

Augmented Bi-path Network for Few-shot Learning

Baoming Yan^{*†‡}, Chen Zhou^{*†}, Bo Zhao[§], Kan Guo[‡], Jiang Yang[‡], Xiaobo Li[‡], Ming Zhang[†] and Yizhou Wang[†]
[†] Peking University, [‡] Alibaba Group, [§] The University of Edinburgh

{bmyan, zhouch18, yizhou.wang}@pku.edu.cn; mzhang@net.pku.edu.cn; {guokan.gk, yangjiang.yj, xiaobo.libx}@alibaba-inc.com; bo.zhao@ed.ac.uk

Abstract—Few-shot Learning (FSL) which aims to learn from few labeled training data is becoming a popular research topic, due to the expensive labeling cost in many real-world applications. One kind of successful FSL method learns to compare the testing (query) image and training (support) image by simply concatenating the features of two images and feeding it into the neural network. However, with few labeled data in each class, the neural network has difficulty in learning or comparing the local features of two images. Such simple image-level comparison may cause serious mis-classification. To solve this problem, we propose Augmented Bi-path Network (ABNet) for learning to compare both global and local features on multi-scales. Specifically, the salient patches are extracted and embedded as the local features for every image. Then, the model learns to augment the features for better robustness. Finally, the model learns to compare global and local features separately, *i.e.*, in two paths, before merging the similarities. Extensive experiments show that the proposed ABNet outperforms the state-of-the-art methods. Both quantitative and visual ablation studies are provided to verify that the proposed modules lead to more precise comparison results.

I. INTRODUCTION

In recent years, Deep Learning methods have achieved significant progress in computer vision by applying deeper architectures [1]–[4] to bigger datasets [5]–[7]. When training a deep neural network, the performance heavily depends on the amount of labeled training data. However, in many real-world tasks, it is time-consuming even prohibitive to collect and annotate enough data for training the popular deep networks. For example, annotating some fine-grained categories [8] or medical data [9] are restricted by not only the few available samples but also the few domain specialists. Therefore, how to get rid of cumbersome labelling and train a good classification model with few labeled data, *i.e.*, Few-shot Learning (FSL) [10]–[14], is a valuable research problem.

Many methods have been proposed to deal with the Few-shot Learning problem in past decades. Early studies [10], [15] use a small number of samples to directly construct a model for classifying new samples. However, the model is not able to learn the real data distribution or generalize to testing data. Recently, Meta-learning based methods [12], [13], [16] are proposed, in which many episodes (basic tasks) are sampled from the training data. The model is trained on many sampled tasks for learning meta-knowledge that is generalized to a distribution of tasks. Methods in this framework differ in the design of the classifier for basic tasks.

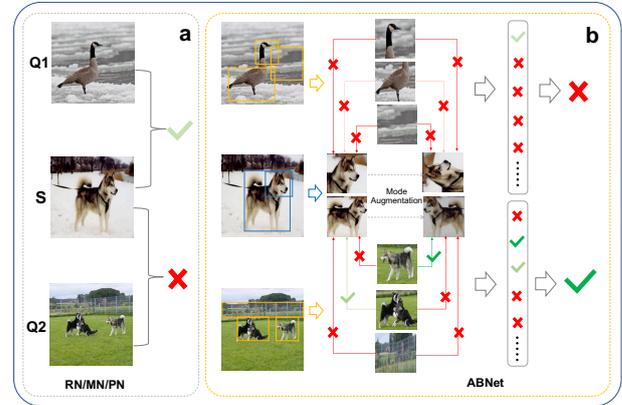


Fig. 1: The effectiveness of salient patches. S stands for the training (support) image. Q1 and Q2 are testing (query) images. In previous methods, only the global features (of the whole images) are compared, as shown in (a). However, our method compares both global and local features by extracting the salient image patches. Hence, the comparison result is more precise.

[12] used a prototype (class center) based nearest neighbor classifier to classify testing (query) images based on those training (support) images. [13] proposed to learn to compare the similarity between the query and support images. With the shared feature embedding network, the features of two images are concatenated and fed it into the comparison neural network which outputs the similarity between two images.

However, with few labeled data in every class, the neural network has difficulty in learning or comparing the local features of two images. Such simple image-level comparison may cause serious mis-classification. As illustrated in Fig. 1, given one support image (S) and two query images (Q1 and Q2), the global features of S and Q1 are more similar than those of S and Q2. However, S and Q2 belongs to the class “dog”, while Q1 is an image of a “bird”. If we extract and compare some salient image patches, *e.g.*, head and legs, we can easily find that S and Q2 have more similar patches, and they should be classified to the same class.

For more precise comparison, we propose Augmented Bi-path Network (ABNet) to learn to compare both global and local features on multi-scales. Our method includes four main modules. First, the salient patches that contains informative parts are extracted for every image. Second, the original image and its salient patches are embedded by the shared feature embedding module. Third, we learn to augment both global and local (salient patch) features of support images for better robustness. Fourth, we learn to compare the features of support and query images and output the similarity between

* Equal contribution.

the two images. Instead of calculating the similarity based on concatenated features directly [13], we generate the similarity maps based on concatenated features and learn to re-weight them. Then the global and local similarity maps are merged to produce the overall similarity. Extensive experiments on three challenging benchmarks show that our method outperforms the state-of-the-art methods with a large margin. Ablation studies are provided for verifying that the proposed modules are important for achieving better performance. Visual analysis is given for better illustration of how feature augmentation and local features lead to correct classification.

The main contributions of this paper includes two folds:

- We propose Augmented Bi-path Network (ABNet) to learn to compare both global and local features of support and query images, which includes two novel modules, *namely*, “Learning to Augment” and “Learning to Compare”.
- We evaluate our approach on three challenging Few-shot Learning benchmarks, miniImageNet, Caltech-256 and tieredImageNet. Our ABNet outperforms the state-of-the-art by a large margin.

II. RELATED WORK

A. Meta-learning Based Methods

Meta-learning based methods [11], [12], [16]–[19] aim to learn a more generalized model by meta-training on many sampled FSL tasks. Typically, Finn et al. [20] propose an model-agnostic algorithm for meta-learning that trains a models parameters such that a small number of gradient updates will lead to fast learning on a new task In addition, Ravi et al. [21] describe an LSTM-based model for meta-learning. The model is trained to discover a good initialization of the learners parameters, as well as a successful mechanism for updating the parameters to new task.

The idea of “learn to compare” is also widely used. Gregory Koch et al. [11] first propose siamese neural networks for One-shot image recognition. Through a shared network structure, deep features are learnt and compared to decide whether the two inputs belong to the same class. Oriol Vinyals et al. [16] build a matching network, which learns a LSTM encoder to embedding the deep features conditioned on the specific support set and query set. Inspired by the evaluation process of Few-shot Learning, they construct similar few-shot classification task in the training data. Following the same training strategy, ProtoNet [12], RelationNet [13] and many other superior networks [21]–[23] are proposed.

Our ABNet also belongs to meta-learning based methods. The main difference between aforementioned methods and ours is that we learn to compare both global and local similarities of the support and query images instead of only using the global features. We also learn to augment features instead of using hand-craft augmentation strategies.

B. Data Augmentation Based Methods

One of the important difficulties in Few-shot Learning is the small number of samples. Data augmentation methods [22], [24]–[28] aim to build a generation model to enhance the variety of the input. Eli Schwartz et al. [25] design an auto-encoder [29] for data augmentation. The core idea of their approach is learning to reconstruct a sample via another one, hence the abundant labeled samples can be used to augment the few-shot input. Combining a meta-learner with a hallucinatory, Wang et al. [22] present a method to augment the input samples by producing additional imaginary data. Chen et al. [28] propose to augment the features by semantic information. Creatively in our method, we learn meaningful affine transformation augmentations from training categories in feature map space, and apply them to augment the feature maps of support image.

C. Salient Patches

Salient patches, which contains rich vision details of the object, is widely used to improve the discrimination of features in various computer vision tasks, such as fine-grained image classification [30], [31], clothes retrieval [32], person re-identification [33], image caption [34] et al. For these application scenarios, predefined salient patches for specific categories could be extracted by the part detector or key-point detector. While in the few-shot scenario, extracting salient patches for new categories is challenging in the situation of few labeled images and lack of prior knowledge. Wang et al. [35] employ the semantic embedding of class tag to generate various local features, and combine them into an image-level feature for Few-shot Learning, which could be limited by the performance of unified visual-semantic embedding. Zhang et al. [36] employed a saliency network pre-trained on MSRA-B to obtain the foregrounds and backgrounds of images, then hallucinated additional datapoints by foreground-background combinations. Chu et al. [37] propose a sampling method based on maximum entropy reinforcement learning to extract various sequences of patches, and aggregate the extracted features for classification. Training the sampler from scratch without effective supervision would extract many background patches and degrades the performance of salient patches. In this paper, we learn to compare by exploiting more *class-relevant* salient patches, and learn to augment the samples by learn *intra-class variations* in feature space, which is significantly different from those related works.

III. AUGMENTED BI-PATH NETWORK

For few-shot classification, the dataset contains a meta-train set and a meta-test set, which have disjoint categories ($C_{meta-train} \cap C_{meta-test} = \emptyset$). The meta-train set, in which each category contains many labeled samples, is used to train a generalized model. Then, it is evaluated on the meta-test set with only few labeled images in each category. In the popular setting, many testing episodes are performed during evaluation process. Each testing episode contains N categories and each category has K labeled images. These labeled data form the

support set $S = \{(\mathbf{x}_{1,1}, y_{1,1}), (\mathbf{x}_{1,2}, y_{1,2}), \dots, (\mathbf{x}_{N,K}, y_{N,K})\}$, where $(\mathbf{x}_{i,j}, y_{i,j})$ is j th datum of i th class. The rest unlabeled testing data of these N categories form the query set $Q = \{\mathbf{x}_{1,K+1}, \mathbf{x}_{1,K+2}, \dots, \mathbf{x}_{N,K+1}, \mathbf{x}_{N,K+2}, \dots\}$. Such setting is called N -way K -shot learning.

We follow the widely used episode-based training strategy [16], which mimics the evaluation process. Similar as each testing episode, a training episode is constructed by random sampling N classes from $C_{meta-train}$. K labeled data of every class sever as the support set, and M labeled data from the rest form the query set. The popular methods [12], [13], [16] train a deep embedding model by feeding the whole image into the deep neural network. To fully exploit the few labeled images, we propose Augmented Bi-path Network (ABNet) for Few-shot Learning, the framework is illustrated in Fig. 2. In total, ABNet includes four modules: 1) Given the support s and query q images, a patch extraction module is first applied to obtain N salient patches for each image. 2) Then, both the whole image and patches are fed into a shared convolution neural network (CNN) $f(\cdot)$ for feature embedding. 3) To enhance the robustness, a learnable feature augmentation module is utilized to augment the features of support images by mimicking the diversity of query images. 4) We learn to compare the features of the support and query images. The similarity maps between two features from the support and query images are computed and re-weighted by a learnable attention module. Finally, combining the global path and local path, the merging module learns to merge the global and local similarity and regress the similarity score of the support and query images. The following subsections are details of the proposed method.

A. Salient Patch Extraction

Salient patches are vital for the comprehensive description of an object, especially in the few-shot scenario. Hence, we develop a salient patch extraction module. Our extraction module starts with the patches sampled by selective search (SS) method [38], which applies bottom-up grouping procedure to generate good object locations capturing all scales. Then elaborate selection method is utilized to distill salient patches.

We measure the importance of patches from two aspects, namely, geometry property and visual salience. Geometry property is referred to the area and aspect ratio of the patch. As small patches usually lack enough discriminate features and large patches are close to the global image, only patches with moderate scale and aspect ratio are significant and will be kept. We use rectangular function $\Pi(x)$ as the geometry property, if the geometry property satisfies the requirement, $\Pi(r_i) = 1$ else $\Pi(r_i) = 0$. Visual salience measures the attraction of pixels to human attention, which is the principal element of importance. The Minimum Barrier Distance (MBD) [39], which represents the connectivity to the background regions, is utilized to measure the visual salience. Hence, the salience of specific patch could be represented as the average of the pixels in it. Taking both geometry property and visual salience

into consideration, the total importance of patch r_i is defined as follows [39]:

$$S(r_i) = \Pi(r_i) \cdot \frac{1}{K} \sum_{j=0}^K v(p_j) \quad (1)$$

$$v(p_j) = \min_{\pi \in \mathbb{S}} [\max_{t=0}^T I(\pi(t)) - \min_{t=0}^T I(\pi(t))]$$

where $v(p_j)$ is the visual salience of the pixel, K is the total number of pixels in patch r_i , $I(\cdot)$ is the pixel value, π is a sequence of pixels where consecutive pairs of pixels are adjacent and the total number of pixels is T , \mathbb{S} is the set of all sequences that connect p_j and seed pixels from background.

Then, the patches could be ranked by the importance, and top N patches could be selected to balance efficiency and effectiveness. Our salient patch extraction is not time consuming as the procedure is performed only once and offline before training. As shown in Fig. 2, the whole image accompanied with the extracted salient patches introduces multi-scale inputs to the feature embedding module (CNN), which improves its representation ability.

B. Feature Embedding

Both the whole image and salient patches are fed into the shared CNN backbone f for feature embedding. Following the classic setting, f consists of four basic convolution blocks. Each block includes a 2D convolution layer with 64 (3x3) kernels, a batch normalization layer and a ReLU nonlinear function. 2x2 max-pooling layer is added after the first two basic convolution blocks. Given image I , we concatenate the feature maps in the last convolution layer and construct a 3D feature map:

$$f(I) = \omega_1 \oplus \omega_2 \oplus \dots \oplus \omega_n, \quad (2)$$

where ω_i is the i th feature map in the output of CNN backbone and \oplus denotes the concatenate operator, $f(I) \in \mathbb{R}^{H \times W \times C}$ with width W , height H and C channels. Different from existing methods [12], [40], where feature embedding is reshaped into one dimensional vector as the input of classifiers, we keep the spatial information in the feature map by leveraging the 3D feature map. As training samples are limited in Few-shot Learning, the additional spatial information will benefit the learning with few samples.

C. Learning to Augment

In Few-shot Learning, the number of training samples is quite limited, while the testing images varies considerably in orientation, viewpoint and clutter background. Hence, even the feature embedding of the same object will be diverse, which results in mis-classification due to the extremely large intra-class distance. To alleviate the problem, we manage to augment the feature embedding of every whole support image I_s and its salient patches in the support set by learning intra-class variations from the categories with sufficient labeled training data. The feature augmentation module aims to learn to transform the support feature to new ones which capture the

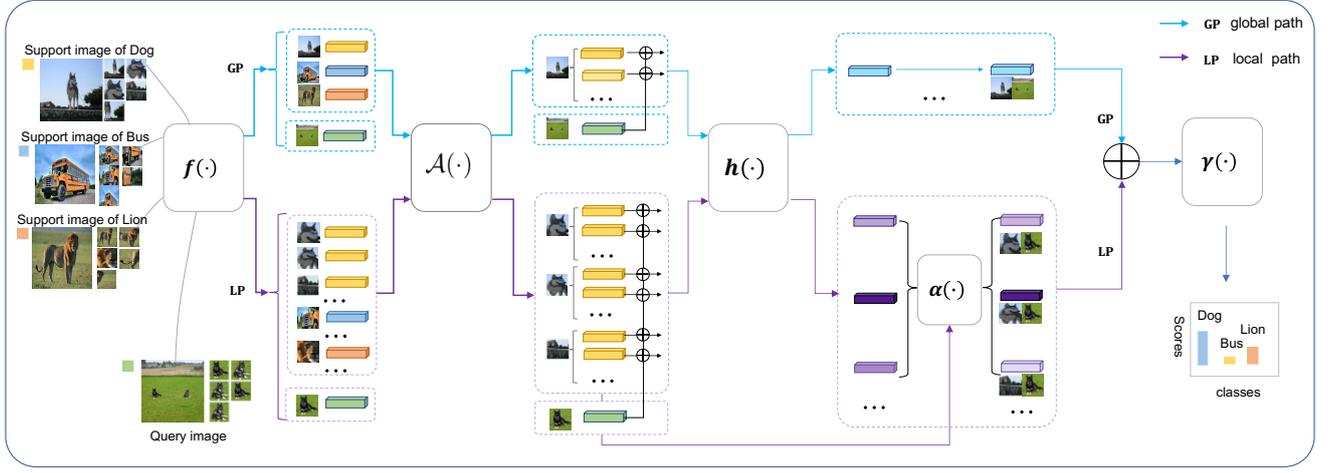


Fig. 2: Illustration of the proposed Augmented Bi-path Network, which contains a global path (GP) for comparing the whole images and a local path (LP) for comparing the extracted salient patches. The whole image accompanied with its salient patches are fed into the network $f(\cdot)$ for feature embedding. To enhance the robustness, a learnable feature augmentation module $\mathcal{A}(\cdot)$ is utilized to augment the support features. Then, similarity maps between augmented support features and the query feature are computed through $h(\cdot)$. Specially for the local path, a re-weighting module $\alpha(\cdot)$ is employed to suppress irrelevant or meaningless patches. Finally, combining GP and LP, the merging module $\gamma(\cdot)$ learns to merge the global and local similarity maps and regress the overall similarity. The label of the query image is predicted as the one with the largest similarity.

diversity of query images. Hence, we introduce a regularization between the generated features and the query features. We factorize the complicated variations into several independent modes, each mode is represented by an affine transformation matrix. For a specific point (x_i, y_i, z_i) in a feature map $f(I)$, the transformed location (x_i^a, y_i^a, z_i^a) in the augmented feature map is defined in the equation:

$$\begin{pmatrix} x_i^a \\ y_i^a \\ z_i^a \end{pmatrix} = \mathcal{A} \begin{pmatrix} x_i \\ y_i \\ z_i \\ 1 \end{pmatrix} = \begin{bmatrix} \mathcal{A}_{11} & \mathcal{A}_{12} & \mathcal{A}_{13} & \mathcal{A}_{14} \\ \mathcal{A}_{21} & \mathcal{A}_{22} & \mathcal{A}_{23} & \mathcal{A}_{24} \\ \mathcal{A}_{31} & \mathcal{A}_{32} & \mathcal{A}_{33} & \mathcal{A}_{34} \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{pmatrix} x_i \\ y_i \\ z_i \\ 1 \end{pmatrix} \quad (3)$$

where \mathcal{A} is an affine transformation matrix with learnable parameters. Different from existing methods [16], [18], we learn to augment feature instead of using hand-crafted augmentation strategy. All the feature maps in support set share the same transformation defined by \mathcal{A} . A new group of transformed output feature maps $f_k(I_s)$ could be generated, which could be translated, scaled, rotated or affined, whatever.

Hence, employing K affine transformation matrices, the original feature embedding $f(I_s)$ could be greatly expanded into $K + 1$ variants defined as Ω :

$$\begin{aligned} \Omega = \mathbb{A}f(I_s) &= \{f_0(I_s), f_1(I_s), \dots, f_K(I_s)\} \\ \mathbb{A} &= \{\mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_K\} \end{aligned} \quad (4)$$

where \mathbb{A} is a group of learned affine transformation matrices (\mathcal{A}_0 is an identity matrix, indicating the original feature embedding, and $f_0(I_s)$ represents $f(I_s)$).

D. Learning to Compare

Inspired by [13], we classify a query image by learning to compare its feature with those of (labeled) support images. Instead of only comparing the support and query image, we propose to compare both the whole image and salient patches. First, we compare the whole image of the support and query

image. Then, we compare the N salient patches of the support and query image. Hence, there are $1 + N^2$ comparison results.

a) *Generating Similarity Maps.*: For comparing two features from support and query images (I_s and I_q), we first concatenate the support feature and every feature variant (Ω_i) of the support image in the channel dimension, which is defined as follows:

$$p_{q,s}^k = f(I_q) \oplus f_k(I_s) \quad (5)$$

Then, a convolution block $g(\cdot)$ is utilized to calculate the similarity and generate the similarity map $g(p_{q,s}^k)$. Each feature embedding in Ω is processed independently in this stage. All of the $K + 1$ similarity maps are concatenated in the channel dimension, and another convolution block $h(\cdot)$ is applied to accumulate the results for calculating the total similarity maps $m_{s,q}$ for I_s and I_q :

$$m_{s,q} = h[g_\phi(p_{q,s}^0) \oplus g(p_{q,s}^1) \oplus \dots \oplus g(p_{q,s}^K)] \quad (6)$$

Similar as the global feature comparison, the comparison between every local patches of support and query image is calculated. In total, $N^2 + 1$ similarity maps are generated:

$$\mathbb{G} = \{m_{s,q}\} \cup \{m_{s,q}^{i,j} | i, j \leq N\} \quad (7)$$

where i, j is positive integer indicating the index of patches, and N is total number of salient patches. In total, we have one global similarity map and N^2 local similarity maps.

b) *Learning to Re-weight*: Instead of using the original similarity maps directly, we learn to re-weight the local similarity maps for emphasizing class-relevant patches and suppressing background patches. The re-weighting value is adaptively predicted by an attention module $\alpha(\cdot)$, which outputs a normalized weight between 0 and 1.

The input of the attention network is a triplet defined as $(f(I_s), f(I_q), m_{s,q}^{i,j})$, which combines the features of

the support and query images and corresponding similarity map. For simply, we use $\alpha(i, j)$ to denote the weight $\alpha(f(I_s), f(I_q), m_{s,q}^{i,j})$. By multiplying the weight $\alpha(i, j)$ to corresponding local similarity maps $m_{s,q}^{i,j}$, a re-weighted similarity group \mathbb{G}^* could be obtained:

$$\mathbb{G}^* = \{m_{s,q}\} \cup \{\alpha(i, j) \cdot m_{s,q}^{i,j} | i, j \leq N\}. \quad (8)$$

c) Learning to Merge: To obtain the overall similarity score between the support and query images, we merge both the global and local similarity maps. All the elements from the weighted similarity group \mathbb{G}^* are concatenated in channel dimension, and regarded as the input to the merge module $\gamma(\cdot)$. The module is utilized to merge these similarity maps and develop a similarity metric, which contains two convolution blocks and two fully-connected layers. The overall similarity $o(s, q)$ between I_q and I_s is defined as:

$$o(s, q) = \gamma\left(\left[\bigoplus_{i,j=1}^{i,j \leq N} \alpha(i, j) \cdot m_{s,q}^{i,j}\right] \oplus m_{s,q}\right), \quad (9)$$

where \oplus denotes the concatenate operator. For few-shot classification, we compute the overall similarities between a query image and all support images. Then, the similarities are averaged for every class. The label is predicted as the one with the largest averaged similarity.

E. Loss Function

To train all the parameters in our model, we minimize the classification loss between the predicted similarity score and ground-truth score. The predicted score $o(i, j)$ is first converted to probabilistic score by Sigmoid function:

$$P(i, j) = \frac{1}{1 + e^{-o(i, j)}} \quad (10)$$

Then the classification loss function is computed as :

$$\mathcal{L}_{cls} = \frac{1}{B \times C} \sum_{i=1}^B \sum_{j=1}^C (P(i, j) - \mathbb{I}(y_i == y_j))^2, \quad (11)$$

where B and C are the numbers of query and support images respectively. y_i and y_j are the class labels of the query (I_q) and support (I_s) images. $\mathbb{I}(x)$ is an indicator function, $\mathbb{I}(x) = 1$ if x is true and 0 otherwise.

To learn sparse attention values for salient patches, we add the L_1 regularization. The loss is defined as follows:

$$\mathcal{L}_{att} = \frac{S_{att}}{N^2} \sum_{i=1}^N \sum_{j=1}^N |\alpha(i, j)|, \quad (12)$$

where N is the number of salient patches, S_{att} is scaling factor to make the loss in the same scale.

Besides, we introduce the augmentation loss to regularize the feature augmentation. The feature augmentation module tries to learn the diversity of query images and generate new features with such diversity. Hence, we introduce a regularization between the generated features and the query features. During training, we have labels of both support and query

images. For the query image I_q , if it has the same label as the support image I_s , we compute the mean square error (MSE) between $f(I_q)$ and augmented features $\{f_k(I_s)\}_{1 \leq k \leq K}$:

$$\mathcal{L}_{aug} = \frac{S_{aug}}{B \times C} \sum_{i=1}^B \sum_{j=1}^C \sum_{k=1}^K (f(I_q) - f_k(I_s))^2 \mathbb{I}(y_i == y_j), \quad (13)$$

where K is the total number of augmented features and S_{aug} is the scaling factor for augmentation loss.

Finally, the total loss is computed as the weighted sum of the above losses:

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda_{att} \cdot \mathcal{L}_{att} + \lambda_{aug} \cdot \mathcal{L}_{aug}, \quad (14)$$

where λ_{att} and λ_{aug} are the weights of the attention loss and augmentation loss.

IV. EXPERIMENT

In this section, we first compare to the state-of-the-art methods on three popular datasets. Then, we show that our method is generalized and can achieve good results on different learning settings. Latter, blation study with visualization is provided to illustrate that every proposed module is important.

A. Datasets

We do experiments on three datasets, namely, miniImageNet [16], Caltech-256 [41] and tieredImageNet [26]. miniImageNet is the most popular benchmark for few-shot classification, Caltech-256 is a larger dataset with more categories, and tieredImageNet is a more challenging dataset with test classes that are less similar to training ones.

miniImageNet. The dataset contains 100 classes with 600 images per class [16]. Objects in images have variable appearances, positions, viewpoints, poses as well as background clutter and occlusion. We follow the popular split [21], where 64 classes are for training, 16 classes are for validation and the rest 20 classes are for testing. All images are resized to 84×84 size.

Caltech-256. This dataset [41] contains 30,607 images from 256 object categories. These categories are diverse, ranging from grasshopper to tuning fork. We follow the split for FSL which is provided by [19]. The training, validation and testing sets include 150, 56 and 50 classes respectively. The same to miniImageNet, all images in Caltech-256 are resized to 84×84 size.

tieredImageNet. This dataset [26] contains 608 classes (779,165 images) grouped into 34 higher-level nodes from the ImageNet human-curated hierarchy. This set of nodes is partitioned into 20, 6, and 8 disjoint sets of training, validation, and testing nodes, and the corresponding classes constitute the respective meta-sets. All images are resized to 84×84 size.

B. Implementation Details

During salient patches extraction, the size requirement is set to be 1%~50% and the aspect ratio requirement is set to be 1/3~3. To further remove duplicate patches, Non-Maximum Suppression (NMS) technique is applied to keep smaller ones.

Method	miniImageNet		Caltech-256		tieredImageNet		
	5way1shot	5way5shot	5way1shot	5way5shot	5way1shot	5way5shot	
MatchingNet [16]	NIPS'16	43.56±0.84	55.31±0.73	45.59±0.77	54.61±0.73	54.02	70.11
MetaLSTM [21]	ICLR'17	43.44±0.77	60.60±0.71	-	-	-	-
MAML [20]	ICML'17	48.70±0.84	55.31±0.73	48.09±0.83	57.45±0.84	51.67±1.81	70.30±0.08
MetaNet [17]	ICML'17	49.21±0.96	-	-	-	-	-
ProtoNet [12]	NIPS'17	49.42±0.87	68.20±0.70	-	-	54.28±0.67	71.42±0.61
RelationNet [13]	CVPR'18	50.44±0.82	65.32±0.77	56.12±0.94	73.04±0.72	54.48±0.93	71.32±0.78
CTM [42]	CVPR'19	41.62	58.77	-	-	-	-
Spot&Learn [37]	CVPR'19	51.03±0.78	67.96±0.71	-	-	-	-
MetaOptNet [43]	CVPR'19	52.87±0.57	68.76±0.48	-	-	54.71±0.67	71.79±0.59
ABNet		58.12±0.94	72.02±0.75	63.20±0.99	78.42±0.69	62.10±0.96	75.11±0.78

TABLE I: Few-shot classification accuracy (%) on miniImageNet, Caltech-256 and tieredImageNet datasets. The results are averaged over 600 testing episodes, and the 95% confidence intervals are reported. We compare to methods using the same **4-layer** feature embedding module, *i.e.*, ($64 \times 64 \times 64 \times 64$).

Then, the extracted patches are ranked by the visual saliency, and top five salient patches are selected for each image. For those images with less than five extracted salient patches, we pad the number to five by duplicating.

For fair comparison, we employ the most widely used 4-layer convolution module [12], [13], [21] with 64 filters in each convolution layer as the backbone. The architecture is ($64 \times 64 \times 64 \times 64$). This embedding module generates 64 (19×19) feature maps for each input image or patch. Unless specified, all experiments are implemented with this 4-layer backbone. We also provide the results with ResNet backbone on miniImageNet.

In feature augmentation module, four affine transformation matrices are utilized to learn augmentations from the training data. The module g and h for generating similarity maps contains the same basic convolution block as the backbone, and an additional 2×2 max-pooling layer is utilized in h . The triplet input of attention module α is global-average pooled in spatial dimension and further concatenated into an one dimensional vector. Then three fully-connected layers are employed to learn the attention value. Finally, in the merging module γ , one basic convolution block followed by two fully-connected layers are applied to regress the similarity score. Besides, ReLU is used as the default activation function in all fully-connected layers except the output layer of attention module and merging module, where Sigmoid is used to normalize the score to be (0, 1). When training shot > 1 , *e.g.* 5-shot learning, we average the similarity scores of all training shots.

We train the model from scratch, and no extra data or pretrained models are used. Adam [44] optimizer is used during training with initial learning rate 0.01 for feature augmentation module and 0.001 for the others. The learning rate is decayed in half every 100,000 episodes. λ_{att} and λ_{aug} are both set to be 0.1. Following [13], we report the averaged testing accuracy over 600 episodes, and each episode contains 15 query images of every class.

C. Comparison to the State-of-the-art

We compare our approach to several state-of-the-art methods under 5-way-1-shot and 5-way-5-shot settings. The results are shown in Table I. As a deeper backbone with higher resolution input image will always increase the classification

Method	Backbone	miniImageNet	
		5way1shot	5way5shot
RelationNet [13]	CVPR'18 ResNet-18	58.21	74.29
CTM [42]	CVPR'19 ResNet-18	62.05±0.55	78.63±0.06
MetaOptNet [43]	CVPR'19 ResNet-12	62.64±0.61	78.63±0.46
ABNet	ResNet-18	63.15±0.63	78.85±0.56

TABLE II: Few-shot classification accuracy (%) on miniImageNet with ResNet backbones.

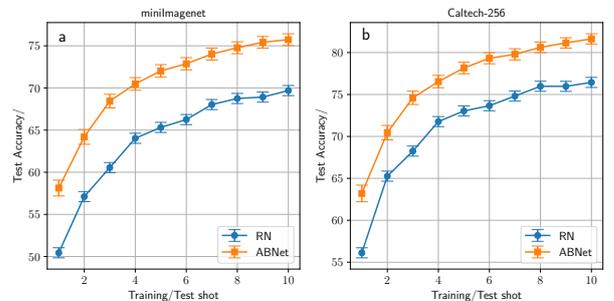


Fig. 3: The 5-way n -shot classification accuracy on miniImageNet (left) and Caltech-256 (right). The 4-layer backbone is used. Orange points are results of our ABNet, while cyan points are results of RelationNet.

performance by a large margin [45], [46], we compare to methods using the same backbone ($64 \times 64 \times 64 \times 64$) for fair comparison.

According to Table I, our ABNet achieves the best performance in all settings. Our method outperforms the runner-up (MetaOptNet [43]) by 5.25% (1-shot) and 3.26% (5-shot) on miniImageNet and 7.39% (1-shot) and 3.32% (5-shot) on tieredImageNet. In addition, our method also shows significant improvements on Caltech-256 dataset. Compared to RelationNet [13], our method achieves 7.08% and 5.38% improvements on 1-shot and 5-shot learning settings respectively.

Some works employ deeper networks, *e.g.* ResNet, to extract features. We also provide the results on miniImageNet with ResNet-18 as the backbone in Table II. Obviously, our method achieves the best performances on two settings. It is interesting that the gaps between different methods (CTM [42], MetaOptNet [43] and ours) with ResNet backbone are not large. The main reason may be that the features extracted by deep networks are already good enough for few-shot classification.

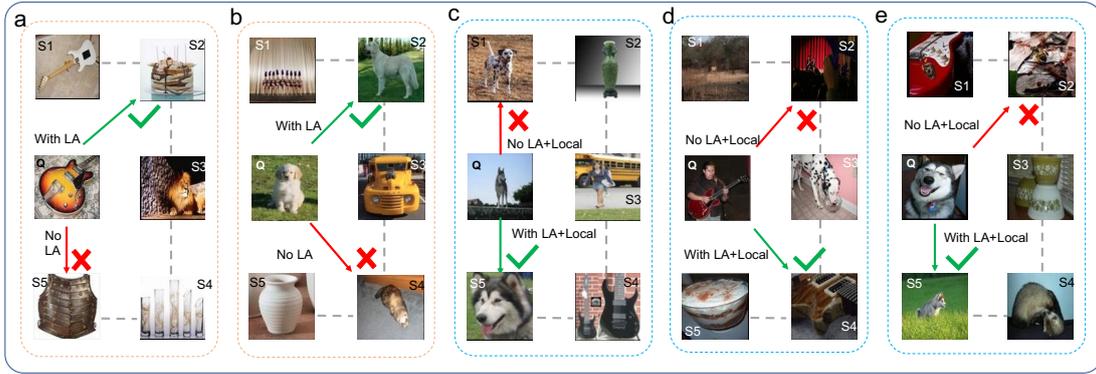


Fig. 4: Visualization of 5-way 1-shot classification results on miniImageNet. LA denotes Learning to Augmentation module. Local stands for local features which is produced by Salient Patch module. Correct predictions are made by the model with feature augmentation when the orientation (a) and viewpoint (b) varies. When the scale of object varies (c), even both scale and viewpoint/orientation varies in (d)/(e), the combination of feature augmentation and local features in ABNet leads to the correct classification.

D. Generalization Ability

To verify that our method is generalized to settings with different training shots, we provide the detailed comparison to RelationNet [13] under 5-way- n -shot setting where $n \in [1, 10]$. As illustrated in Fig. 3, along with the increase of training shot, the performance of our ABNet stably improves. On both two datasets, our method always significantly outperforms RelationNet under all settings. It demonstrates that our method has good generalization ability on n -shots tasks.

E. Ablation Studies

To study the effectiveness of Learning to Augment (LA), Salient Patch (SP) and Learning to Re-weight (LR) modules independently, we perform quantitative comparison and visual comparison on miniImageNet dataset.

Quantitative comparison. We start with a baseline model, which only employs global features with handcrafted feature augmentations, e.g. horizontal flip, rotation of $\pm\pi/2$ and π . Then we gradually add LA, SP and LR modules. Then, we evaluate the four variants of our method.

- Baseline: global features with handcrafted augmentations
- Baseline+LA: Baseline with Learning to Augment
- Baseline+LA+SP: Baseline with Learning to Augment and Salient Patch
- Baseline+LA+SP+LR: Baseline with Learning to Augment, Salient Patch and Learning to Re-weight

The classification accuracies for 1-shot and 5-shot learning are evaluated and shown in Table III. Compared to baseline model, the accuracy with Learn to Augment module (Baseline+LA) improves by 1.83% in 1-shot learning and 1.52% in 5-shot learning. It means learnable feature augmentations could acquire more appropriate variations from query images. When Salient Patch model is further introduced (Baseline+LA+SA), the performance is boosted by 1.49% in 1-shot learning and 1.68% in 5-shot learning. Hence, the importance of local features is verified. Finally, a significant improvement of 2.36% and 2.32% is observed after introducing Learning

to Re-weight module (Baseline+LA+SA+LR), which demonstrates the importance of re-weighting local similarity maps before merging. As class-irrelevant patches or background patches could be introduced inevitably by the unsupervised salient patch extraction method, the novel Learning to Re-weight module could emphasize more relevant local similarities by further taking features into consideration.

Model	miniImageNet	
	5way1shot	5way5shot
Baseline	52.44 \pm 0.91	66.50 \pm 0.77
Baseline+LA	54.27 \pm 0.91	68.02 \pm 0.77
Baseline+LA+SP	55.76 \pm 0.89	69.70 \pm 0.72
Baseline+LA+SP+LR	58.12\pm0.94	72.02\pm0.75

TABLE III: Few-shot classification accuracy (%) for ablation studies. The results are averaged over 600 test episodes with the 95% confidence intervals. The baseline model is trained with handcrafted feature augmentation (horizontal flip and rotation), LA: Learning to Augment, SP: Salient Patch, LR: Learning to Re-weight.

Visual Comparison. The effectiveness of Learning to Augment and Salient Patches is also demonstrated by visualizing the 1-shot classification results in Fig. 4. As the orientation and viewpoint of specific object varies in the wild, it is very difficult to predict correct category merely based on one glance (shot) of the object. For example, when the “guitar” image in the query set is placed in totally different orientation from the one in support set (see Fig. 4(a)), one-shot classifier without feature augmentation fails to recognize it and predicts the wrong category. Similar as the “dog” image with different viewpoints in Fig. 4(b). However, our method can learn to augment the support image feature and predict the correct category. Moreover, when the scale of object in the query image is significantly different from the one in support image, classifier with merely global features no longer works well and local features are required to make correct prediction. Taking the “dog” image as an example, when close-shot image is compared to long-shot image in Fig. 4(c), the image of a different specie (category) but similar scale is mis-matched by the method using only global features. In contrast, ABNet with local features could easily recognize the correct specie

(category) even the two objects have significantly different scale. Even with different scales, viewpoints (see Fig. 4(d)) and orientations (see Fig. 4(e)), ABNet can predict the correct category. The above improvements benefit from the meta-learning ability of Learning to Augment and Learning to Compare modules with salient patches.

V. CONCLUSION

We propose a novel meta-learning based method, namely Augmented Bi-path Network, for Few-shot Learning. The proposed method extends the previous “learn-to-compare” based methods by introducing both global and local features on multi-scales. Experimental results show that our method significantly outperforms the state-of-the-art on three challenging datasets under all settings. Ablation studies verify the importance of the proposed Learning to Augment, Salient Patch and Learning to Re-weight modules. We also provide visual comparison to illustrate how these modules can improve FSL performance.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *NeurIPS*, 2012.
- [2] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *CVPR*, 2015, pp. 1–9.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *CVPR*, 2009, pp. 248–255.
- [6] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *ECCV*. Springer, 2014, pp. 740–755.
- [7] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci, T. Duerig *et al.*, “The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale,” *arXiv:1811.00982*, 2018.
- [8] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona, “Caltech-UCSD Birds 200,” California Institute of Technology, Tech. Rep. CNS-TR-2010-001, 2010.
- [9] V. Gulshan, L. Peng, M. Coram, M. C. Stumpe, D. Wu, A. Narayanaswamy, S. Venugopalan, K. Widner, T. Madams, J. Cuadros, R. Kim, R. Raman, P. C. Nelson, J. L. Mega, and D. R. Webster, “Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs,” *JAMA*, vol. 316, no. 22, pp. 2402–2410, 12 2016.
- [10] L. Fei-Fei *et al.*, “A bayesian approach to unsupervised one-shot learning of object categories,” in *CVPR*, 2003, pp. 1134–1141.
- [11] G. Koch, R. Zemel, and R. Salakhutdinov, “Siamese neural networks for one-shot image recognition,” in *ICML Workshop*, 2015.
- [12] J. Snell, K. Swersky, and R. Zemel, “Prototypical networks for few-shot learning,” in *NeurIPS*, 2017, pp. 4077–4087.
- [13] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, “Learning to compare: Relation network for few-shot learning,” in *CVPR*, Jun 2018.
- [14] B. Zhao, X. Sun, Y. Fu, Y. Yao, and Y. Wang, “Msplit lbi: Realizing feature selection and dense estimation simultaneously in few-shot and zero-shot learning,” in *ICML*, 2018.
- [15] L. Fei-Fei, R. Fergus, and P. Perona, “One-shot learning of object categories,” *PAMI*, vol. 28, no. 4, 2006.
- [16] O. Vinyals, C. Blundell, T. Lillicrap, k. kavukcuoglu, and D. Wierstra, “Matching networks for one shot learning,” in *NeurIPS*, 2016.
- [17] T. Munkhdalai and H. Yu, “Meta networks,” in *ICML*, 2017.
- [18] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, “Meta-learning with memory-augmented neural networks,” in *ICML*, Jun 2016, pp. 1842–1850.
- [19] F. Zhou, B. Wu, and Z. Li, “Deep meta-learning: Learning to learn in the concept space,” *arXiv preprint arXiv:1802.03596*, 2018.
- [20] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *ICML*, 2017, pp. 1126–1135.
- [21] S. Ravi and H. Larochelle, “Optimization as a model for few-shot learning,” in *ICLR*, 2017.
- [22] Y.-X. Wang, R. Girshick, M. Hebert, and B. Hariharan, “Low-shot learning from imaginary data,” in *CVPR*, 2018, pp. 7278–7286.
- [23] A. Nichol, J. Achiam, and J. Schulman, “On first-order meta-learning algorithms,” *arXiv preprint arXiv:1803.02999*, 2018.
- [24] Y.-X. Wang and M. Hebert, “Learning from small sample sets by combining unsupervised meta-training with cnns,” in *NeurIPS*, 2016.
- [25] E. Schwartz, L. Karlinsky, J. Shtok, S. Harary, M. Marder, R. Feris, A. Kumar, R. Giryes, and A. M. Bronstein, “Delta-encoder: an effective sample synthesis method for few-shot object recognition,” in *NeurIPS*, 2018.
- [26] M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. B. Tenenbaum, H. Larochelle, and R. S. Zemel, “Meta-learning for semi-supervised few-shot classification,” in *ICLR*, 2018. [Online]. Available: <http://arxiv.org/pdf/1803.00676v1>
- [27] R. Zhang, T. Che, Z. Ghahramani, Y. Bengio, and Y. Song, “Metagan: An adversarial approach to few-shot learning,” in *NeurIPS*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., 2018, pp. 2365–2374.
- [28] Z. Chen, Y. Fu, Y. Zhang, Y.-G. Jiang, X. Xue, and L. Sigal, “Multi-level semantic feature augmentation for one-shot learning,” *TIP*, vol. 28, pp. 4594–4605, 2018.
- [29] C.-Y. Liou, J.-C. Huang, and W.-C. Yang, “Modeling word perception using the elman network,” *Neurocomputing*, vol. 71, no. 16, 2008.
- [30] N. Zhang, J. Donahue, R. Girshick, and T. Darrell, “Part-based r-cnns for fine-grained category detection,” in *ECCV*. Springer, 2014.
- [31] H. Zhang, T. Xu, M. Elhoseiny, X. Huang, S. Zhang, A. Elgammal, and D. Metaxas, “Spda-cnn: Unifying semantic part detection and abstraction for fine-grained recognition,” in *CVPR*, 2016, pp. 1143–1152.
- [32] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang, “Deepfashion: Powering robust clothes recognition and retrieval with rich annotations,” in *CVPR*, 2016, pp. 1096–1104.
- [33] L. Zheng, Y. Huang, H. Lu, and Y. Yang, “Pose-invariant embedding for deep person re-identification,” *TIP*, vol. 28, no. 9, Sep. 2019.
- [34] Z. Gan, C. Gan, X. He, Y. Pu, K. Tran, J. Gao, L. Carin, and L. Deng, “Semantic compositional networks for visual captioning,” *CVPR*, 2016.
- [35] P. Wang, L. Liu, C. Shen, Z. Huang, A. van den Hengel, and H. T. Shen, “Multi-attention network for one shot learning,” in *CVPR*, 2017.
- [36] H. Zhang, J. Zhang, and P. Koniusz, “Few-shot learning via saliency-guided hallucination of samples,” in *CVPR*, 2019, pp. 2770–2779.
- [37] W.-H. Chu, Y.-J. Li, J.-C. Chang, and Y.-C. F. Wang, “Spot and learn: A maximum-entropy patch sampler for few-shot image classification,” in *CVPR*, 2019, pp. 6251–6260.
- [38] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders, “Selective search for object recognition,” *IJCV*, vol. 104, no. 2, 2013.
- [39] J. Zhang, S. Sclaroff, Z. Lin, X. Shen, B. Price, and R. Mech, “Minimum barrier salient object detection at 80 fps,” in *CVPR*, 2015.
- [40] L. Bertinetto, J. F. Henriques, P. Torr, and A. Vedaldi, “Meta-learning with differentiable closed-form solvers,” in *ICLR*, 2019.
- [41] G. Griffin, A. Holub, and P. Perona, “Caltech-256 object category dataset,” Tech. Rep., 2007.
- [42] H. Li, D. Eigen, S. Dodge, M. Zeiler, and X. Wang, “Finding task-relevant features for few-shot learning by category traversal,” in *CVPR*, 2019, pp. 1–10.
- [43] K. Lee, S. Maji, A. Ravichandran, and S. Soatto, “Meta-learning with differentiable convex optimization,” in *CVPR*, 2019, pp. 10 657–10 665.
- [44] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *ICLR*, 2015.
- [45] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel, “A simple neural attentive meta-learner,” *arXiv preprint arXiv:1707.03141*.
- [46] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *ICML*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. Long Beach, California, USA: PMLR, 09–15 Jun 2019, pp. 6105–6114.