (Ours)</b>	<b>67.31 ± 0.51%</b>	<b>47.36 ± 0.48%</b>	<b>68.14 ± 0.58%</b>	<b>52.28 ± 0.52%</b>
RNet (Sung et al., 2018)	62.13 ± 0.74%	40.64 ± 0.54%	64.34 ± 0.57%	46.29 ± 0.56%
RNet+LFT (Tseng et al., 2020)	64.99 ± 0.54%	43.44 ± 0.59%	67.35 ± 0.54%	50.39 ± 0.52%
<b>RNet+MemREIN (Ours)</b>	<b>68.39 ± 0.48%</b>	<b>46.92 ± 0.50%</b>	<b>69.87 ± 0.54%</b>	<b>58.64 ± 0.50%</b>
GNN (Satorras & Estrach, 2018)	69.26 ± 0.68%	48.91 ± 0.67%	72.59 ± 0.67%	58.36 ± 0.68%
GNN+LFT (Tseng et al., 2020)	73.11 ± 0.68%	49.88 ± 0.67%	77.05 ± 0.65%	58.84 ± 0.66%
<b>GNN+MemREIN (Ours)</b>	<b>77.54 ± 0.62%</b>	<b>56.78 ± 0.66%</b>	<b>78.84 ± 0.66%</b>	<b>65.44 ± 0.64%</b>

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda(\mathcal{L}_{rcl}^+ + \mathcal{L}_{rcl}^-), \quad (13)$$

where  $\lambda$  is a hyper-parameter to control the balance of these two terms in the training phase.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

**Baselines:** We make extensive experiments on three existing metric-based few-shot learning methods: MatchingNet (Vinyals et al., 2016), RelationNet (Sung et al., 2018), and GNN (Satorras & Estrach, 2018). We compare our proposed method with several existing cross-domain few-shot learning methods: FT (Tseng et al., 2020), LRP (Sun et al., 2021), and ATA (Wang & Deng, 2021) to demonstrate the advantages of our method. In addition, we also make comparisons with several existing state-of-the-art few-shot learning methods: TADAM (Oreshkin et al., 2018), DC (Lifchitz et al., 2019), DC+IMP (Lifchitz et al., 2019), MetaOptNet (Lee et al., 2019), EGNN (Kim et al., 2019), TPN (Yanbin et al., 2019), DPGN (Yang et al., 2020), MCGN (Tang et al., 2021), ECKPN (Chen et al., 2021), FromActi (Qiao et al., 2018), and LEO (Rusu et al., 2018) for further demonstration. More quantitative results and visualizations are provided in Appendix B and C.

**Datasets:** We conduct experiments on five public datasets that are widely used for few-shot classification task<sup>1</sup>: mini-ImageNet (Ravi & Larochelle, 2016), CUB (Wah et al., 2011), Cars (Krause et al., 2013), Places (Zhou et al., 2017), and Plantae (Van Horn et al., 2018). Detailed introduction and statistical analysis of these five datasets could be found in Appendix A.

**Setting:** To ensure the fair comparison with other methods, we apply the exactly same two cross-domain settings that applied in the baselines. The first setting is to train on the dataset mini-ImageNet and test on other four datasets, which means there is only one source domain and one target domain. The second setting is the leave-one-out setting by selecting one dataset among CUB, Cars, Places, and Plantae as the target domain for testing, and using the remaining three datasets along with dataset mini-ImageNet as the source domains for training. The second setting is more difficult since there are multiple source domains with only one target domain, which results in much larger domain shift.

<sup>1</sup>Detailed introductions of these datasets are provided in the appendix.

Table 3: Classification Accuracy (%) of our proposed method comparing with the state-of-the-art few-shot learning methods. The model is trained/evaluated only on the dataset mini-ImageNet.

Backbone	Method	Classification Accuracy (%)	
		5-way 1-shot	5-way 5-shot
ResNet-12	TADAM (Oreshkin et al., 2018)	58.50 ± 0.30%	76.70 ± 0.30%
	DC (Lifchitz et al., 2019)	62.53 ± 0.19%	78.95 ± 0.13%
	DC+IMP (Lifchitz et al., 2019)	–	79.77 ± 0.19%
	MetaOptNet (Lee et al., 2019)	64.09 ± 0.62%	80.00 ± 0.45%
	EGNN (Kim et al., 2019)	59.63 ± 0.52%	76.34 ± 0.48%
	TPN (Yanbin et al., 2019)	55.51 ± 0.86%	69.86 ± 0.65%
	DPGN (Yang et al., 2020)	66.14 ± 0.43%	81.23 ± 0.41%
	MCGN (Tang et al., 2021)	67.32 ± 0.43%	83.03 ± 0.54%
WRN-28	ECKPN (Chen et al., 2021)	<b>70.48 ± 0.38%</b>	<b>85.42 ± 0.46%</b>
	FromActi (Qiao et al., 2018)	59.60 ± 0.41%	77.74 ± 0.19%
	LEO (Rusu et al., 2018)	61.76 ± 0.08%	77.59 ± 0.12%
ResNet-10	MNet (Vinyals et al., 2016)	59.10 ± 0.64%	70.96 ± 0.65%
	MNet+LFT (Tseng et al., 2020)	58.76 ± 0.61%	72.53 ± 0.69%
	<b>MNet+MemREIN (Ours)</b>	60.03 ± 0.60%	74.72 ± 0.66%
ResNet-10	RNet (Sung et al., 2018)	57.80 ± 0.88%	71.00 ± 0.69%
	RNet+LFT (Tseng et al., 2020)	58.64 ± 0.85%	73.78 ± 0.64%
	<b>RNet+MemREIN (Ours)</b>	61.64 ± 0.74%	75.98 ± 0.56%
ResNet-10	GNN (Satorras & Estrach, 2018)	60.77 ± 0.75%	80.87 ± 0.56%
	GNN+LFT (Tseng et al., 2020)	66.32 ± 0.80%	81.98 ± 0.55%
	<b>GNN+MemREIN (Ours)</b>	<b>70.64 ± 0.72%</b>	<b>85.48 ± 0.52%</b>

**Implementation details:** We use the public implementation<sup>2</sup> to train both the MatchingNet and the RelationNet method, and we use the implementation<sup>3</sup> to train the GNN method. In all the experiments, we adopt the ResNet-10 (He et al., 2016) as the backbone network for our feature encoder  $E$ . We insert our proposed MemREIN method after the last batch normalization layer of all the residual blocks in the feature encoder  $E$  at the training stage. Instead of optimizing from the scratch, we apply a strategy that pre-trains the feature extractor by minimizing the standard cross-entropy classification loss on the 64 training categories from the dataset mini-ImageNet and this strategy is also applied in all the baselines. In the training phase, we set  $\lambda = 0.1$  and train 1000 trials for all the methods. In each trial, we randomly sample  $N_w$  categories with  $N_s$  randomly selected images for each support set, and 16 images for the query set. We use the Adam optimizer with the learning rate 0.001.

## 4.2 EXPERIMENTAL RESULTS

### 4.2.1 GENERALIZATION FROM ONE SOURCE DOMAIN

Table 1 shows the results of the first cross-domain setting that the model is trained on the dataset mini-ImageNet, and tested on the other four datasets CUB, Cars, Places and Plantae, respectively. The results demonstrate that with our proposed MemREIN method, the performance of all three metric-based few-shot learning methods makes obvious improvements, which validates the effectiveness of our proposed MemREIN method to enhance the generalization ability to the unseen domain. In addition, our proposed method consistently outperforms other existing cross-domain few-shot learning methods, which indicates the the superiority of our method over previous best methods.

### 4.2.2 GENERALIZATION FROM MULTIPLE SOURCE DOMAINS

Table 2 shows the results under the leave-one-out setting. We first select out one dataset as the unseen domain for testing and use the remaining three datasets as well as the dataset mini-ImageNet for training since we already use the dataset mini-ImageNet for pre-training. Note that the baseline (Tseng et al., 2020) has two different training strategies, one is the “learn to learn” strategy and another is using fixed hyper-parameters. We consider the better results for comparison here, which is denoted as “+LFT” in the Table 2. The results demonstrate that our proposed MemREIN method can greatly

<sup>2</sup><https://github.com/wyharveychen/CloserLookFewShot>.

<sup>3</sup><https://github.com/hytseng0509/CrossDomainFewShot>.

Table 4: Ablation study on our method and the objective function. We consider the dataset CUB as the unseen domain. “GNN+IN” indicates that we only employ the instance normalization strategy, “w/o  $\mathcal{L}_{rcl}^-$ ” indicates that we remove the  $\mathcal{L}_{rcl}^-$  term, and “w/o  $\mathcal{L}_{rcl}^+$ ” indicates that we remove the  $\mathcal{L}_{rcl}^+$  term.

5-way 5-shot		Classification Accuracy (%)	
Variant ID	Method	CUB	
1	GNN (Satorras & Estrach, 2018)	69.26 $\pm$ 0.68%	
2	GNN+IN	67.34 $\pm$ 0.66%	
3	GNN+MemREIN	<b>77.54 <math>\pm</math> 0.62%</b>	
4	w/o $\mathcal{L}_{rcl}^-$	75.38 $\pm$ 0.63%	
5	w/o $\mathcal{L}_{rcl}^+$	73.02 $\pm$ 0.62%	

improve the performance of all three metric-based few-shot learning methods, which reflects that our method has the capability of mitigating the domain gap problem. In addition, results show that our method consistently outperforms the “+LFT” method, which validates that our proposed method can better capture the variation of feature distributions across multiple domains than the “+LFT” method, thus the generalization ability of extracted features are better enhanced.

#### 4.2.3 COMPARE WITH SOTA METHODS

We also make comparisons to existing SOTA few-shot learning methods to further demonstrate the superiority of our proposed method. The results are shown in Table 3. To ensure fair comparison, all the methods are trained and evaluated only on the dataset mini-ImageNet. It is the conventional few-shot learning setting and there is no unseen domain. The results indicate that our proposed method can achieve competitive performance compared with SOTA few-shot learning methods.

#### 4.2.4 ABLATION STUDY

We make ablation study on our propose method to demonstrate the functionality of two main parts: (1) instance normalization and (2) memorized restitution. We take the GNN baseline under the leave-one-out setting (5-way 5-shot) on the dataset CUB as the example. Meanwhile, we also make ablation study on our proposed reverse contrastive loss to demonstrate the functionality of two terms:  $\mathcal{L}_{rcl}^+$  and  $\mathcal{L}_{rcl}^-$ . Table 7 indicates the results of the ablation study. Comparing the results of Variant 1 and 2, it indicates that only applying the instance normalization operation results in the decrease of the accuracy. It is reasonable because the instance normalization operation is too “strong” to inevitably remove some discriminative useful information. It also validates the necessity and effectiveness of our proposed memorized restitution approach. Comparing the results of Variant 3, 4, and 5, it indicates that both  $\mathcal{L}_{rcl}^+$  and  $\mathcal{L}_{rcl}^-$  contribute to the final results, which means our proposed reverse contrastive loss is capable of promoting the disentanglement of features effectively.

## 5 CONCLUSION

In this paper, we investigated the cross-domain few-shot classification problem where exists the domain gap issue. We propose a novel framework, MemREIN, which considers Memorized, Restitution, and Instance Normalization to address this issue. We first alleviate feature dissimilarity across sample features via an instance normalization algorithm to enhance the overall generalization ability. In order to avoid the loss of fine-grained discriminative knowledge between different classes, a memorized restitution approach is further proposed to adaptively remember the long-term refined knowledge and retribute the discrimination ability. Finally, A novel reverse contrastive learning strategy is proposed to stabilize the distillation process. Extensive experiments on five popular benchmark datasets demonstrate that MemREIN well addresses the domain shift challenge, and significantly improves the performance up to 16.37% compared with state-of-the-art baselines.

## REFERENCES

- Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification tasks to a new unlabeled sample. *Proceedings of the Advances in Neural Information Processing Systems*, 24:2178–2186, 2011.
- Karsten M Borgwardt, Arthur Gretton, Malte J Rasch, Hans-Peter Kriegel, Bernhard Schölkopf, and Alex J Smola. Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics*, 22(14):e49–e57, 2006.
- John Cai, Bill Cai, and Sheng Mei Shen. Sb-mtl: Score-based meta transfer-learning for cross-domain few-shot learning. *arXiv preprint arXiv:2012.01784*, 2020.
- Zhangjie Cao, Mingsheng Long, Jianmin Wang, and Michael I Jordan. Partial transfer learning with selective adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2724–2732, 2018a.
- Zhangjie Cao, Lijia Ma, Mingsheng Long, and Jianmin Wang. Partial adversarial domain adaptation. In *Proceedings of the European Conference on Computer Vision*, pp. 135–150, 2018b.
- Chaofan Chen, Xiaoshan Yang, Changsheng Xu, Xuhui Huang, and Zhe Ma. ECKPN: Explicit class knowledge propagation network for transductive few-shot learning. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 6596–6605, June 2021.
- Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *The International Conference on Learning Representations*, 2019.
- Zhengming Ding, Sheng Li, Ming Shao, and Yun Fu. *Graph Adaptive Knowledge Transfer for Unsupervised Domain Adaptation: 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part II*. Springer, Cham, 2018.
- Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. In *The International Conference on Learning Representations*, 2016.
- Bo Fu, Zhangjie Cao, Mingsheng Long, and Jianmin Wang. Learning to detect open classes for universal domain adaptation. In *Proceedings of the European Conference on Computer Vision*, pp. 567–583. Springer, 2020.
- Yuqian Fu, Yanwei Fu, and Yu-Gang Jiang. Meta-fdmixup: Cross-domain few-shot learning guided by labeled target data. In *Proceedings of the ACM International Conference on Multimedia*, pp. 5326–5334, 2021.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proceedings of International Conference on Machine Learning*, pp. 1180–1189. PMLR, 2015.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 999–1006, 2011.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2961–2969, 2017.
- Nathan Hilliard, Lawrence Phillips, Scott Howland, Artëm Yankov, Courtney D Corley, and Nathan O Hodas. Few-shot learning with metric-agnostic conditional embeddings. *arXiv preprint arXiv:1802.04376*, 2018.

- Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. CyCADA: Cycle-consistent adversarial domain adaptation. In *Proceedings of International Conference on Machine Learning*, pp. 1989–1998. PMLR, 2018.
- Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 7132–7141, 2018.
- Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1501–1510, 2017.
- Ashrafal Islam, Chun-Fu Chen, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, and Richard J Radke. Dynamic distillation network for cross-domain few-shot recognition with unlabeled data. In *Proceedings of the Advances in Neural Information Processing Systems*, 2021.
- Xin Jin, Cuiling Lan, Wenjun Zeng, Zhibo Chen, and Li Zhang. Style normalization and restitution for generalizable person re-identification. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 3143–3152, 2020.
- Xin Jin, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Style normalization and restitution for domain generalization and adaptation. *arXiv preprint arXiv:2101.00588*, 2021.
- Jongmin Kim, Taesup Kim, Sungwoong Kim, and Chang D Yoo. Edge-labeling graph neural network for few-shot learning. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 11–20, 2019.
- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pp. 554–561, 2013.
- Abhishek Kumar, Prasanna Sattigeri, Kahini Wadhawan, Leonid Karlinsky, Rogério Schmidt Feris, Bill Freeman, and Gregory W Wornell. Co-regularized alignment for unsupervised domain adaptation. In *Proceedings of the Advances in Neural Information Processing Systems*, 2018.
- Jogendra Nath Kundu, Nishank Lakkakula, and R Venkatesh Babu. UM-Adapt: Unsupervised multi-task adaptation using adversarial cross-task distillation. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1436–1445, 2019.
- Jogendra Nath Kundu, Naveen Venkat, R Venkatesh Babu, et al. Universal source-free domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4544–4553, 2020.
- Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- Kwonjoon Lee, Subhransu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex optimization. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 10657–10665, 2019.
- Hanwen Liang, Qiong Zhang, Peng Dai, and Juwei Lu. Boosting the generalization capability in cross-domain few-shot learning via noise-enhanced supervised autoencoder. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 9424–9434, 2021.
- Yann Lifchitz, Yannis Avrithis, Sylvaine Picard, and Andrei Bursuc. Dense classification and implanting for few-shot learning. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 9258–9267, 2019.
- Mingsheng Long, Jianmin Wang, Guiguang Ding, Jiaguang Sun, and Philip S Yu. Transfer feature learning with joint distribution adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2200–2207, 2013.
- Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *Proceedings of the Proceedings of Proceedings of International Conference on Machine Learning*, pp. 97–105. PMLR, 2015.

- Saeid Motiian, Marco Piccirilli, Donald A Adjeroh, and Gianfranco Doretto. Unified deep supervised domain adaptation and generalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5715–5725, 2017.
- Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *Proceedings of the International Conference on Machine Learning*, pp. 10–18. PMLR, 2013.
- Boris N Oreshkin, Pau Rodriguez, and Alexandre Lacoste. TADAM: Task dependent adaptive metric for improved few-shot learning. In *Proceedings of the Advances in Neural Information Processing Systems*, pp. 719–729, 2018.
- Pau Panareda Busto and Juergen Gall. Open set domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 754–763, 2017.
- Cheng Perng Phoo and Bharath Hariharan. Self-training for few-shot transfer across extreme task differences. In *The International Conference on Learning Representations*, 2020.
- Siyuan Qiao, Chenxi Liu, Wei Shen, and Alan L Yuille. Few-shot image recognition by predicting parameters from activations. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 7229–7238, 2018.
- Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. 2016.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149, 2016.
- Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *The International Conference on Learning Representations*, 2018.
- Kuniaki Saito and Kate Saenko. Ovanet: One-vs-all network for universal domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, 2021.
- Kuniaki Saito, Donghyun Kim, Stan Sclaroff, and Kate Saenko. Universal domain adaptation through self supervision. *Proceedings of the Advances in Neural Information Processing Systems*, 33, 2020.
- Victor Garcia Satorras and Joan Bruna Estrach. Few-shot learning with graph neural networks. In *The International Conference on Learning Representations*, 2018.
- Shiv Shankar, Vihari Piratla, Soumen Chakrabarti, Siddhartha Chaudhuri, Preethi Jyothi, and Sunita Sarawagi. Generalizing across domains via cross-gradient training. In *The International Conference on Learning Representations*, 2018.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Proceedings of the Advances in Neural Information Processing Systems*, pp. 4080–4090, 2017.
- Baochen Sun and Kate Saenko. Deep CORAL: Correlation alignment for deep domain adaptation. In *Proceedings of the European Conference on Computer Vision*, pp. 443–450. Springer, 2016.
- Jiamei Sun, Sebastian Lapuschkin, Wojciech Samek, Yunqing Zhao, Ngai-Man Cheung, and Alexander Binder. Explanation-guided training for cross-domain few-shot classification. In *Proceedings of the International Conference on Pattern Recognition*, pp. 7609–7616, 2021.
- Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE on Computer Vision and Pattern Recognition*, pp. 1199–1208, 2018.
- Shixiang Tang, Dapeng Chen, Lei Bai, Kaijian Liu, Yixiao Ge, and Wanli Ouyang. Mutual crf-gnn for few-shot learning. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 2329–2339, June 2021.

- Hung-Yu Tseng, Hsin-Ying Lee, Jia-Bin Huang, and Ming-Hsuan Yang. Cross-domain few-shot classification via learned feature-wise transformation. In *The International Conference on Learning Representations*, 2020.
- Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Simultaneous deep transfer across domains and tasks. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4068–4076, 2015.
- 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*, pp. 7167–7176, 2017.
- Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(11), 2008.
- Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE on Computer Vision and Pattern Recognition*, pp. 8769–8778, 2018.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *Proceedings of the Advances in Neural Information Processing Systems*, pp. 3637–3645, 2016.
- Riccardo Volpi, Hongseok Namkoong, Ozan Sener, John Duchi, Vittorio Murino, and Silvio Savarese. Generalizing to unseen domains via adversarial data augmentation. In *Proceedings of the Advances in Neural Information Processing Systems*, pp. 5339–5349, 2018.
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- Haoqing Wang and Zhi-Hong Deng. Cross-domain few-shot classification via adversarial task augmentation. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 1075–1081, 8 2021.
- Liu Yanbin, Lee Juho, Park Minseop, Kim Saehoon, Yang Eunho, Hwang Sungju, and Yang Yi. Learning to propagate labels: Transductive propagation network for few-shot learning. In *The International Conference on Learning Representations*, 2019.
- Baoyao Yang, Andy J. Ma, and Pong C. Yuen. Learning domain-shared group-sparse representation for unsupervised domain adaptation. *Pattern Recognition*, pp. S0031320318301614, 2018.
- Ling Yang, Liangliang Li, Zilun Zhang, Xinyu Zhou, Erjin Zhou, and Yu Liu. Dpgn: Distribution propagation graph network for few-shot learning. In *Proceedings of the IEEE Computer Vision and Pattern Recognition*, pp. 13390–13399, 2020.
- Kaichao You, Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Universal domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2720–2729, 2019.
- Jing Zhang, Zewei Ding, Wanqing Li, and Philip Ogunbona. Importance weighted adversarial nets for partial domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8156–8164, 2018.
- Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1452–1464, 2017.
- Yixiong Zou, Shanghang Zhang, Jianpeng Yu, Yonghong Tian, and José MF Moura. Revisiting mid-level patterns for cross-domain few-shot recognition. In *Proceedings of the ACM International Conference on Multimedia*, pp. 741–749, 2021.

## A DATASET DETAILS

We make evaluations on five public datasets that are commonly used in few-shot classification task: mini-ImageNet (Vinyals et al., 2016), CUB (Wah et al., 2011), Cars (Krause et al., 2013), Places (Zhou et al., 2017), and Plantae (Van Horn et al., 2018). Note that the procedure of all the datasets are exactly the same as baselines to ensure the fair comparison. The origins of these five datasets are introduced as follows:

Table 5: Statistical summary of five datasets.

Number of Category	Dataset				
	mini-ImageNet	CUB	Cars	Places	Plantae
<b>Training</b>	64	100	98	183	100
<b>Validation</b>	16	50	49	91	50
<b>Testing</b>	20	50	49	91	50

- **mini-ImageNet:** It was proposed by (Vinyals et al., 2016), it consists of 60,000 colour images of size  $84 \times 84$  with 100 classes, each class has 600 instances. The dataset process is the same with (Ravi & Larochelle, 2016).
- **CUB**<sup>4</sup>: It is a challenging dataset of 200 bird species (mostly North American). The same with the setting in (Hilliard et al., 2018), this dataset is divided into 100 classes for training, 50 for validation, and 50 for testing. Each image is again resized to  $84 \times 84$  pixels and put through the data augmentation process described previously to reduce over-fitting. The original of these five datasets are introduced as follows:
- **Car**<sup>5</sup>: It contains 16,185 images of 196 classes of cars. Classes are typically at the level of Make, Model, Year, e.g. 2012 Tesla Model S or 2012 BMW M3 coupe. The dataset split setting applied in cross-domain few-shot learning task is the same with (Tseng et al., 2020).
- **Places:** It contains 1.8 million images from 365 scene categories, where there are at most 5000 images per category. The dataset setting applied in cross-domain few-shot learning task is the same with (Tseng et al., 2020).
- **Plantae:** It is a sub-category of dataset iNaturalist<sup>6</sup> that was originally developed for species classification and detection. The dataset process setting is the same with (Tseng et al., 2020).

## B DIFFERENT NUMBERS OF WAYS

We consider a more practical situation that  $N_w$  may be different from that at the training stage. It also reflects the generalization ability of the model and results are shown in Table 6. Note that GNN (Satorras & Estrach, 2018) requires the number of ways to be the same while the training and testing, thus we evaluate with method MatchingNet (Vinyals et al., 2016) and RelationNet (Sung et al., 2018) (MNet and RNet for short). The model is trained on the datasets mini-ImageNet, Cars, Places, and Plantae and evaluated on the dataset CUB with different number of ways  $N_w$ . The results indicate that our proposed method are still capable of improving the generalization ability to the unseen domain with various numbers of ways. In addition, our proposed method consistently outperforms the baseline (Tseng et al., 2020) that has considered the domain-shift issue, which validates the superiority of our method.

## C EXTRA MODEL ANALYSIS

As illustrated in Figure 2, we employ the t-SNE algorithm (Van der Maaten & Hinton, 2008) to visualize features that obtained by the feature encoder “before/within/after” our MemREIN method, where each color represents one class. We take the GNN baseline under the leave-one-out setting

<sup>4</sup><http://www.vision.caltech.edu/visipedia/CUB-200.html>

<sup>5</sup>[https://ai.stanford.edu/~jkrause/cars/car\\_dataset.html](https://ai.stanford.edu/~jkrause/cars/car_dataset.html)

<sup>6</sup><https://www.inaturalist.org/>.

Table 6: Classification Accuracy (%) of our proposed method with different  $N_w$ . We consider the CUB dataset as the unseen domain under the leave-one-out setting.

5-shot	Classification Accuracy (%)			
	2-way	5-way	10-way	20-way
MNet (Vinyals et al., 2016)	78.46 $\pm$ 0.78%	51.92 $\pm$ 0.80%	38.22 $\pm$ 0.38%	26.17 $\pm$ 0.24%
MNet+LFT (Tseng et al., 2020)	83.88 $\pm$ 0.72%	61.41 $\pm$ 0.57%	45.69 $\pm$ 0.39%	32.81 $\pm$ 0.23%
<b>MNet+MemREIN (Ours)</b>	<b>88.68 <math>\pm</math> 0.68%</b>	<b>67.31 <math>\pm</math> 0.51%</b>	<b>49.22 <math>\pm</math> 0.34%</b>	<b>33.99 <math>\pm</math> 0.22%</b>
RNet (Sung et al., 2018)	84.25 $\pm$ 0.72%	62.13 $\pm$ 0.74%	47.15 $\pm$ 0.40%	34.52 $\pm$ 0.24%
RNet+LFT (Tseng et al., 2020)	85.44 $\pm$ 0.72%	64.99 $\pm$ 0.54%	49.90 $\pm$ 0.40%	37.20 $\pm$ 0.25%
<b>RNet+MemREIN (Ours)</b>	<b>89.12 <math>\pm</math> 0.66%</b>	<b>68.39 <math>\pm</math> 0.48%</b>	<b>52.85 <math>\pm</math> 0.32%</b>	<b>42.82 <math>\pm</math> 0.20%</b>

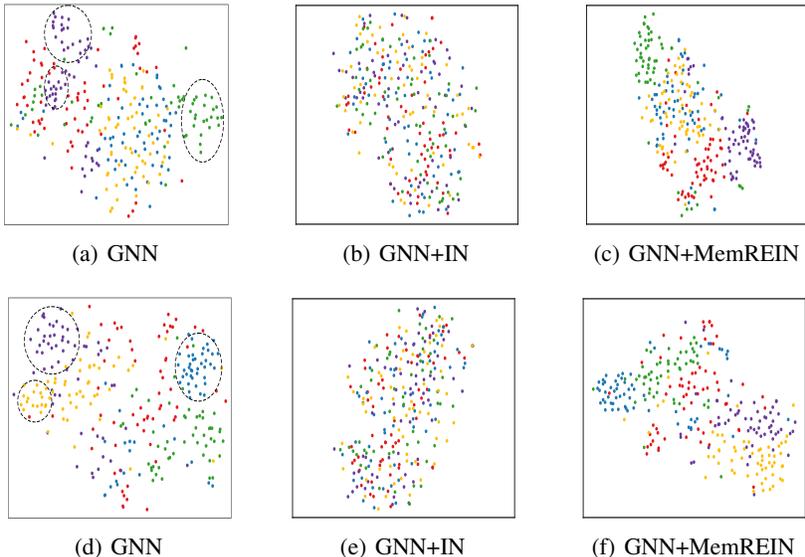


Figure 2: t-SNE visualization of sample features extracted by encoder.

on the dataset CUB as the example. We randomly select 5 categories with 60 samples of each category in the testing split of the dataset CUB. The first column indicates two examples of the features from conventional GNN baseline, The second column indicates the features that only applied the instance normalization operation, and the third column indicates the features that applied our proposed MemREIN method. As shown in the first column, there exists several rough clusters and but the boundaries are unclear. After instance normalization, the overall model generalization ability of features is enhanced. In comparison with the second column and the third column, the features learned by our method are more clustered and separable, which validates the effectiveness of our novel memorized restitution approach.

## D ABLATION STUDY

We carry out ablation studies to validate the effectiveness of different components in our proposed method. We compare with the GNN baseline under the leave-one-out setting (5-way 5-shot) and results are shown in Table 7. Comparing the results of Variant 1 and 2, it indicates that only applying the instance normalization operation results in the decrease of the accuracy. It is reasonable because the instance normalization operation will inevitably remove some discriminative useful information. In comparison with Variant 3 and 6, it validates the effectiveness of employing the memory bank on feature  $D$ . Comparing Variant 3, 6, and 7, it indicates that when employing memory bank on feature  $G$ , it would cause performance decrease. Experimentally, when applying the memory bank on the feature  $D$  and directly using feature  $G$  has the best performance.

Table 7: Ablation study on our method. ‘‘GNN+IN’’ indicates that we only employ the instance normalization strategy, ‘‘w/o  $\mathcal{L}_{rcl}^-$ ’’ indicates that we remove the  $\mathcal{L}_{rcl}^-$  term, and ‘‘w/o  $\mathcal{L}_{rcl}^+$ ’’ indicates that we remove the  $\mathcal{L}_{rcl}^+$  term, ‘‘GNN+MemREIN w/o MB’’ represents that we remove memory bank and directly use the feature map  $D$ , and ‘‘GNN+MemREIN ( $D\&G$ )’’ represents that the memory bank is operated both on feature map  $D$  and  $G$  (not shared).

5-way 5-shot		Classification Accuracy (%)			
Variant ID	Method	CUB	Cars	Places	Plantae
1	GNN (Satorras & Estrach, 2018)	69.26 $\pm$ 0.68%	48.91 $\pm$ 0.67%	72.59 $\pm$ 0.67%	58.36 $\pm$ 0.68%
2	GNN+IN	67.34 $\pm$ 0.66%	42.76 $\pm$ 0.75%	67.82 $\pm$ 0.73%	54.04 $\pm$ 0.69%
3	GNN+MemREIN	<b>77.54 <math>\pm</math> 0.62%</b>	<b>56.78 <math>\pm</math> 0.66%</b>	<b>78.84 <math>\pm</math> 0.66%</b>	<b>65.44 <math>\pm</math> 0.64%</b>
4	w/o $\mathcal{L}_{rcl}^-$	75.38 $\pm$ 0.63%	55.34 $\pm$ 0.72%	78.03 $\pm$ 0.68%	65.22 $\pm$ 0.64%
5	w/o $\mathcal{L}_{rcl}^+$	73.02 $\pm$ 0.62%	51.45 $\pm$ 0.64%	73.26 $\pm$ 0.66%	62.22 $\pm$ 0.64%
6	GNN+MemREIN w/o MB	75.98 $\pm$ 0.62%	54.64 $\pm$ 0.66%	74.86 $\pm$ 0.68%	64.08 $\pm$ 0.68%
7	GNN+MemREIN ( $D\&G$ )	76.02 $\pm$ 0.66%	55.26 $\pm$ 0.69%	78.08 $\pm$ 0.66%	64.84 $\pm$ 0.68%

Table 8: Performance study on the hyper-parameter  $\lambda$ .

5-way 5-shot		Classification Accuracy (%)	
GNN+MemREIN	CUB	Cars	
$\lambda = 0.01$	77.02 $\pm$ 0.62%	56.12 $\pm$ 0.66%	
$\lambda = 0.1$	<b>77.54 <math>\pm</math> 0.62%</b>	<b>56.78 <math>\pm</math> 0.66%</b>	
$\lambda = 0.5$	77.34 $\pm$ 0.62%	56.66 $\pm$ 0.66%	
$\lambda = 1$	76.78 $\pm$ 0.64%	56.22 $\pm$ 0.66%	

## E PERFORMANCE STUDY OF $\lambda$

We carry out performance study on the hyper-parameter  $\lambda$ . We take our method under the leave-one-out setting (5-way 5-shot) and dataset CUB and Cars as the example. We set four different values  $\lambda = \{0.01, 0.1, 0.5, 1\}$  and the results are shown in Table 8. It can be observed that when setting  $\lambda = 0.1$ , it can achieve the best performance.