

# Adaptive Weighted Co-Learning for Cross-Domain Few-Shot Learning

Abdullah Alchihabi<sup>1</sup>  
abdullahalchihabi@gmail.carleton.ca  
Marzi Heidari<sup>1</sup>  
marziheidari@gmail.carleton.ca  
Yuhong Guo<sup>1,2</sup>  
yuhong.guo@carleton.ca

<sup>1</sup> School of Computer Science  
Carleton University  
Ottawa, Canada  
<sup>2</sup> Canada CIFAR AI Chair  
Amii, Canada

---

## Abstract

Due to the availability of only a few labeled instances for the novel target prediction task and the significant domain shift between the well annotated source domain and the target domain, cross-domain few-shot learning (CDFSL) induces a very challenging adaptation problem. In this paper, we propose a simple Adaptive Weighted Co-Learning (AWCoL) method to address the CDFSL challenge by adapting two independently trained source prototypical classification models to the target task in a weighted co-learning manner. The proposed method deploys a weighted moving average prediction strategy to generate probabilistic predictions from each model, and then conducts adaptive co-learning by jointly fine-tuning the two models in an alternating manner based on the pseudo-labels and instance weights produced from the predictions. Moreover, a negative pseudo-labeling regularizer is further deployed to improve the fine-tuning process by penalizing false predictions. Comprehensive experiments are conducted on multiple benchmark datasets and the empirical results demonstrate that the proposed method produces state-of-the-art CDFSL performance.

## 1 Introduction

Deep learning has achieved great success on a wide set of computer vision tasks, ranging from image classification to image segmentation and object detection. Such success however has been contingent on the availability of a large set of annotated training instances, which induces significant data annotation cost. To alleviate this annotation burden, Few-Shot Learning (FSL) methods have been developed to exploit a base set of classes with a large number of labeled instances to help train deep prediction models for target tasks over novel classes, where only a few instances are labeled for each class. Standard FSL methods are designed for an in-domain setting, where data for the base classes and the novel task are from the same domain, and their performance degrades as the dissimilarity between the base dataset (source domain) and the novel target task (target domain) increases [9], which raises significant demands for cross-domain few-shot learning.

Cross-Domain Few-Shot Learning (CDFSL) expands the FSL study by leveraging labeled base datasets from a remote source domain to help FSL in an annotation-strenuous target domain. CDFSL is much more challenging than its in-domain counterpart due to the limited few-shot labeled instances for the novel target prediction task and the significant domain shift between the source and target domains, and has just started gaining attention from the research community. Some approaches have been developed to tackle CDFSL by employing data augmentation and data generation techniques to increase the number of available labeled instances [13, 56], while several others have deployed ensemble learning or hierarchical variational memory techniques to learn diverse features and facilitate model adaptation from the source domain to the target domain [10, 9]. However, these methods either have limited scalability for higher-shot learning problems [56], induce significant computational cost [10, 9], or require a large set of unlabeled target-domain samples to be available during source-domain training [13].

In this paper, we propose a simple but novel Adaptive Weighted Co-Learning (AWCoL) method for CDFSL. The proposed method jointly fine-tunes two simple prototypical classification models that were independently pre-trained on the source-domain dataset for the target few-shot prediction task in a weighted adaptive co-learning manner. Specifically, for each model we propose to use a weighted moving average (WMA) strategy to generate robust probabilistic predictions for the query instances of the target task, where the class prototypes are computed using the support instances. Pseudo-labels and the corresponding confidence weights for the query instances can be determined based on the average probabilistic predictions produced by the two models. The two prototypical models are then fine-tuned in an alternating fashion on the pseudo-labeled query instances by minimizing a weighted cross-entropy loss, aiming to effectively exploit the more diverse unlabeled queries for stable co-learning. Moreover, we also adopt a negative pseudo-label regularizer to further assist the fine-tuning process by pushing the predictions away from the negative pseudo-labels. To evaluate the proposed method, we conduct comprehensive experiments on eight benchmark datasets for CDFSL. The proposed AWCoL method demonstrates superior performance compared to the existing state-of-the-art CDFSL methods.

## 2 Related Works

### 2.1 Few-Shot Learning

Standard in-domain few-shot learning (FSL) methods can be broadly categorized into the following four main groups: generative and augmentation-based methods, transfer learning methods, meta-learning methods, and metric learning methods. The generative and augmentation-based methods generate new instances to increase the training set diversity [10, 19, 24, 76, 39]. The transfer learning methods focus on adapting models pre-trained on the base dataset to the novel target task with fine-tuning techniques [8, 7, 8, 14, 55]. Meta-learning approaches learn to perform few-shot learning by simulating few-shot tasks on the annotated base dataset [6, 16]. Metric learning methods aim to exploit the base dataset to induce good similarity/distance metrics [20, 25, 27, 28, 54]. Among these methods, ProtoNet [27] represents each class using a prototype and learns a metric space in which instances are classified based on their distances to the class prototypes. RelationNet [28], GNN [25] and Transductive Propagation Network (TPN) [20] utilize similarities between the support instances and query instances to perform few-shot learning.

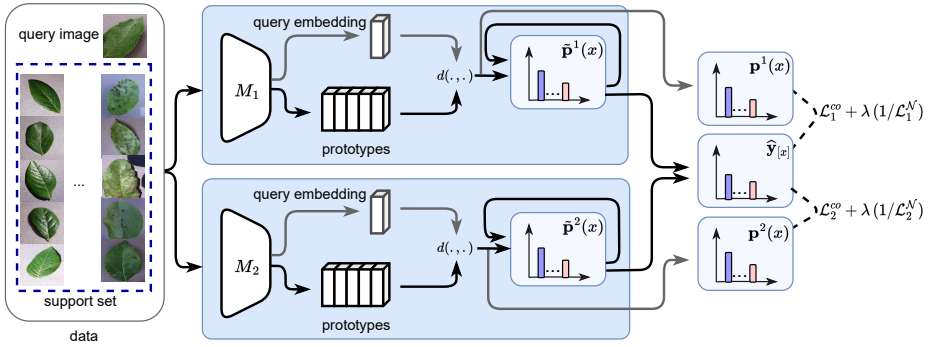


Figure 1: The framework of the proposed Adaptive Weighted Co-Learning (AWCoL) method with two prototypical models,  $M_1$  and  $M_2$ , that are pre-trained independently in the source domain: At each fine-tuning iteration in the target domain, each model generates probabilistic predictions for the query instances using a weighted moving average strategy. The predictions from both models are then combined to determine the positive pseudo-labels, negative pseudo-labels, and adaptive weights for the query instances, which are used to fine-tune the two models in the target domain in an adaptive co-learning manner.

## 2.2 Cross-Domain Few-Shot Learning

By broadening the applicable domains, cross-domain few-shot learning (CDFSL) has recently started drawing more attention [9, 23]. Most works try to induce better generalizable representations for CDFSL. Li *et al.* proposed a ranking distance calibration with fine-tuning (RDC-FT) method that constructs a non-linear subspace to reduce task-irrelevant features [17]. The Hierarchical Variational neural Memory framework (HVM) learns hierarchical prototypes with variational inference [4]. Data augmentation and data generation techniques have also been used for CDFSL [13, 37]. Adversarial Task Augmentation (ATA) has been deployed to improve meta-learning models for cross-domain few-shot classification [36]. Fu *et al.* proposed a meta Style Adversarial training (StyleAdv) method that perturbs the original style of samples using a novel style attack method, enabling the model to generalize to challenging styles and inducing cross-domain robustness [6]. In addition, self-supervised learning based methods have also been developed for CDFSL [2]. Das *et al.* proposed a Contrastive learning and feature selection system (ConFeSS) to bridge the domain shift with self-supervised representation learning [2]. Zhou *et al.* proposed a two-branch Local-global Distillation Prototypical Network (LDP-net) that leverages cross-domain transferable semantic features by enforcing local-global semantic consistency between the two branches through knowledge distillation [11].

A few other works [18, 21, 29, 32] have employed alternative CDFSL settings such as Meta-Dataset [29] and Meta-Album [32] CDFSL settings where models are trained on several source-domain datasets and tested on multiple target-domain datasets. In this work, we focus on the CDFSL setting in [9] as it is the most widely studied CDFSL setting.

## 3 The Proposed Method

**Problem Setup** CDFSL aims to train a model in an annotation-abundant source domain, and then adapt/fine-tune the model for the novel few-shot prediction task in the target do-

main. The source domain and the target domain are assumed to have different distributions in the input feature space ( $\mathcal{P}_s \neq \mathcal{P}_t$ ) and disjoint classes ( $\mathcal{Y}_s \cap \mathcal{Y}_t = \emptyset$ ) in the output label space. In the target domain, the model has access to a labelled support set  $S = \{(x_i, y_i)\}_{i=1}^{N_s}$  and is evaluated on a query set  $Q = \{(x_i, y_i)\}_{i=1}^{N_q}$ , where  $N_s$  and  $N_q$  are the sizes of the support and query sets respectively. The support set typically consists of  $K$  instances from each of the  $N$  classes, resulting a total of  $N_s = N \times K$  labeled support instances. This setup is commonly dubbed as N-way K-shot learning.

In this section, we present the proposed Adaptive Weighted Co-Learning (AWCoL) method to address the CDFSL problem. The overall framework of AWCoL is illustrated in Figure 1. The method provides a simple adaptive co-learning mechanism for fine-tuning two models that were independently pre-trained in the source domain. In particular, we consider pre-training two simple prototypical FSL models in the source domain. The models are then fine-tuned for the target FSL task in the target domain using an adaptive weighted co-learning procedure. In each iteration of the fine-tuning, a weighted moving average (WMA) strategy is applied to each of the two models independently to generate probabilistic predictions for the unlabeled query instances. The predictions from the two models are then combined to produce pseudo-labels and adaptive weights for the query instances, which are subsequently used to fine-tune the models in an adaptive co-learning manner. In addition, negative pseudo-labels generated from the predictions can also be used to enhance the fine-tuning process by pushing the models away from the predicted negative class labels. Overall, the two models are fine-tuned in an alternating fashion to yield a stable co-learning process, which will be elaborated below.

### 3.1 Source-domain Pre-training

We use a simple prototypical FSL model [27] as our base classification model, due to its simplicity and efficiency. In the well-annotated source domain, we pre-train two classic prototypical few-shot learning models independently using two different sets of randomly sampled N-way K-shot FSL tasks, which we refer to as  $M_1$  and  $M_2$ . Each of the two models has its own feature encoder  $f$  parameterized by  $\theta$  that maps an input image  $x$  to a feature embedding vector. For each model, given each sampled training FSL task with a support set  $S$  and a query set  $Q$ , a prototype embedding vector  $\mathbf{c}_n$  can be computed from the support instances for each class  $n$  using the feature extractor  $f$ , such that:

$$\mathbf{c}_n = \frac{1}{K} \sum_{(x,y) \in S_n} f_{\theta}(x), \quad (1)$$

where  $S_n$  denotes the set of  $K$  support instances from class  $n$ . With the class prototypes, the prototypical classifier predicts the class probability vector  $\mathbf{p}(x) = [\mathbf{p}_1(x), \dots, \mathbf{p}_N(x)]$  for each query instance  $x$  by calculating the negative distances between  $f_{\theta}(x)$  and each class prototype vector and then normalizing them using the softmax function, such that:

$$\mathbf{p}_j(x) = \frac{\exp(-d(f_{\theta}(x), \mathbf{c}_j))}{\sum_{n=1}^N \exp(-d(f_{\theta}(x), \mathbf{c}_n))}, \quad (2)$$

where  $d(\cdot, \cdot)$  is a distance function and  $\mathbf{p}_j(x)$  is the predicted probability of instance  $x$  belonging to class  $j$ .

Each model can be trained consecutively on the set of randomly sampled FSL tasks in the source domain by minimizing the following cross-entropy loss over the query instances:

$$\mathcal{L}_{CE}(Q) = \sum_{(x,y) \in Q} \ell_{CE}(\mathbf{1}_y, \mathbf{p}(x)), \quad (3)$$

where  $\ell_{CE}$  is the cross-entropy loss function,  $\mathbf{p}(x)$  and  $\mathbf{1}_y$  are the predicted class probability vector and the ground-truth one-hot label indicator vector (with value 1 at the  $y$ -th entry) respectively for the query instance  $x$ . The output of the source-domain pre-training step is two independently trained prototypical classification models,  $M_1$  and  $M_2$ , with separate feature encoders  $f_1$  and  $f_2$  parameterized by  $\theta_1$  and  $\theta_2$  respectively.

### 3.2 Weighted Moving Average based Probabilistic Prediction Generation

Given the significant domain shift between the source domain and target domain in the CDFSL problem, it is important to adapt the models pre-trained in the source domain to the novel FSL task in the target domain through fine-tuning. Due to the limited number of labeled support instances,  $\mathcal{S}$ , available for the FSL target task, we propose to exploit the unlabeled query instances,  $\mathcal{Q}$ , with pseudo-labels for transductive model fine-tuning, aiming to increase training data diversity and avoid over-adaptation.

The pseudo-labels for query instances will be produced based on the predictions of the two models,  $M_1$  and  $M_2$ . In order to generate robust predictions across fine-tuning iterations, we adopt a weighted moving average (WMA) strategy to produce probabilistic predictions for the query instances using each of the two models separately. Specifically, let  $M_m$  denote either one of the two models, such that  $m \in \{1, 2\}$ . Then at each fine-tuning iteration, we calculate the prototype vector for each class from the support instances with the feature extractor  $f_m$  in the given model  $M_m$  using Eq. (1). Next for each query instance  $x$ , we calculate its class prediction probability vector  $\mathbf{p}^m(x)$  using Eq. (2) under the current model  $M_m$ . Instead of using this predicted probability vector directly, we maintain a WMA class prediction probability vector  $\tilde{\mathbf{p}}^m(x)$  for each query instance  $x$  across the fine-tuning iterations, which can be updated in each iteration using the new predictions  $\mathbf{p}^m(x)$  as follows:

$$\tilde{\mathbf{p}}^m(x) = (1 - \alpha_m) \tilde{\mathbf{p}}^m(x) + \alpha_m \mathbf{p}^m(x) \quad (4)$$

where  $\alpha_m \in (0, 1)$  is a trade-off hyper-parameter that controls the combination weights between the class prediction probability vector  $\mathbf{p}^m$  of the current iteration and the WMA class prediction probability vector  $\tilde{\mathbf{p}}^m$  from the previous iteration, and hence determines the updating degree. With such a WMA update, we expect to maintain stable improvements for pseudo-label predictions over the query instances and avoid local oscillations.

To further improve the stability of the WMA prediction generation, we employ a rectified annealing schedule to update the hyper-parameter  $\alpha_m$  in each iteration as follows:

$$\alpha_m = \max(\alpha_{\min}, \gamma \alpha_m), \quad (5)$$

where  $\alpha_{\min}$  provides a lower bound for  $\alpha_m$ , and  $\gamma \in (0, 1)$  is a reduction ratio hyper-parameter that gradually reduces the  $\alpha_m$  value across iterations. This annealing schedule allows the method to perform larger updates to the WMA class prediction probability vectors with larger  $\alpha_m$  values in the early fine-tuning iterations, while decreasing the updating degree with smaller  $\alpha_m$  in later iterations of fine-tuning, aiming to help the fine-tuning process reach convergence.

### 3.3 Weighted Adaptive Co-Learning

With the WMA class prediction probability vectors for the query instances produced by  $M_1$  and  $M_2$ , we propose to fine-tune the two models together using a weighted adaptive co-

learning procedure by effectively integrating their predictions. Specifically, we simply take a softmax rescaled average of the WMA predictions from the two models to produce the co-prediction probability vector  $\tilde{\mathbf{p}}^{\text{co}}(x)$  for each query instance  $x$ , such that:

$$\tilde{\mathbf{p}}^{\text{co}}(x) = \text{softmax}(\text{avg}(\tilde{\mathbf{p}}^1(x), \tilde{\mathbf{p}}^2(x))). \quad (6)$$

Here we apply the softmax function to rescale the average probability vector towards less skewed predictions, aiming to avoid early stage poor local optima for the co-learning procedure and help stabilize the fine-tuning process.

The pseudo-labels for the query instances can then be determined based on the integrated co-prediction probability vectors. Specifically, for a query instance  $x$ , its pseudo-label  $\hat{y}_{[x]}$  will be a scalar index indicating the class with the highest prediction probability:

$$\hat{y}_{[x]} = \text{argmax}_{n \in \{1, \dots, N\}} \tilde{\mathbf{p}}_n^{\text{co}}(x) \quad (7)$$

where  $N$  is the total number of classes for the target FSL task. For convenience, we further transform the pseudo-label  $\hat{y}_{[x]}$  into a one-hot pseudo-label indicator vector  $\hat{\mathbf{y}}_{[x]} \in \{0, 1\}^N$  that has a single 1 value at its  $\hat{y}_{[x]}$ -th entry.

By producing pseudo-labels from the integrated predictions of the two models, we expect to yield more accurate pseudo-labels than using each individual model alone. Nevertheless, noisy labels, i.e., mistakes, are still unavoidable among the predicted pseudo-labels. To alleviate the negative impact of the potential pseudo-label noise, we further propose to calculate a dynamic and adaptive weight for each unlabeled query instance  $x$  based on the current co-predictions:

$$w_{[x]} = \max_{n \in \{1, \dots, N\}} \tilde{\mathbf{p}}_n^{\text{co}}(x) \quad (8)$$

This adaptive weighting mechanism allows us to assign different weights to different query instances based on the co-prediction confidence of the generated pseudo-labels, such that the query instances with more confident pseudo-labels—less likely to be noise—get larger weights.

With the co-predicted pseudo-labels and adaptive weights for the query instances, we fine-tune the two models using an alternating co-learning procedure; i.e., alternatively updating each model across co-learning iterations. In a given iteration, for the chosen model  $M_m$  ( $m \in \{1, 2\}$ ), we update its model parameters  $\theta_m$  by minimizing the following adaptive weighted cross-entropy loss on the pseudo-labeled query instances:

$$\mathcal{L}_m^{\text{co}} = \frac{1}{\sum_{x \in Q} w_{[x]}} \sum_{x \in Q} w_{[x]} \ell_{CE}(\hat{\mathbf{y}}_{[x]}, \mathbf{p}^m(x)) \quad (9)$$

where  $\mathbf{p}^m(x)$  is the prediction probability vector for instance  $x$  using the prototypical model  $M_m$  via Eq. (2).

### 3.4 Negative Pseudo-Label Regularization

In addition to exploiting the predicted positive pseudo-labels, we further propose to employ the negative pseudo-labels—i.e., the other class labels except the predicted ones—to regularize the co-learning process. Specifically, for a query instance  $x$ , in addition to determining the predicted positive pseudo-label  $\hat{y}_{[x]}$  using Eq. (7) from the integrated co-prediction probability vector  $\tilde{\mathbf{p}}^{\text{co}}(x)$ , we treat all the other classes  $\{1, \dots, N\} \setminus \{\hat{y}_{[x]}\}$  as candidate negative

labels for  $x$  and randomly choose one to use as a negative pseudo-label:

$$\hat{y}_{[x]}^N = \text{Rand}(\{1, \dots, N\} \setminus \{\hat{y}_{[x]}\}) \quad (10)$$

Same as before, we can transform this negative pseudo-label index into a one-hot negative pseudo-label indicator vector  $\hat{y}_{[x]}^N \in \{0, 1\}^N$  with a single 1 at its  $\hat{y}_{[x]}^N$ -th entry.

For a multi-class classification problem, only one class label is the true label for each given query instance. Our hypothesis is that even when the predicted positive pseudo-label is a false positive label, a randomly selected negative pseudo-label can still possibly be a true negative label. Hence by minimizing the likelihood of the chosen negative pseudo-labels on query instances, one can push the model away from making true negative predictions and help alleviate the negative impact of the prediction noise in the positive pseudo-labels. To encode this hypothesis, we incorporate the selected negative pseudo-labels into the co-learning process by *maximizing* the following adaptive weighted cross-entropy loss on the query instances for each model  $M_m$  ( $m \in \{1, 2\}$ ):

$$\mathcal{L}_m^N = \frac{1}{\sum_{x \in Q} w_{[x]}} \sum_{x \in Q} w_{[x]} \ell_{CE}(\hat{y}_{[x]}^N, \mathbf{p}^m(x)) \quad (11)$$

where the adaptive weight  $w_{[x]}$  and the label prediction probability vector  $\mathbf{p}^m(x)$  are same as the ones used for positive pseudo-labels in Eq. (9).

### 3.5 Alternating Co-Learning Procedure

By integrating both the co-learning minimization loss in Eq. (9) and the negative pseudo-label based maximization loss in Eq. (11) together, we deploy the following adaptive co-learning loss to fine-tune each model  $M_m$  ( $m \in \{1, 2\}$ ):

$$\mathcal{L}_m = \mathcal{L}_m^{co} + \lambda (1/\mathcal{L}_m^N) \quad (12)$$

where  $\lambda$  is a trade-off hyper-parameter that controls the relative contribution of the two loss terms. Minimizing this total loss will simultaneously minimize the positive pseudo-label based weighted adaptive co-learning loss  $\mathcal{L}_m^{co}$  and maximize the negative pseudo-label based weighted cross-entropy loss  $\mathcal{L}_m^N$ . In order to maintain a stable fine-tuning process and prevent oscillating co-prediction updates between the two models, we utilize an alternating updating mechanism to fine-tune the two models with the integrated total loss. In particular, we consecutively update each model for  $\beta$  iterations while keeping the other model fixed, and then switch to update the other model similarly. As a result, when calculating the co-prediction probability vector  $\tilde{\mathbf{p}}^{co}(x)$  in each iteration, only one model’s WMA predictions will be updated while anchoring on the fixed WMA predictions of the other model.

## 4 Experiments

### 4.1 Experimental Setup

**Datasets** We conducted comprehensive experiments to evaluate the performance of AW-CoL on eight CDFSL benchmark datasets. Following the standard CDFSL setting, we used Mini-ImageNet [54] as the source domain dataset in all experiments and employed the following eight datasets separately as the target domain datasets: ChestX [37], CropDiseases[2], ISIC [60], EuroSAT [12], Places [40], Planatae [33], Cars [15] and CUB [35]. We used the same train/validation/test splits as in [9] to ensure fair evaluation and comparison.

Table 1: Mean classification accuracy (95% confidence interval in brackets) on 8 target domain datasets for cross-domain 5-way 5-shot classification. “\*” and “†” indicate the results reported in [14] and [36] respectively. The best results are indicated in bold font and the second best results are underlined.

	ChestX	CropDisea.	ISIC	EuroSAT	Places	Planatae	Cars	CUB
MatchingNet* [14]	22.40 <sub>(0.70)</sub>	66.39 <sub>(0.78)</sub>	36.74 <sub>(0.53)</sub>	64.45 <sub>(0.63)</sub>	—	—	—	—
MAML* [14]	23.48 <sub>(0.96)</sub>	78.05 <sub>(0.68)</sub>	40.13 <sub>(0.58)</sub>	71.70 <sub>(0.72)</sub>	—	—	—	—
ProtoNet* [22]	24.05 <sub>(1.01)</sub>	79.72 <sub>(0.67)</sub>	39.57 <sub>(0.57)</sub>	73.29 <sub>(0.71)</sub>	58.54 <sub>(0.68)</sub>	46.80 <sub>(0.65)</sub>	41.74 <sub>(0.72)</sub>	55.51 <sub>(0.68)</sub>
MetaOpt* [16]	22.53 <sub>(0.91)</sub>	68.41 <sub>(0.73)</sub>	36.28 <sub>(0.50)</sub>	64.44 <sub>(0.73)</sub>	—	—	—	—
RelationNet† [28]	24.07 <sub>(0.20)</sub>	72.86 <sub>(0.40)</sub>	38.60 <sub>(0.30)</sub>	65.56 <sub>(0.40)</sub>	64.25 <sub>(0.40)</sub>	42.71 <sub>(0.30)</sub>	40.46 <sub>(0.40)</sub>	56.77 <sub>(0.40)</sub>
GNN† [25]	23.87 <sub>(0.20)</sub>	83.12 <sub>(0.40)</sub>	42.54 <sub>(0.40)</sub>	78.69 <sub>(0.40)</sub>	70.91 <sub>(0.50)</sub>	48.51 <sub>(0.40)</sub>	43.70 <sub>(0.40)</sub>	62.87 <sub>(0.50)</sub>
TPN † [20]	22.17 <sub>(0.20)</sub>	81.91 <sub>(0.50)</sub>	45.66 <sub>(0.30)</sub>	77.22 <sub>(0.40)</sub>	71.39 <sub>(0.40)</sub>	50.96 <sub>(0.40)</sub>	44.54 <sub>(0.40)</sub>	63.52 <sub>(0.40)</sub>
ATA† [36]	24.43 <sub>(0.20)</sub>	90.59 <sub>(0.30)</sub>	45.83 <sub>(0.30)</sub>	83.75 <sub>(0.40)</sub>	75.48 <sub>(0.40)</sub>	55.08 <sub>(0.40)</sub>	49.14 <sub>(0.40)</sub>	66.22 <sub>(0.50)</sub>
STARTUP [23]	26.94 <sub>(0.44)</sub>	93.02 <sub>(0.45)</sub>	47.22 <sub>(0.61)</sub>	82.29 <sub>(0.60)</sub>	—	—	—	—
HVM [9]	<b>27.15</b> <sub>(0.45)</sub>	87.65 <sub>(0.35)</sub>	42.05 <sub>(0.34)</sub>	74.88 <sub>(0.45)</sub>	—	—	—	—
ConFeSS [8]	<u>27.09</u>	88.88	48.85	84.65	—	—	—	—
RDC-FT [17]	25.48 <sub>(0.20)</sub>	93.55 <sub>(0.30)</sub>	49.06 <sub>(0.30)</sub>	84.67 <sub>(0.30)</sub>	74.65 <sub>(0.40)</sub>	60.63 <sub>(0.40)</sub>	53.75 <sub>(0.50)</sub>	67.77 <sub>(0.40)</sub>
LDP-net [18]	26.88 <sub>(0.46)</sub>	91.89 <sub>(0.50)</sub>	48.44 <sub>(0.67)</sub>	84.05 <sub>(0.66)</sub>	75.47 <sub>(0.73)</sub>	59.64 <sub>(0.77)</sub>	53.06 <sub>(0.82)</sub>	<u>73.34</u> <sub>(0.75)</sub>
StyleAdv [6]	26.24 <sub>(0.35)</sub>	<u>96.51</u> <sub>(0.28)</sub>	<u>53.05</u> <sub>(0.54)</sub>	<u>91.64</u> <sub>(0.43)</sub>	<u>79.35</u> <sub>(0.61)</sub>	<u>64.10</u> <sub>(0.64)</sub>	<u>56.44</u> <sub>(0.68)</sub>	70.90 <sub>(0.63)</sub>
AWCoL	24.50 <sub>(0.33)</sub>	<b>99.59</b> <sub>(0.10)</sub>	<b>58.75</b> <sub>(0.60)</sub>	<b>96.76</b> <sub>(0.27)</sub>	<b>92.56</b> <sub>(0.36)</sub>	<b>67.31</b> <sub>(0.64)</sub>	<b>62.94</b> <sub>(0.75)</sub>	<b>86.23</b> <sub>(0.52)</sub>

**Implementation Details** We use ResNet10 [14] as the backbone network for the two prototypical classification models in AWCoL. The distance function employed in AWCoL is the squared L2 norm. We train each model in the source domain for 400 iterations with 100 randomly sampled FSL tasks and 15 query instances per class, using the Adam optimizer with a learning rate of 1e-3. We evaluate our proposed AWCoL method on 600 randomly selected few-shot tasks in each target domain. For each task, we randomly sample 5 classes and randomly select 15 images per class as the query set. We fine-tune the models for a total of 100 iterations for each novel target task using the Adam optimizer with a learning rate of 1e-3. Hyper-parameters  $\lambda$ ,  $\alpha_{min}$ ,  $\alpha_0$  (the initial value for  $\alpha_m$ ),  $\gamma$ , and  $\beta$  take the values of 1e-2, 0.1, 0.5, 0.99 and 5 respectively.

## 4.2 Comparison Results

### 4.2.1 Learning with Few Shots

We evaluate the proposed AWCoL method on the standard cross-domain 5-way 5-shot tasks. We compare AWCoL with a set of representative FSL baselines (MatchingNet [14], MAML [9], ProtoNet [22], RelationNet [28], MetaOpt [16], GNN [25] and TPN [20]) as well as 7 state-of-the-art CDFSL methods (ATA [36], STARTUP [23], HVM [9], ConFeSS [8], RDC-FT [17], StyleAdv [6] and LDP-net [18]). The comparison results are reported in Table 1, where the top section of the table presents the results of the standard FSL methods and the bottom section presents the results of the CDFSL methods.

The table shows that most CDFSL methods naturally outperform the FSL methods on all the datasets given their ability to handle the significant domain shift between the source and target domains. The proposed method AWCoL outperforms all the FSL and CDFSL methods on all the datasets except ChestX. Only on ChestX, AWCoL is slightly outperformed by a few other CDFSL methods. Nevertheless, the performance gains produced by AWCoL over all the other methods are notable, exceeding 5%, 13% and 12% on EuroSAT, Places and CUB



Table 2: Mean classification accuracy (95% confidence interval in brackets) on 4 target domain datasets for cross-domain 5-way 20-shot and 50-shot classification. “\*\*\*” indicates results reported in [9]. The best results are in bold font and the second best results are underlined.

	ChestX		CropDiseases		ISIC		EuroSAT	
	20-shot	50-shot	20-shot	50-shot	20-shot	50-shot	20-shot	50-shot
MatchingNet* [34]	23.61 (0.86)	22.12 (0.88)	76.38 (0.67)	58.53 (0.73)	45.72 (0.53)	54.58 (0.65)	77.10 (0.57)	54.44 (0.67)
MAML* [6]	27.53 (0.43)	—	89.75 (0.42)	—	52.36 (0.57)	—	81.95 (0.55)	—
ProtoNet* [27]	28.21 (1.15)	29.32 (1.12)	88.15 (0.51)	90.81 (0.43)	49.50 (0.55)	51.99 (0.52)	82.27 (0.57)	80.48 (0.57)
MetaOpt* [16]	25.53 (1.02)	29.35 (0.99)	82.89 (0.54)	91.76 (0.38)	49.42 (0.60)	54.80 (0.54)	79.19 (0.62)	83.62 (0.58)
RelationNet* [28]	26.63 (0.92)	28.45 (1.20)	80.45 (0.64)	85.08 (0.53)	41.77 (0.49)	49.32 (0.51)	74.43 (0.66)	74.91 (0.58)
MatchingNet+FWT* [34]	23.23 (0.37)	23.01 (0.34)	74.90 (0.71)	75.68 (0.78)	32.01 (0.48)	33.17 (0.43)	63.38 (0.69)	62.75 (0.76)
ProtoNet+FWT* [34]	26.87 (0.43)	30.12 (0.46)	85.82 (0.51)	87.17 (0.50)	43.78 (0.47)	49.84 (0.51)	75.74 (0.70)	78.64 (0.57)
RelationNet+FWT* [34]	26.75 (0.41)	27.56 (0.40)	78.43 (0.59)	81.14 (0.56)	43.31 (0.51)	46.38 (0.53)	69.40 (0.64)	73.84 (0.60)
STARTUP [23]	33.19 (0.46)	<u>36.91</u> (0.50)	<u>97.51</u> (0.21)	<u>98.45</u> (0.17)	58.63 (0.58)	64.16 (0.58)	89.26 (0.43)	91.99 (0.36)
HVM [9]	30.54 (0.47)	32.76 (0.46)	95.13 (0.35)	97.83 (0.33)	54.97 (0.35)	61.71 (0.32)	84.81 (0.34)	87.16 (0.35)
ConFeSS [9]	<u>33.57</u>	<b>39.02</b>	95.34	97.56	<u>60.10</u>	<u>65.34</u>	<u>90.40</u>	<u>92.66</u>
AWCoL	<b>34.13</b> (0.47)	36.72 (0.48)	<b>99.93</b> (0.03)	<b>99.98</b> (0.02)	<b>72.40</b> (0.52)	<b>73.40</b> (0.51)	<b>98.98</b> (0.12)	<b>98.73</b> (0.16)

respectively. This highlights the significant performance gains obtained by our proposed method over the existing state-of-the-art CDFSL methods.

## 4.2.2 Learning with Higher Shots

We also conducted experiments to evaluate the performance of the proposed AWCoL method on higher-shot tasks in the target domain against a set of representative FSL baselines (MatchingNet [34], MAML [6], ProtoNet [27], RelationNet [28] and MetaOpt [16]) as well as 4 state-of-the-art CDFSL methods (FWT [34], STARTUP [23], HVM [9] and ConFeSS [9]). FWT has been applied jointly with three standard FSL methods: MatchingNet, ProtoNet and RelationNet. In particular, we evaluate our proposed AWCoL on 5-way 20-shot and 5-way 50-shot learning tasks on four target-domain datasets: ChestX, CropDiseases, ISIC and EuroSAT. The results are presented in Table 2, where the top section of the table reports results of the standard FSL methods and the bottom section reports results of the CDFSL methods.

Similar to the comparisons on 5-shot tasks, most CDFSL methods outperform the in-domain FSL methods across all datasets for both 20-shot and 50-shot tasks. It is also clear from the table that AWCoL outperforms all the other FSL and CDFSL methods on the CropDiseases, ISIC and EuroSAT datasets for both 20-shot and 50-shot learning tasks. The performance gains are notable, exceeding 12% and 8% on the ISIC dataset for the 20-shot and 50-shot learning tasks respectively. AWCoL is slightly outperformed on ChestX by ConFess and STARTUP in the case of 50-shot tasks. Nevertheless our proposed method still obtains the best result and third best result on ChestX for the cases of 20-shot and 50-shot tasks respectively.

## 4.3 Ablation Study

In order to investigate the contribution of each component of the proposed method, we conducted an ablation study to compare AWCoL with its seven variants: (1) “–w/o Co-Learn.”, which drops the adaptive co-learning component and fine-tunes each model independently using its own WMA predictions. (2) “–w/o Alt. Update”, which updates the two models

Table 3: Ablation study results in terms of mean classification accuracy (95% confidence interval within brackets) for cross-domain 5-way 5-shot classification tasks.

	ChestX	CropDisea.	ISIC	EuroSAT	Places	Planatae	Cars	CUB
AWCoL	24.50 <sub>(0.33)</sub>	<b>99.59</b> <sub>(0.10)</sub>	<b>58.75</b> <sub>(0.60)</sub>	96.76 <sub>(0.27)</sub>	92.56 <sub>(0.36)</sub>	<b>67.31</b> <sub>(0.64)</sub>	<b>62.94</b> <sub>(0.75)</sub>	<b>86.23</b> <sub>(0.52)</sub>
–w/o Co-Learn.	24.08 <sub>(0.42)</sub>	81.65 <sub>(0.62)</sub>	42.93 <sub>(0.60)</sub>	72.72 <sub>(0.72)</sub>	69.15 <sub>(0.73)</sub>	46.16 <sub>(0.65)</sub>	42.60 <sub>(0.72)</sub>	56.35 <sub>(0.73)</sub>
–w/o Alt. Update	<b>24.79</b> <sub>(0.40)</sub>	84.31 <sub>(0.63)</sub>	44.42 <sub>(0.61)</sub>	75.34 <sub>(0.69)</sub>	68.91 <sub>(0.70)</sub>	49.40 <sub>(0.67)</sub>	44.53 <sub>(0.77)</sub>	58.63 <sub>(0.72)</sub>
–w/o WMA	24.47 <sub>(0.44)</sub>	84.19 <sub>(0.62)</sub>	43.96 <sub>(0.63)</sub>	75.31 <sub>(0.70)</sub>	68.85 <sub>(0.71)</sub>	48.19 <sub>(0.67)</sub>	40.05 <sub>(0.68)</sub>	58.05 <sub>(0.70)</sub>
–w/o $\mathcal{L}_m^{co}$	23.45 <sub>(0.41)</sub>	84.74 <sub>(0.62)</sub>	41.83 <sub>(0.62)</sub>	71.16 <sub>(0.69)</sub>	69.77 <sub>(0.73)</sub>	47.72 <sub>(0.69)</sub>	43.30 <sub>(0.75)</sub>	56.07 <sub>(0.71)</sub>
–w/o $\mathcal{L}_m^N$	23.22 <sub>(0.33)</sub>	99.50 <sub>(0.12)</sub>	57.82 <sub>(0.65)</sub>	<b>96.86</b> <sub>(0.24)</sub>	<b>93.13</b> <sub>(0.33)</sub>	66.90 <sub>(0.65)</sub>	62.90 <sub>(0.69)</sub>	85.28 <sub>(0.51)</sub>
–w/o Adapt. Weight	23.62 <sub>(0.34)</sub>	98.52 <sub>(0.09)</sub>	57.45 <sub>(0.62)</sub>	<b>96.65</b> <sub>(0.27)</sub>	<b>93.06</b> <sub>(0.36)</sub>	62.23 <sub>(0.75)</sub>	61.54 <sub>(0.77)</sub>	84.61 <sub>(0.52)</sub>
–with $\mathcal{L}_m^S$	22.28 <sub>(0.26)</sub>	99.22 <sub>(0.12)</sub>	54.98 <sub>(0.63)</sub>	93.97 <sub>(0.36)</sub>	86.94 <sub>(0.49)</sub>	46.38 <sub>(0.64)</sub>	37.28 <sub>(0.70)</sub>	77.77 <sub>(0.57)</sub>

simultaneously instead of using the alternating update strategy. (3) “–w/o WMA”, which drops the WMA prediction generating strategy by setting  $\alpha_m = 1$ . (4) “–w/o  $\mathcal{L}_m^{co}$ ”, which drops the adaptive co-learning loss term ( $\mathcal{L}_m^{co}$ ) from the total loss function. (5) “–w/o  $\mathcal{L}_m^N$ ”, which drops the negative pseudo-labeling regularization loss term ( $\mathcal{L}_m^N$ ) from the total loss function. (6) “–w/o Adapt. Weight”, which drops the adaptive weighting component of AWCoL by setting all  $w_{[x]} = 1$ . (7) “–with  $\mathcal{L}_m^S$ ”, which adds a cross-entropy loss term on the labeled support instances to the total loss function for each model.

We compare AWCoL with all of its seven variants using the cross-domain 5-way 5-shot learning tasks on all the 8 datasets and report the results in Table 3. From the table, we can see that dropping any component from the proposed full model results in performance degradation in all the 48 cases except for only 4 cases: the “–w/o Alt. Update” variant on ChestX, the “–w/o Adapt. Weight” variant on Places, and the “–w/o  $\mathcal{L}_m^N$ ” variant on EuroSAT and Places. Even with the exception on ChestX, the “–w/o Alt. Update” variant demonstrates substantial performance drops from the full AWCoL method on all the other datasets, ranging from 14% (on ISIC) to 27% (on CUB). Such consistent performance drops across many datasets validate the essential contribution of each corresponding component of the proposed AWCoL. Meanwhile, with an additional loss term on the support instances, the “–with  $\mathcal{L}_m^S$ ” variant produces results that are much inferior to the proposed AWCoL across all the 8 datasets. The possible reason is that the labeled support instances have already been used to produce the class prototypes for the cross-entropy loss on the query instances in AWCoL, hence an additional cross-entropy loss  $\mathcal{L}_m^S$  may overfit the limited support instances.

## 5 Conclusion

In this paper, we proposed a weighted adaptive co-learning method to address the challenging cross-domain few-shot learning problem. The proposed method fine-tunes two prototypical classification models independently pre-trained in the source domain for the target FSL task. In each co-learning iteration, a weighted moving average strategy is deployed to generate probability predictions for the query instances using each model separately. The predictions of the two models are then combined to produce positive pseudo-labels, negative pseudo-labels, and adaptive weights for the query instances. We adopted an alternating fine-tuning mechanism to update each model separately by minimizing the weighted cross-entropy loss over the pseudo-labeled query instances while maximizing a similar cross-entropy loss with negative pseudo-labels to penalize false predictions. We conducted extensive experiments on eight CDFSL benchmark datasets. The results demonstrated the effectiveness of the proposed simple method compared to the state-of-the-art baselines in the literature.

## References

- [1] Thomas Adler, Johannes Brandstetter, Michael Widrich, Andreas Mayr, David Kreil, Michael Kopp, Günter Klambauer, and Sepp Hochreiter. Cross-domain few-shot learning by representation fusion. In *arXiv preprint arXiv:2010.06498*, 2020.
- [2] Debasmit Das, Sungrack Yun, and Fatih Porikli. ConfeSS: A framework for single source cross-domain few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2022.
- [3] Guneet Singh Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. In *International Conference on Learning Representations (ICLR)*, 2019.
- [4] Yingjun Du, Xiantong Zhen, Ling Shao, and Cees G. M. Snoek. Hierarchical variational memory for few-shot learning across domains. In *International Conference on Learning Representations (ICLR)*, 2022.
- [5] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, 2017.
- [6] Yuqian Fu, Yu Xie, Yanwei Fu, and Yu-Gang Jiang. Styleadv: Meta style adversarial training for cross-domain few-shot learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [7] Weifeng Ge and Yizhou Yu. Borrowing treasures from the wealthy: Deep transfer learning through selective joint fine-tuning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [8] Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogerio Feris. Spottune: transfer learning through adaptive fine-tuning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [9] Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris. A broader study of cross-domain few-shot learning. In *European Conference on Computer Vision (ECCV)*, 2020.
- [10] Bharath Hariharan and Ross Girshick. Low-shot visual recognition by shrinking and hallucinating features. In *International Conference on Computer Vision (ICCV)*, 2017.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [12] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. In *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019.

- [13] Ashraful Islam, Chun-Fu Richard Chen, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, and Richard J Radke. Dynamic distillation network for cross-domain few-shot recognition with unlabeled data. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [14] Taewon Jeong and Heeyoung Kim. Ood-maml: Meta-learning for few-shot out-of-distribution detection and classification. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [15] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *International Conference on Computer Vision workshops*, 2013.
- [16] Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex optimization. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [17] Pan Li, Shaogang Gong, Chengjie Wang, and Yanwei Fu. Ranking distance calibration for cross-domain few-shot learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [18] Wei-Hong Li, Xialei Liu, and Hakan Bilen. Cross-domain few-shot learning with task-specific adapters. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7161–7170, 2022.
- [19] Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim, and Sungwoong Kim. Fast auto-augment. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [20] Yanbin Liu, Juho Lee, Minseop Park, Saehoon Kim, Eunho Yang, Sungju Hwang, and Yi Yang. Learning to propagate labels: Transductive propagation network for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2019.
- [21] Yanbin Liu, Juho Lee, Linchao Zhu, Ling Chen, Humphrey Shi, and Yi Yang. A multi-mode modulator for multi-domain few-shot classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8453–8462, 2021.
- [22] Sharada P Mohanty, David P Hughes, and Marcel Salathé. Using deep learning for image-based plant disease detection. In *Frontiers in plant science*, 2016.
- [23] Cheng Perng Phoo and Bharath Hariharan. Self-training for few-shot transfer across extreme task differences. In *International Conference on Learning Representations (ICLR)*, 2020.
- [24] Scott Reed, Yutian Chen, Thomas Paine, Aäron van den Oord, SM Ali Eslami, Danilo Rezende, Oriol Vinyals, and Nando de Freitas. Few-shot autoregressive density estimation: Towards learning to learn distributions. In *International Conference on Learning Representations (ICLR)*, 2018.
- [25] Victor Garcia Satorras and Joan Bruna Estrach. Few-shot learning with graph neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.

- [26] Eli Schwartz, Leonid Karlinsky, Joseph Shtok, Sivan Harary, Mattias Marder, Abhishek Kumar, Rogerio Feris, Raja Giryes, and Alex Bronstein. Delta-encoder: an effective sample synthesis method for few-shot object recognition. *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [27] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- [28] 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 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [29] Eleni Triantafyllou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, et al. Meta-dataset: A dataset of datasets for learning to learn from few examples. *arXiv preprint arXiv:1903.03096*, 2019.
- [30] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions. In *Scientific data*, 2018.
- [31] Hung-Yu Tseng, Hsin-Ying Lee, Jia-Bin Huang, and Ming-Hsuan Yang. Cross-domain few-shot classification via learned feature-wise transformation. In *International Conference on Learning Representations (ICLR)*, 2020.
- [32] Ihsan Ullah, Dustin Carrión-Ojeda, Sergio Escalera, Isabelle Guyon, Mike Huisman, Felix Mohr, Jan N van Rijn, Haozhe Sun, Joaquin Vanschoren, and Phan Anh Vu. Meta-album: Multi-domain meta-dataset for few-shot image classification. *arXiv preprint arXiv:2302.08909*, 2023.
- [33] Grant Van Horn, Oisín 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 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [34] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2016.
- [35] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. California Institute of Technology, 2011.
- [36] Haoqing Wang and Zhi-Hong Deng. Cross-domain few-shot classification via adversarial task augmentation. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.
- [37] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

- [38] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems (NIPS)*, 2014.
- [39] Ruixiang Zhang, Tong Che, Zoubin Ghahramani, Yoshua Bengio, and Yangqiu Song. Metagan: An adversarial approach to few-shot learning. In *Advances in Neural information processing systems (NeurIPS)*, 2018.
- [40] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. In *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 2017.
- [41] Fei Zhou, Peng Wang, Lei Zhang, Wei Wei, and Yanning Zhang. Revisiting prototypical network for cross domain few-shot learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.