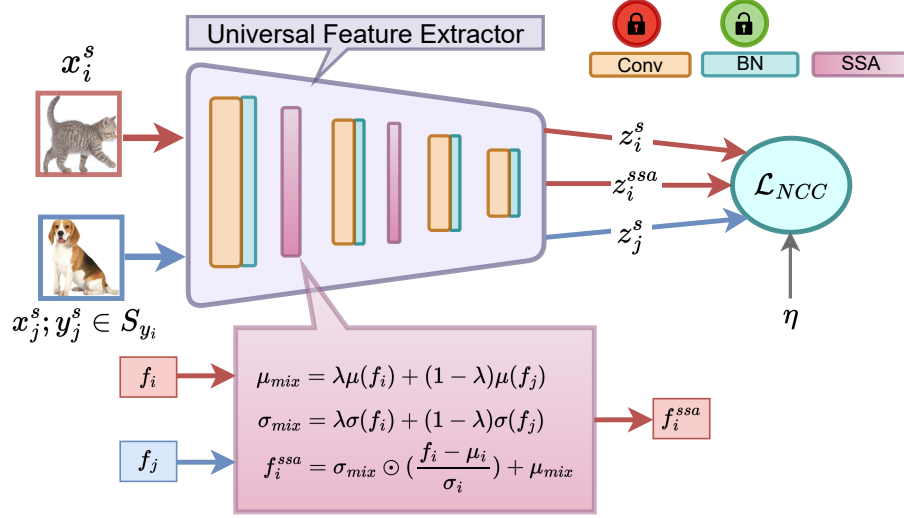


Graphical Abstract

Similar Class Style Augmentation for Efficient Cross-Domain Few-Shot Learning

Manogna Sreenivas, Soma Biswas



Cross-Domain Few-Shot Learning (CD-FSL) aims to recognize new classes from unseen domains, given limited training samples. Majority of the state-of-the-art approaches for this task introduce new task-specific additional parameters for adapting to the novel task, which involves changing the trained model architecture, in addition to increasing the number of model parameters. The first contribution of this work is to revisit the existing approaches like modifying the Batch Normalization affine parameters and the scale hyperparameter in cosine similarity based softmax loss for adapting the trained model to the new tasks, without changing the model architecture. Secondly, to aid the model learning with few examples per class, we propose to augment the data of each class with the styles of the semantically similar classes. Extensive evaluation on the challenging Meta-Dataset shows that this simple framework is very effective for the CD-FSL task. We also show that the Similar- class Style Augmentation module can be seamlessly integrated with existing approaches to further improve their performance, thus establishing the state-of-the-art in this challenging area.

Highlights

Similar Class Style Augmentation for Efficient Cross-Domain Few-Shot Learning

Manogna Sreenivas, Soma Biswas

- We address the important and challenging problem of Cross-Domain Few Shot Learning task.
- We show the effectiveness of adaptation of Batch Normalization layers for this task.
- We analyze the role of cosine similarity scaling factor which is scarcely studied in prior works.
- We propose a simple and effective similar class style-augmentation module to aid the learning in low-data regime.
- Extensive experiments on the large-scale MetaDataset show that the proposed framework in addition to performing well, is also efficient in terms of the parameters compared to the existing techniques.

Similar Class Style Augmentation for Efficient Cross-Domain Few-Shot Learning

Manogna Sreenivas^{a,**}, Soma Biswas^b

^aIndian Institute of Science, Bangalore, 560012, Karnataka, India

ARTICLE INFO

Keywords:
Few Shot Learning
Domain Shift
Data augmentation

ABSTRACT

Cross-Domain Few-Shot Learning (CD-FSL) aims to recognize new classes from unseen domains, given limited training samples. Majority of the state-of-the-art approaches for this task introduce new task-specific additional parameters for adapting to the novel task, which involves changing the trained model architecture, in addition to increasing the number of model parameters. The first contribution of this work is to revisit the existing approaches like modifying the Batch Normalization affine parameters and the scale hyperparameter in cosine similarity based softmax loss for adapting the trained model to the new tasks, without changing the model architecture. Secondly, to aid the model learning with few examples per class, we propose to augment the data of each class with the styles of the semantically similar classes. Extensive evaluation on the challenging Meta-Dataset shows that this simple framework is very effective for the CD-FSL task. We also show that the Similar-class Style Augmentation module can be seamlessly integrated with existing approaches to further improve their performance, thus establishing the state-of-the-art in this challenging area.

1. Introduction

Remarkable performance has been obtained in several computer vision tasks including image classification, object detection, image segmentation [6, 17, 21, 10, 8, 3] etc. in the past decade. However, though humans are capable of recognizing new objects by looking at just a handful of samples, the performance of the state-of-the-art deep neural networks usually degrades significantly in such low-data regime. In real world scenarios, collecting large-scale data can be infeasible and further annotating them would be very expensive, and thus it is necessary that machines can learn from limited amount of labelled data. In addition, in real-world, the test data can come from a distribution different from that of the training data. The objective of Cross-Domain Few-Shot Learning (CD-FSL) is to learn a feature extractor on large scale annotated data, and then adapt it to a new task with unseen classes from unseen domains, given a few labelled samples per class.


Due to the wide spread applicability of CD-FSL, several approaches have been recently proposed for addressing this task. Most of the recent state-of-the-art frameworks for CD-FSL address this task in a two-step process, namely (i) learning a universal feature extractor using training samples from multiple source datasets; and (ii) adapting the trained model to the novel tasks from the new domains using few labelled samples per class. After learning the universal feature extractor in the first step, the recent state-of-the-art approaches [24, 15, 16] usually incorporate task-specific learnable layers to adapt to the new task. Though they give impressive performance for the CD-FSL task, these frameworks result in changing the model architecture, also introducing additional model parameters. Since this may

not be desirable for many practical applications, here we propose a simple, yet effective framework for the CD-FSL task, without changing the model architecture or increasing the model parameters. In this work, we make the following contributions: First, starting with the universal feature extractor as in [15, 16], we revisit existing well known techniques for adapting the model to unseen domains by only modifying a few parameters of the trained model, as limited training data is available. Specifically, we propose to modify the Batch Normalization (BN) affine parameters to handle the unseen domain data. We also analyze the importance of the scale hyperparameter in the cosine similarity based softmax loss for the CD-FSL task. Next, since the novel tasks have less training data, as is usually the case in real-scenarios, we propose to augment the labelled data using styles from other semantically similar classes. For example, augmenting a *rose* class with the different variations present in a different flower class, say *sunflower*, is more intuitive as compared to using styles of a semantically dissimilar class, say *cat*. The proposed framework is termed **Similar-class Style Augmentation** with appropriate **Batch Norm** and **Scale** parameters (**SSA-BNS**). Extensive evaluation of this simple framework on the Meta-Dataset benchmark shows that it is capable of achieving close to state-of-the-art performance without altering the model architecture learnt in the first stage, thus requiring much lesser number of parameters when compared to the state-of-the-art approaches [15, 16]. We also show that the proposed Similar-class Style Augmentation (SSA) module can be integrated with existing CD-FSL frameworks to improve their performance.

2. Related work

Here, we discuss some of the related works in the areas of Few-Shot Learning (FSL), Cross-Domain Few-Shot learning (CD-FSL) and Data augmentation.

*Corresponding author

 manognas@iisc.ac.in (M. Sreenivas); somabiswas@iisc.ac.in (S.

Biswas)

ORCID(s):

Few-Shot Learning: FSL approaches [22, 26] often employ a Nearest Centroid Classifier (NCC), where the class centroids are obtained using the labelled support set samples. These centroids are then used to classify the samples from query set based on a distance metric like cosine similarity or Euclidean distance. Several prior works [22, 26, 9] formulate FSL as a learning to learn or meta-learning problem. In this perspective, the model is trained over a distribution of few-shot tasks sampled from the training data. Having the experience of performing well on the few-shot tasks, the model is now equipped to effectively learn from limited data of a novel test task. [22] uses Euclidean distance based NCC to meta-learn from training data. Matching Networks [26] classify query samples based on a distance-weighted linear combination of the support labels. In MAML [9], the authors propose to learn the model such that its weights act as a good initialization for future tasks, so that it can adapt using just a few gradient updates with only limited data. Such learning strategies, having proven effective for in-domain FSL are also extended to the CD-FSL setting. But the additional challenge here is to actually learn image representations that are effective across domains.

Cross-Domain Few-Shot Learning: We now discuss prior methods for CD-FSL task that involves learning multi-domain feature extractor(s). SUR [7] proposes to train multiple feature extractors, one for each domain. During test time, these feature extractors provide universal representations, over which a task-specific feature selection mechanism is employed along with NCC. URT [18] alternatively proposes a meta-learning based selection mechanism to adapt the universal representations. As these are computationally complex, requiring a test sample to forward pass through multiple feature extractors, later works [24, 15] propose effective ways to learn a single network using multiple domain data. In FLUTE [24], the authors propose to learn domain specific FiLM [20] layers while keeping the rest of the backbone shared across all training domains. In URL [15], given the features from the universal feature extractor, a linear transformation is used which is task specific and learnable, to obtain the adapted features. In TSA [16], task-specific adapters are included in the backbone, which are then fine-tuned using support set samples. These approaches are evaluated on the large-scale Meta-Dataset [25] benchmark, which establishes several baselines, extending the FSL methods for CD-FSL.

Data Augmentation: In general, augmenting the data improves the generalization abilities of a model, as it introduces more variations in the training data. Since data augmentation has been very successful in addressing several tasks including semi-supervised learning [23], unsupervised representation learning, contrastive learning [28, 2, 4, 1], we explore whether appropriate data augmentation can help to improve the adaptation performance for the CD-FSL task. Popular image based augmentations [5] include random crop, scale, flip, rotation, contrast enhancement, color distortions etc. Recently, several works [31, 29] propose to use Generative Adversarial Networks to augment their data. Style transfer

is another interesting data augmentation technique as it preserves the semantic content of the image. [12] propose to train a Style Transfer Network (STN) based on an Adaptive Instance Normalization layer. AdaIN layer transfers the feature statistics (channel-wise mean and variance) of one image to another and then train a decoder network to obtain the style transferred image. For CDFSL, since the model has access to few support set samples to update the base network, it is difficult to utilise some of these approaches like GAN or STN. In addition to these label preserving transformations, MixUp [11] is a very successful data augmentation approach where virtual samples and labels are created as convex combinations of two random training samples. This encourages linear behaviour for in-between samples and is observed to improve model robustness.

3. Problem definition and notations

First, we describe the CD-FSL problem and the notations used in this work. In general, a few shot learning task is described as N-way K-shot classification problem, where N represents the number of classes and K represents the number of examples per class available for training. We follow the Meta-Dataset [25] benchmark where the tasks are sampled such that the number of classes N and samples per class K are varied. Each task $\mathcal{T} = (\mathcal{S}, \mathcal{Q})$ comprises of a support set \mathcal{S} and a query set \mathcal{Q} . The model has to learn from a small labelled support set $\mathcal{S} = \{(x_i^s, y_i^s), i = 1 \dots n_S\}$ and is then evaluated on the query set $\mathcal{Q} = \{x_i^q, i = 1 \dots n_Q\}$ containing the same classes as in the corresponding support set. Here, x_i^s and y_i^s denotes the support sample and its corresponding label, while x_i^q refers to a query sample. n_S and n_Q denote the number of samples in the support and query set respectively. We denote the universal feature extractor as F and the feature representation of a sample x_i^s as $z_i^s = F(x_i^s)$.

In the cross-domain scenario, the training \mathcal{D}^{train} and test \mathcal{D}^{test} domains comprises of M_{train} and M_{test} datasets respectively. The objective is to use labelled data from \mathcal{D}^{train} to learn a feature extractor. This feature extractor is then adapted for novel few-shot tasks \mathcal{T} sampled from an unseen domain \mathcal{D}^{test} , comprising of previously unseen classes.

4. Proposed SSA-BNS Framework

With this background, we now describe the proposed framework (Fig. 1) for handling the CD-FSL task. First, we briefly describe the training of the feature extractor and NCC used in earlier works [15, 16] which we also use for fair comparison.

Universal Feature Extractor: We use the universal feature extractor proposed in URL [15], which is also used in TSA [16]. They learn the network in two stages: 1) Firstly, M_{train} single domain feature extractors are trained, one for each of the training domains, using cross entropy loss. 2) For efficient task adaptation in the future, knowledge distillation is performed from the single domain feature

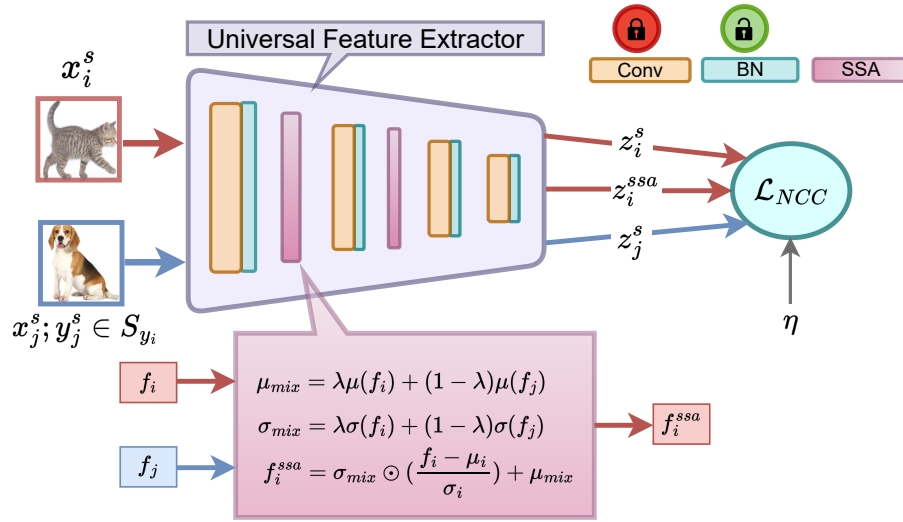


Figure 1: Proposed SSA-BNS framework: A support set instance x_i^s is augmented using a similar class sample $x_j^s, y_j^s \in S_{y_i}$ through SSA modules inserted in the universal feature extractor. Few shot task adaptation is done minimizing the \mathcal{L}_{NCC} loss over the actual features along with SSA augmented features.

extractors to a universal feature extractor, all of them sharing the same architecture. The domain specific features are mapped to a common space using domain specific adaptors, by minimizing Kullback-Liebler (KL) divergence between the model predictions of domain-specific features and the universal features. Alongside, the distance between these two sets of features are also minimized through Centered Kernel Alignment [13].

Nearest Centroid Classifier (NCC): Given a few-shot task $\mathcal{T} = (S, Q)$ sampled from \mathcal{D}^{test} , class centroids are obtained by averaging the feature representations $z_i^s = F(x_i^s)$ of the support set samples $x_i^s \in S$, for each class $k = 1 \dots C$ present in this task as:

$$\mathbf{c}_k = \frac{1}{|S_k|} \sum_{x_i^s \in S_k} z_i^s; \quad S_k = \{x_i^s | y_i^s = k\} \quad (1)$$

The NCC loss based on cosine similarity, with scaling factor η is obtained as

$$\mathcal{L}_{NCC}(z_i^s, y_i^s; \eta) = -\log \frac{e^{\eta \cos \theta_{i,y_i^s}}}{\sum_{j=1}^C e^{\eta \cos \theta_{i,j}}} \quad (2)$$

$$\text{where } \cos \theta_{i,y_i^s} = \frac{z_i^{sT} \mathbf{c}_{y_i^s}}{\|z_i^s\| \|\mathbf{c}_{y_i^s}\|}.$$

Now, we describe each of the components of the proposed framework, namely updating the BN affine and scale parameters and also the similar class style augmentation module.

4.1. Revisiting BN Adaptation for CD-FSL

In this work, we investigate how well we can adapt the trained model for a test task without modifying its architecture, or adding additional parameters. Specifically, we only adapt the BN affine parameters present in the network.

Batch Normalization: Let f^l denote the feature activations at layer l . The batch normalized feature activations f_{BN}^l are obtained as

$$f_{BN}^l = \gamma^l \hat{f}^l + \beta^l; \quad \text{where } \hat{f}^l = \frac{f^l - \mu^l}{\sqrt{(\sigma^l)^2 + \epsilon}} \quad (3)$$

Here, μ^l and σ^l are the running mean and standard deviation of the features estimated from the training data.

Recently, adaptation of BN parameters (γ^l, β^l) have been successfully used for many applications. In TENT [27], the BN scale and shift parameters are adapted to minimize the test entropy and has been very effective in mitigating performance degradation due to distribution shifts for Test-Time Adaptation (TTA) task. CD-FSL task is very different from TTA, as the classes in test tasks are unseen and there are only limited samples to learn from. In this work, we investigate the effectiveness of adapting BN parameters for CD-FSL task. The combined BN scale and shift parameters, which we denote as $\{\gamma, \beta\}$, are optimized to minimize the NCC loss over the support set as

$$\min_{\gamma, \beta} \frac{1}{n_S} \sum_{(x_i^s, y_i^s) \in S} \mathcal{L}_{NCC}(z_i^s, y_i^s; \eta) \quad (4)$$

where $z_i^s = F(x_i^s)$. We now study the impact of cosine similarity scale factor η on the CD-FSL performance.

4.2. Effect of Scale Factor on CD-FSL

Following several prior works in the few-shot learning regime [22, 7, 24], in this work, we use a Nearest Centroid Classifier based on cosine similarity metric. Using this classifier, the posterior probabilities over C classes are computed using softmax as

$$p(y = k | z_i^s; \eta) = \frac{e^{\eta \cos \theta_{i,k}}}{\sum_{j=1}^C e^{\eta \cos \theta_{i,j}}}; \quad k = \{1, \dots, C\} \quad (5)$$

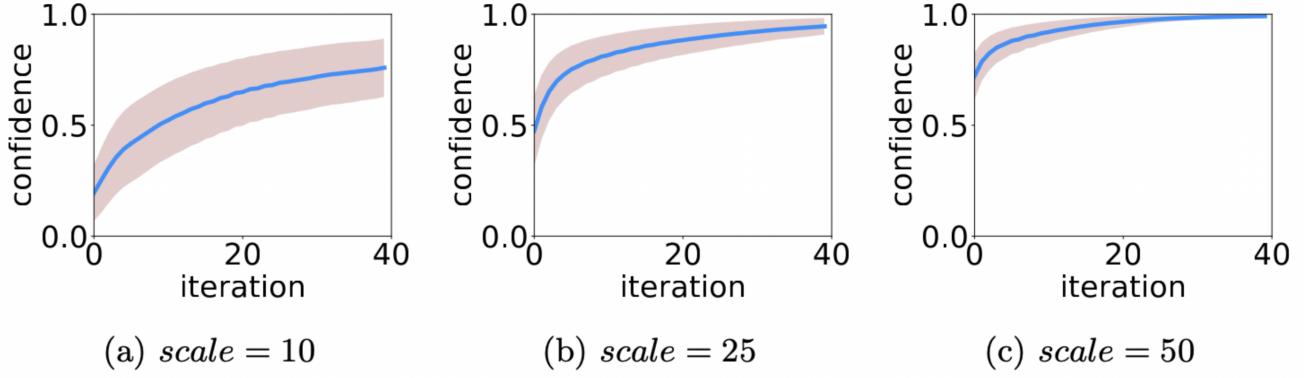


Figure 2: The mean and standard deviation of confidence scores of the predicted class over 600 tasks from Aircraft dataset.

where η is the scale hyper-parameter. The softmax scores being a function of the logits, is dependent on the range of values taken by the logits. The cosine similarity values, being the logits here, are in the range $[-1,1]$. Thus, it is a standard practice to scale these logits with a factor η in order to expand the input range for the softmax function, thereby assigning reasonable probability scores. This scale factor although plays a significant role in the training process of cosine similarity based metric learning, has only been scarcely studied, especially in the few shot regime. AdaCos [30] proposed an approach of identifying an optimal scale parameter based on the number of classes, for a large scale face recognition task. However, the CD-FSL task significantly differs from that addressed in [30], as the number of samples available per class is very low here. In addition, each test task can have varying number of classes and can come from an unseen domain, which can be very different from the training domains. Thus, the analysis in [30], although insightful, may not be directly applicable to the CD-FSL task. In the previous CD-FSL works URL [15] and TSA [16], the scale parameter η was fixed to 10. In this work, we investigate the impact of this scale factor on the CD-FSL task, where we only update the BN parameters using the labelled support set. Here, since the target domain can be very different from the training domains with few labeled examples per class, we expect the model to be less confident in classifying the training data, as compared to other applications where the test domain is similar to the training domains or there are more training examples available. This implies that for majority of the tasks and training examples, the cosine similarities will be lower, resulting in low probability values. Since the scale factor controls the range of probability values that the cosine similarity gets mapped to, we expect that a higher scale factor will be beneficial for the CD-FSL task, so that the model does not incur any loss even when a training sample is correctly classified.

For the CD-FSL task, using the universal feature extractor, we perform experiments by varying scale parameter η as 10, 25 and 50. We use one of the training domains *Aircraft*, to analyse the impact of scale factor on the training process.

To understand the confidence scores of the predicted class of support set samples, we plot the confidence scores averaged over all the support samples and the tasks, as training progresses in Fig. 2. We observe that for a scale of 10, the model is only able to assign confidence scores in the range of 0.6-0.8 even towards the end of training, although the support set samples get correctly classified. Using a scale of 25, the model learns to correctly classify samples while also gradually assigning high confidence. On the other hand, using a scale of 50, the model predictions attain confidence greater than 0.8 for most of the samples in only about 10 iterations. Correspondingly, the loss converges in less than 25 iterations resulting in overfitting on the support set as shown in Fig. 3, resulting in degraded performance when compared to using a scale of 25.

As we later show in Section 2, empirically also, we observe that using a higher scale ($\eta = 25$) improves the CD-FSL performance significantly as compared to setting it to 10, which was the default in URL and TSA. This verifies our hypothesis, that because of the challenging nature of the CD-FSL task, a higher scale factor is desirable. Needless to mention, that if the scale factor is increased to a very high value, the model assigns a high probability score irrespective of the cosine similarity value, which in turn reduces the discriminability between samples and can hurt the training process.

4.3. Similar Class Style Augmentations

Our objective is to learn to recognize new classes in unseen domains using limited training data. Here, we propose to augment the support set in the feature space using the styles of similar class samples. For example, suppose we have a *cat* training class which has examples of *sitting cat*, *sleeping cat*, we can augment it with the styles of a neighbouring class, say *dog*, which might have examples of *standing dog*, *running dog*. It may not be meaningful to augment it with the styles of semantically dissimilar classes, say *bird*, which may have examples of *flying birds*. The early layer feature statistics are representative of the domains/styles of the input sample [12, 32].

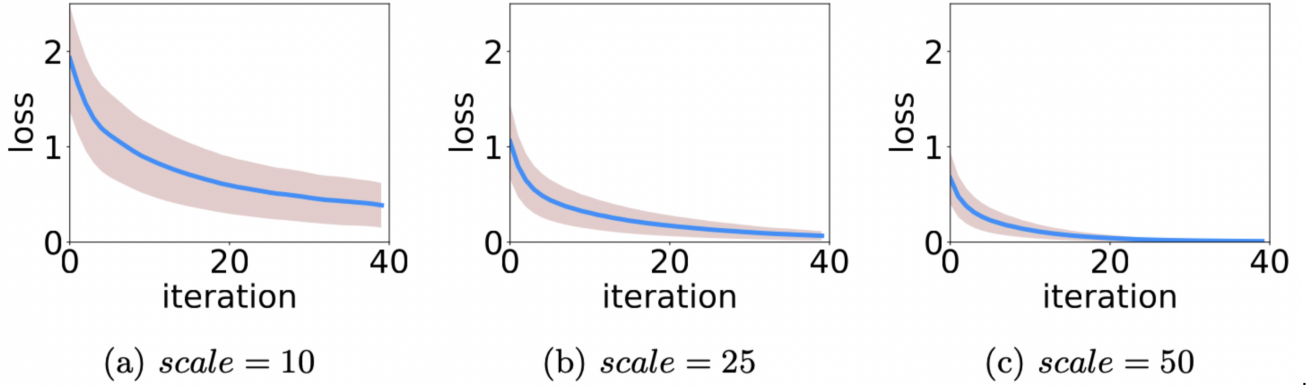


Figure 3: The mean and standard deviation of loss curves over 600 tasks for 40 training iterations from Aircraft dataset.

In this section, we simplify the notation, by dropping the superscript s referring to support set samples for clarity. Given a sample $(x_i, y_i) \in \mathcal{S}$, we propose to mix its feature statistics (representative of style) with that of another sample x_j to synthesize a style augmented feature of the same class y_i . Although, this is a label preserving transform, as explained earlier, arbitrarily picking samples could adversely affect the training process. Given the small support set, in order to effectively style augment while also avoiding irrelevant feature perturbations, we propose to condition the mixing of styles to only relevant classes based on the similarity of the class representations. We now formulate the SSA (Similar-class Style Augmentation) module.

Firstly, we obtain the pairwise class similarities based on their centroids. The cosine similarity of class i with that of class j can be obtained from their centroids $\mathbf{c}_i, \mathbf{c}_j$ as

$$\text{sim}(\mathbf{c}_i, \mathbf{c}_j) = \frac{\mathbf{c}_i^T \mathbf{c}_j}{\|\mathbf{c}_i\| \|\mathbf{c}_j\|} \quad (6)$$

For a class k , the set of similar classes is determined as:

$$\mathcal{S}_k = \{t | \text{sim}(\mathbf{c}_t, \mathbf{c}_k) > \tau; t = 1, \dots, C\} \quad (7)$$

where τ is a pre-set threshold.

Denoting f_i as the intermediate feature at a layer l , of support sample x_i belonging to class y_i , we randomly select another sample x_j belonging to one of its similar classes \mathcal{S}_{y_i} , i.e., $y_j \in \mathcal{S}_{y_i}$. The style of x_j is used to perturb the style of x_i to get the style augmented sample as follows:

$$\begin{aligned} \mu_{\text{ssa}}(f_i; f_j, y_j \in \mathcal{S}_{y_i}) &= \lambda \mu(f_i) + (1 - \lambda) \mu(f_j) \\ \sigma_{\text{ssa}}(f_i; f_j, y_j \in \mathcal{S}_{y_i}) &= \lambda \sigma(f_i) + (1 - \lambda) \sigma(f_j) \\ f_i^{\text{ssa}} &= \sigma_{\text{ssa}}(f_i; f_j, y_j \in \mathcal{S}_{y_i}) \odot \frac{f_i - \mu(f_i)}{\sigma(f_i)} + \mu_{\text{ssa}}(f_i; f_j) \end{aligned} \quad (8)$$

Here, $\mu(f_i), \mu(f_j)$ and $\sigma(f_i), \sigma(f_j)$ are the channel-wise mean and standard deviation of the features f_i and f_j respectively. Let z_i, z_j and z_i^{ssa} denote the output feature vectors (from last layer) obtained on passing f_i, f_j and f_i^{ssa} respectively through the rest of the network. The SSA augmented

feature z_i^{ssa} is obtained by mixing the styles of x_i and x_j . However the content of x_i is preserved in the feature z_i^{ssa} and hence belongs to the class y_i .

SSA-BNS: To summarize, we minimize the NCC loss over the support set features and their augmentations obtained through SSA as

$$\min_{\gamma, \beta} \frac{1}{2n_S} \sum_{(x_i^s, y_i^s)} L_{\text{NCC}}(z_i^s, y_i^s; \eta) + L_{\text{NCC}}(z_i^{\text{ssa}}, y_i^s; \eta) \quad (9)$$

For evaluation on query set \mathcal{Q} , cosine similarity metric is used to assign the query sample x_i^q to the nearest centroid as:

$$\hat{y}_i^q = \arg \max_k \cos \theta_{i,k}; \quad \cos \theta_{i,k} = \frac{\mathbf{c}_k^T z_i^q}{\|\mathbf{c}_k\| \|z_i^q\|} \quad (10)$$

where z_i^q is the feature of the query sample x_i^q .

Though this work is inspired from MixStyle [32], there are significant differences between the two approaches as discussed below:

1) Mixstyle mixes styles of samples from different source domains during training, primarily to interpolate the domain information. In our work, during training for the new tasks, the source domain data is not available. All the examples that are used for augmentation belong to the same unseen domain.

(2) In MixStyle, the two samples can come from any class. We observe that class-agnostic mixing is not meaningful for this application as we are primarily mixing the styles, and it also hurts the performance as seen empirically. Thus, the style mixing takes place only between semantically similar classes to generate realistic augmentations, which aid the training process.

5. Experimental Evaluation

Here, we describe the datasets used, implementation details and the experimental results.

Table 1

Average accuracy over 600 tasks are reported for all the compared approaches for CD-FSL task. URL, TSA, SSA+BNS, TSA+BNS uses the same universal feature extractor. * indicates that these approaches uses additional parameters as compared to the feature extractor (Table 3).

Dataset	FLUTE	tri-M	URL*	TSA*	SSA-BNS	TSA*+SSA
Imagenet	58.6 \pm 1.0	51.8 \pm 1.1	58.8 \pm 1.1	59.5 \pm 1.0	56.6 \pm 1.0	58.9 \pm 1.1
Omniplot	92.0 \pm 0.6	93.2 \pm 0.5	94.5 \pm 0.4	94.9 \pm 0.4	95.2 \pm 0.5	95.6 \pm 0.4
Aircraft	82.8 \pm 0.7	87.2 \pm 0.5	89.4 \pm 0.4	89.9 \pm 0.4	89.6 \pm 0.4	90.0 \pm 0.5
Birds	75.3 \pm 0.8	79.2 \pm 0.8	80.7 \pm 0.8	81.1 \pm 0.8	81.8 \pm 0.8	82.2 \pm 0.7
Textures	71.2 \pm 0.8	68.8 \pm 0.8	77.2 \pm 0.7	77.5 \pm 0.7	76.4 \pm 0.7	77.6 \pm 0.7
Quick draw	77.3 \pm 0.7	79.5 \pm 0.7	82.5 \pm 0.6	81.7 \pm 0.6	82.8 \pm 0.6	82.7 \pm 0.7
Fungi	48.5 \pm 1.0	58.1 \pm 1.1	68.1 \pm 0.9	66.3 \pm 0.8	66.7 \pm 0.8	66.6 \pm 0.8
VGG Flower	90.5 \pm 0.5	91.6 \pm 0.6	92.0 \pm 0.5	92.2 \pm 0.5	92.8 \pm 0.6	93.0 \pm 0.5
Traffic Sign	63.0 \pm 1.0	58.4 \pm 1.1	63.3 \pm 1.1	82.8 \pm 1.0	77.9 \pm 1.1	84.9 \pm 1.1
MSCOCO	52.8 \pm 1.1	50.0 \pm 1.0	57.3 \pm 1.0	57.6 \pm 1.0	56.1 \pm 0.9	58.1 \pm 1.0
MNIST	96.2 \pm 0.3	95.6 \pm 0.5	94.7 \pm 0.4	96.7 \pm 0.4	98.3 \pm 0.5	98.5 \pm 0.4
CIFAR-10	75.4 \pm 0.8	78.6 \pm 0.7	74.2 \pm 0.8	82.9 \pm 0.7	79.4 \pm 0.7	82.9 \pm 0.7
CIFAR-100	62.0 \pm 1.0	67.1 \pm 1.0	63.5 \pm 1.0	70.4 \pm 0.9	69.0 \pm 0.9	70.8 \pm 0.9
Average seen	74.5	76.2	80.4	80.4	80.2	80.8
Average unseen	69.9	69.9	70.6	78.1	76.1	79.0
Average all	72.7	73.8	76.6	79.5	78.7	80.1

Meta-Dataset Benchmark: Meta-Dataset [25] is a recent benchmark proposed for the CD-FSL task. This benchmark comprises of 10 datasets of which eight act as training domains and the other two as test domains. The training domains comprises of ImageNet, Omniglot, Aircraft, Birds, Textures, Quick draw, Fungi and VGG Flower datasets. Traffic Signs and MSCOCO act as test domains. Additionally, following prior works [24, 15, 16], we also include MNIST, CIFAR-10 and CIFAR-100 as test domains. The few shot tasks are sampled with varying number of classes N , with N varying from 5 to the maximum number of classes available in the dataset. The number of samples per class K , although varying, is capped at 100 samples and also the total support set size is limited to 500 samples in all.

Implementation details: We use a ResNet-18 backbone trained on the eight training domains following the Meta-Dataset protocol [25]. This universal feature extractor is the same as that used in [15, 16], which is trained by distilling knowledge from eight individually trained feature extractors.

For few shot task adaptation, we only finetune the BN affine parameters $\{\gamma, \beta\}$ on the support set of a given task using the NCC. As in TSA [16], we also use Adadelta optimizer with learning rate of 0.001 and train for 40 iterations. We set λ to 0.5 in eqn.(8) to perform SSA, scale parameter η to 25 and the threshold τ to 0.7 in all the experiments. MixStyle layers are inserted after the first and second ResNet blocks.

Experimental Results: Table 1 reports the average accuracy and 95% confidence interval over 600 tasks obtained using the proposed framework, alongside comparisons with the state-of-the-art approaches. The first group of results correspond to the seen domains and the second group to

unseen domains. The results of the other works are directly taken from [16]. For the proposed work, we reports two sets of results. First, we report the results of our SSA-BNS framework, which does not include any change in architecture or increase in number of parameters compared to the originally trained model using the source datasets. We observe that for the unseen domains, the proposed SSA-BNS significantly outperforms FLUTE [24], tri-M [19] and URL [15], and is second only to TSA [16]. For the seen domains, it performs comparably to all the other approaches. On an average, SSA-BNS achieves an accuracy of 78.7% as compared to 76.6% obtained by URL. It is only second to [16], which obtains an average accuracy of 79.5%. Note that all the other approaches uses significantly more parameters and also involves change in the trained model architecture as shown in Table 3. In addition, we also report the results of integrating the proposed SSA module into the state-of-the-art [16] approach. We observe that this simple, yet effective augmentation (with no additional parameters) can be used to further improve the performance of the state-of-the-art approaches for both seen and unseen domains.

6. Additional Analysis

Here, we perform additional analysis to better understand the proposed framework and the usefulness of each of its components.

Ablation Study: We study the importance of each component for two seen and unseen domains and report the results in Table 2. Firstly, we do not adapt the network and simply use a NCC to classify the query samples, which corresponds to the first row in Table 2. Then we study the role of BN adaptation for the CD-FSL task. While earlier works like

Table 2

CD-FSL performance on two seen and unseen domains (averaged over 600 tasks). Effect of BN adaptation, scale factor η and the SSA module.

SSA	BNS	Seen domains		Unseen domains	
		Aircraft	Fungi	CIFAR-100	MSCOCO
No	No	87.0	65.6	59.9	53.1
No	Yes($\eta = 10$)	89.1	66.0	66.9	54.5
No	Yes($\eta = 25$)	89.5	66.4	68.4	55.7
No	Yes($\eta = 50$)	89.5	66.2	67.7	55.4
Yes	Yes($\eta = 25$)	89.6	66.7	69.0	56.1

Table 3

Comparison of computational complexity with SOTA methods

Method	#Additional parameters	# Trainable parameters
URL	262144	262144
TSA	1482752	1482752
TSA+SSA	1482752	1482752
SSA-BNS	None	9600

URL[15], TSA[16] use a default scale factor of 10, we experiment with other values (25 and 50) to study its impact. We observe that using a scale of 25 consistently performs better as compared to using a smaller scale factor like 10. Although using scale of 50 is preferable over 10, the results degrade when compared to using a scale of 25. This happens due to the assignment of high confidence for correctly classified support samples with low cosine similarity values as discussed in Section 4.2. We also observe that the SSA module consistently improves the performance for both the seen and unseen domains. The effect of SSA module can also be seen by comparing the accuracies obtained by TSA and TSA+SSA in Table 1. This analysis shows that each of the proposed components contribute significantly towards improving the performance.

Trainable Parameters: We compare the complexity of the proposed SSA-BNS framework with that of the previous state-of-the-art methods, URL [15] and TSA [16]. Specifically, we report (i) the number of additional parameters incorporated and (ii) the number of trainable parameters in Table 3 in the universal feature extractor for few shot task adaptation. The proposed SSA-BNS framework does not require any change in the model architecture, as it only leverages the existing BN layers in the network for adaptation. Here, only the BN affine parameters (9600 in the ResNet-18 architecture) are trainable. We observe that SSA-BNS framework is very efficient, and significantly outperforms URL having 262,144 additional parameters. It is only second to TSA which uses 154 times more parameters. TSA+SSA, which outperforms all the existing approaches has the same number of parameters as TSA, since the proposed SSA module does not introduce any additional parameters.

Performance with Different Augmentations: Although it is common to use data augmentations for training deep models, it is not always clear which type of augmentation will

Table 4

Comparison of the proposed BNS module with different augmentation techniques.

Type	Seen domains		Unseen domains	
	Aircraft	Fungi	CIFAR-100	MSCOCO
RandAugment [5]	88.8	65.2	66.9	55.2
MixUp [11]	88.4	66.3	67.9	54.6
Feature MixUp [14]	88.9	66.3	68.3	55.3
Random MixStyle [32]	89.6	66.2	68.2	55.3
SSA (Proposed)	89.6	66.7	69.0	56.1

help in which application, especially in the low-data regime. Here, we perform extensive experiments using a variety of image and feature space augmentations for the CD-FSL task. Specifically, we investigate the following successful and popularly used data augmentation techniques for our task: 1) RandAugment [5] sequentially applies image transforms like autocontrast, polarize, shear, posterize, equalize etc. to the input image. 2) In MixUp [11], image augmentations and their labels are obtained as convex combinations of image samples and their labels. For e.g., given two samples x_i and x_j , mixup samples are created as $x_{aug} = \lambda x_i + (1 - \lambda)x_j$ with label $y_{aug} = \lambda y_i + (1 - \lambda)y_j$. 3) Inspired from i-Mix [14], we perform Feature MixUp where a feature is perturbed by slightly shifting it towards another sample. The augmented feature of a sample x_i using x_j is obtained as $f_i^{aug} = \lambda f_i + (1 - \lambda)f_j$ with λ set to 0.9. Such feature mixing modules are inserted after the first two ResNet blocks. 4) Random MixStyle [32]: Here, the styles are augmented as described in Section 4.3, but the samples whose styles are mixed can belong to any class, not necessarily similar classes. 5) SSA (Proposed): In this work, we use similar class style augmentations as described in Section 4.3. We observe that the proposed SSA module outperforms all the commonly used augmentation techniques.

7. Conclusion

In this work, we address the challenging task of Cross-Domain Few-Shot Learning (CD-FSL), where a model learnt using diverse source domains has to be adapted to new tasks (with different classes from different domains) with only a few training examples per class. Since changing the trained model architecture and increasing the number of trainable parameters is often infeasible/not desirable, we propose a novel framework, termed SSA-BNS, which modifies the BN affine parameters and the scale parameter of the cosine similarity based softmax loss, without modifying the feature extractor or adding any parameters to the source trained model. Alongside, we also propose a Similar-class Style augmentation module to aid the training in this low-data regime. Extensive experiments on several datasets show that the proposed framework outperforms several recent works with much lesser parameters and modifications. We also show that the proposed SSA module can be integrated with the current state-of-the-art approach [16] to further improve its performance, giving the state-of-the-art performance for the challenging CD-FSL task.

References

- [1] Berthelot, D., Carlini, N., Goodfellow, I., Papernot, N., Oliver, A., Raffel, C.A., 2019. Mixmatch: A holistic approach to semi-supervised learning. *Advances in Neural Information Processing Systems* 32.
- [2] Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., Joulin, A., 2020. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems* 33, 9912–9924.
- [3] Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L., 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *TPAMI* 40, 834–848.
- [4] Chen, T., Kornblith, S., Norouzi, M., Hinton, G., 2020. A simple framework for contrastive learning of visual representations, in: *International conference on machine learning*, PMLR. pp. 1597–1607.
- [5] Cubuk, E.D., Zoph, B., Shlens, J., Le, Q.V., 2020. Randaugment: Practical automated data augmentation with a reduced search space, in: *CVPRW*, pp. 3008–3017. doi:10.1109/CVPRW50498.2020.00359.
- [6] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L., 2009. Imagenet: A large-scale hierarchical image database, in: *CVPR*, pp. 248–255. doi:10.1109/CVPR.2009.5206848.
- [7] Dvornik, N., Schmid, C., Mairal, J., 2020. Selecting relevant features from a multi-domain representation for few-shot classification, in: *ECCV*, Springer. pp. 769–786.
- [8] Everingham, M., Gool, L., Williams, C.K., Winn, J., Zisserman, A., 2010. The pascal visual object classes (voc) challenge. *IJCV* 88, 303–338. URL: <https://doi.org/10.1007/s11263-009-0275-4>, doi:10.1007/s11263-009-0275-4.
- [9] Finn, C., Abbeel, P., Levine, S., 2017. Model-agnostic meta-learning for fast adaptation of deep networks, in: *ICML*, PMLR. pp. 1126–1135.
- [10] He, K., Gkioxari, G., Dollár, P., Girshick, R., 2017. Mask r-cnn, in: *ICCV*, pp. 2961–2969.
- [11] Hongyi Zhang, Moustapha Cisse, Y.N.D.D.L.P., 2018. mixup: Beyond empirical risk minimization. *ICLR* URL: <https://openreview.net/forum?id=r1Ddp1-Rb>.
- [12] Huang, X., Belongie, S., 2017. Arbitrary style transfer in real-time with adaptive instance normalization, in: *ICCV*.
- [13] Kornblith, S., Norouzi, M., Lee, H., Hinton, G., 2019. Similarity of neural network representations revisited, in: Chaudhuri, K., Salakhutdinov, R. (Eds.), *Proceedings of the 36th International Conference on Machine Learning*, PMLR. pp. 3519–3529. URL: <https://proceedings.mlr.press/v97/kornblith19a.html>.
- [14] Lee, K., Zhu, Y., Sohn, K., Li, C.L., Shin, J., Lee, H., 2021. i-mix: A domain-agnostic strategy for contrastive representation learning, in: *ICLR 2021*.
- [15] Li, W.H., Liu, X., Bilén, H., 2021. Universal representation learning from multiple domains for few-shot classification, in: *ICCV*, pp. 9526–9535.
- [16] Li, W.H., Liu, X., Bilén, H., 2022. Cross-domain few-shot learning with task-specific adapters, in: *CVPR*.
- [17] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L., 2014. Microsoft coco: Common objects in context, in: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (Eds.), *ECCV*, pp. 740–755.
- [18] Liu, L., Hamilton, W., Long, G., Jiang, J., Larochelle, H. (Eds.), 2021a. A Universal Representation Transformer Layer for Few-Shot Image Classification.
- [19] Liu, Y., Lee, J., Zhu, L., Chen, L., Shi, H., Yang, Y., 2021b. A multi-mode modulator for multi-domain few-shot classification, in: *ICCV*.
- [20] Perez, E., Strub, F., de Vries, H., Dumoulin, V., Courville, A.C., 2018. Film: Visual reasoning with a general conditioning layer, in: *AAAI*.
- [21] Ren, S., He, K., Girshick, R., Sun, J., 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *NeurIPS* 28.
- [22] Snell, J., Swersky, K., Zemel, R., 2017. Prototypical networks for few-shot learning. *NeurIPS* 30.
- [23] Sohn, K., Berthelot, D., Carlini, N., Zhang, Z., Zhang, H., Raffel, C.A., Cubuk, E.D., Kurakin, A., Li, C.L., 2020. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in Neural Information Processing Systems* 33, 596–608.
- [24] Triantafillou, E., Larochelle, H., Zemel, R., Dumoulin, V., 2021. Learning a universal template for few-shot dataset generalization, in: *ICML*, PMLR. pp. 10424–10433.
- [25] Triantafillou, E., Zhu, T., Dumoulin, V., Lamblin, P., Evci, U., Xu, K., Goroshin, R., Gelada, C., Swersky, K.J., Manzagol, P.A., Larochelle, H., 2020. Meta-dataset: A dataset of datasets for learning to learn from few examples, in: *ICLR*.
- [26] Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al., 2016. Matching networks for one shot learning. *NeurIPS* 29.
- [27] Wang, D., Shelhamer, E., Liu, S., Olshausen, B., Darrell, T., 2021. Tent: Fully test-time adaptation by entropy minimization, in: *ICLR*. URL: <https://openreview.net/forum?id=uX13bZLkr3c>.
- [28] Xie, Q., Dai, Z., Hovy, E., Luong, T., Le, Q., 2020. Unsupervised data augmentation for consistency training. *Advances in Neural Information Processing Systems* 33, 6256–6268.
- [29] Yamaguchi, S., Kanai, S., Eda, T., 2020. Effective data augmentation with multi-domain learning gans, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 6566–6574.
- [30] Zhang, X., Zhao, R., Qiao, Y., Wang, X., Li, H., 2019. Adacos: Adaptively scaling cosine logits for effectively learning deep face representations, in: *CVPR*.
- [31] Zhao, M., Cong, Y., Carin, L., 2020. On leveraging pretrained gans for generation with limited data, in: *International Conference on Machine Learning*, PMLR. pp. 11340–11351.
- [32] Zhou, K., Yang, Y., Qiao, Y., Xiang, T., 2021. Domain generalization with mixstyle, in: *ICLR*.