# Few-Shot Learning Meets Transformer: Unified Query-Support Transformers for Few-Shot Classification

Xixi Wang, Xiao Wang, *Member, IEEE*, Bo Jiang*, Bin Luo, *Senior Member, IEEE*

arXiv:2208.12398v1 [cs.CV] 26 Aug 2022

*Abstract*—Few-shot classification which aims to recognize unseen classes using very limited samples has attracted more and more attention. Usually, it is formulated as a metric learning problem. The core issue of few-shot classification is how to learn (1) consistent representations for images in both support and query sets and (2) effective metric learning for images between support and query sets. In this paper, we show that the two challenges can be well modeled simultaneously via a unified Query-Support TransFormer (QSFormer) model. To be specific, the proposed QSFormer involves *global* query-support sample Transformer (sampleFormer) branch and *local* patch Transformer (patchFormer) learning branch. sampleFormer aims to capture the dependence of samples in support and query sets for image representation. It adopts the Encoder, Decoder and Cross-Attention to respectively model the Support, Query (image) representation and Metric learning for few-shot classification task. Also, as a complementary to global learning branch, we adopt a local patch Transformer to extract structural representation for each image sample by capturing the long-range dependence of local image patches. In addition, a novel Cross-scale Interactive Feature Extractor (CIFE) is proposed to extract and fuse multi-scale CNN features as an effective backbone module for the proposed few-shot learning method. All modules are integrated into a unified framework and trained in an end-to-end manner. Extensive experiments on four popular datasets demonstrate the effectiveness and superiority of the proposed QSFormer.

*Index Terms*—Few-Shot Learning, Transformer, Metric Learning, Deep Learning.

## I. INTRODUCTION

CURRENT deep neural networks learn from large-scale training samples and achieve good performance on many tasks. However, in many scenarios, data collection and annotation is expensive and it is usually very challenging to collect enough data for the training of deep neural networks. The Few-shot classification aims to recognize unseen/query classes by using very limited seen/support samples has attracted more and more attention.

Many deep learning methods [1]–[3] have been proposed to address few-shot learning problem. These methods can be roughly classified into three types, i.e., generation-based methods, optimization-based methods and metric-based methods. Metric-based methods are derived to distinguish support and query samples by using some image representation and metric learning techniques. As we know, the core issues for metric-based few-shot classification are two aspects: 1) How to learn

The authors are all from School of Computer Science and Technology, Anhui University, Hefei 230601, China
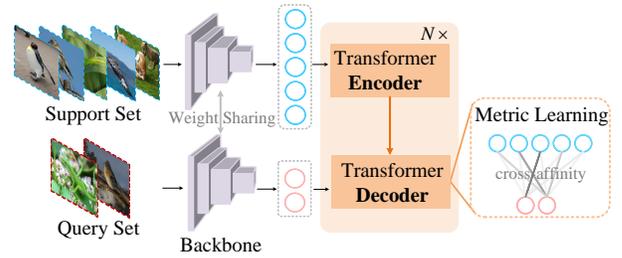Corresponding author: Bo Jiang



Fig. 1. Illustration of our proposed unified Query-Support Transformer for few-shot learning. It models the feature engineering on query/support samples and metric learning simultaneously.

consistent representations for images in both support and query sets. 2) How to conduct effective metric learning for images between support and query sets. According to our observation, existing works [3]–[7] usually first employ Convolution Neural Networks (CNNs) to learn image feature representation and then use a metric function to directly compute the similarities (e.g., cosine) between query and support images for few-shot classification. The good performance can be achieved, however, many recent studies [8], [9] demonstrate that CNN only captures the local relations well due to its limited receptive field. To address this issue, some researchers [10]–[12] propose to combine or replace CNN with Transformer networks to model the long-range relationships of local image patches and obtain better image representation results. However, they may still obtain sub-optimal performance due to the following two reasons: 1) Existing works generally adopt Transformers (or CNN+Transformer) as the backbone network for engineering each image representation, which obviously ignores the inherent relationships among samples in query and support sets for image representation. 2) Existing works generally adopt the two-stage learning scheme, i.e., 'representation learning + metric learning'. Although the two stages are usually learned together in an end-to-end manner, this decoupling way may lead to sub-optimal learning results.

To address these challenges, in this work, we propose a unified Query-Support Transformer architecture for few-shot learning, termed QSFormer. The core of QSFormer is our new design of query-support sample Transformer (named sampleFormer) module, which aims to explore the relationships of samples for coupling **sample representations** and **metric learning of samples** together in a unified module for few-shot classification. To be specific, as shown in Figure 1, we

dexterously adopt the *Encoder, Decoder and Cross-Attention in our sampleFormer architecture to model the Support, Query (image) representation and Metric learning in few-shot classification task*, respectively. For the support branch, we represent all support images as a sequence of image tokens and feed them into the Transformer encoder to enhance the support features. For the query branch, it receives a sequence of query image tokens to learn their representations. Meanwhile, it interacts with the previous support branch via the cross-attention for modeling the similarities/affinities between query and support tokens, therefore, naturally achieving metric learning in the decoding procedure.

Based on our newly proposed sampleFormer, we further extend it by introducing two additional new modules for high-performance few-shot learning, including Cross-scale Interactive Feature Extractor (CIFE) and local patch Transformer (patchFormer) module. Specifically, as shown in Figure 2, given the query and support images, we first use CIFE as the backbone module to extract the image features. Then, the sampleFormer takes the embedded image tokens as input and outputs global metrics. Meanwhile, the local/patch correspondence of query-support image pairs is also considered using the patchFormer. The global and local metrics are combined for few-shot classification. Note that, the whole network can be optimized in an end-to-end way.

To sum up, the contributions of this paper can be summarized as follows:

- We propose a unified Query-Support Transformer (termed QSFormer) for few-shot learning, which models the representation learning and metric learning simultaneously.
- We propose a novel Sample Transformer module (sampleFormer) to capture the sample relationships in few-shot problem setting. Also, we propose a patch Transformer (patchFormer) module for few-shot image representation and metric learning.
- We propose a Cross-scale Interactive Feature Extractor for image representation by considering the interaction of different CNN levels.
- Extensive experiments on four widely used few-shot classification datasets demonstrate the effectiveness and superiority of our proposed method.

## II. RELATED WORK

**Few-shot Learning.** Current few-shot learning algorithms can be broadly divided into two categories: optimization-based approaches [2], [6] and metric-based approaches [3], [4], [13], [14]. Our method is more relevant to the metric-based approaches, which mainly focus on the representation learning and metric learning of samples. Specifically, Sung et al. [15] propose a Relation Network (RN) for few-shot learning, which computes the relation scores between query examples and the few examples of each new class to classify the examples of new classes. Hou et al. [13] develop a Cross Attention Network, which highlights the target object regions to enhance the feature representation by producing cross attention maps for each feature. Zhang et al. [3] introduce Earth Mover's Distance to capture a structural distance between the

local image representations for few-shot classification. Xie et al. [14] introduce a deep Brownian Distance Covariance approach to learn image representations and then use distance metric for classification.

**Transformer for Few-shot Classification.** Transformer [16] has universal modeling capability because its core module self-attention learning mechanism. In recent years, Transformer has been employed by a large number of researchers for various visual tasks, including object tracking [17], [18], object detection [19], [20], object re-identification [21], [22], multi-label classification [23], [24], Medical Image Segmentation [25], [26], and so on. For few-shot learning tasks, some works [10]–[12], [27]–[29] demonstrate that Transformer architecture is also promising. For example, Ye et al. [27] develop a Few-Shot Embedding Adaptation Transformer (FEAT) to instantiate set-to-set transformation and thus make instance embedding task-specific for few-shot learning. Liu et al. [28] propose a Universal Representation Transformer (URT) layer by combining feature representations from multiple domains together for multi-domain few-shot classification. Zhmoginov et al. [12] introduce a transformer-based model, called HyperTransformer (HT), which encodes task-dependent variations in the weights of a small CNN model for few-shot learning. These works mainly employ Transformer architecture for representation learning. Differently, in our work, we develop a Query-Support Transformer (QSFormer) to accomplish both feature representation and metric learning simultaneously.

## III. THE PROPOSED METHOD

The purpose of few-shot classification is to classify the unseen samples when only a small number of samples are available. Many recent approaches [3], [13], [30], [31] indicate that the episode mechanism provides an effective way for few-shot classification task and we follow them in both training and testing phases. Formally, let $\mathcal{D}_{train}$, $\mathcal{D}_{val}$ and $\mathcal{D}_{test}$ respectively represent meta-training, meta-validation and meta-testing set, where $\mathcal{D}_{train} \cap \mathcal{D}_{val} \cap \mathcal{D}_{test} = \emptyset$. Taking $C$-way $K$-shot few-shot classification task as an example, each episode consists of support set $\mathcal{X}^s = \{(X_i^s, Y_i^s)\}_{i=1}^{n_s}$ and query set $\mathcal{X}^q = \{(X_j^q, Y_j^q)\}_{j=1}^{n_q}$. Concretely, we randomly select $C$ classes and $K$ labeled samples per class to form the support set $\mathcal{X}^s$, i.e., $n_s = C \times K$. Meanwhile, we randomly sample $q$ samples per class to form the query set $\mathcal{X}^q$, i.e., $n_q = C \times q$.

As shown in Figure 2, we propose a novel Query-Support Transformer (QSFormer) framework for few-shot learning, which contains the following four parts:

- **Cross-Scale Interactive Feature Extractor (CIFE):** we propose a cross-scale interactive feature extractor as backbone network to obtain the spatial enhanced support/query CNN feature representations.
- **Sample Transformer Module:** we introduce a query-support sample Transformer (sampleFormer) module to couple image sample representation and global metric learning of samples together for few-shot learning.
- **Patch Transformer Module:** we also propose a patch Transformer (patchFormer) module to model the context
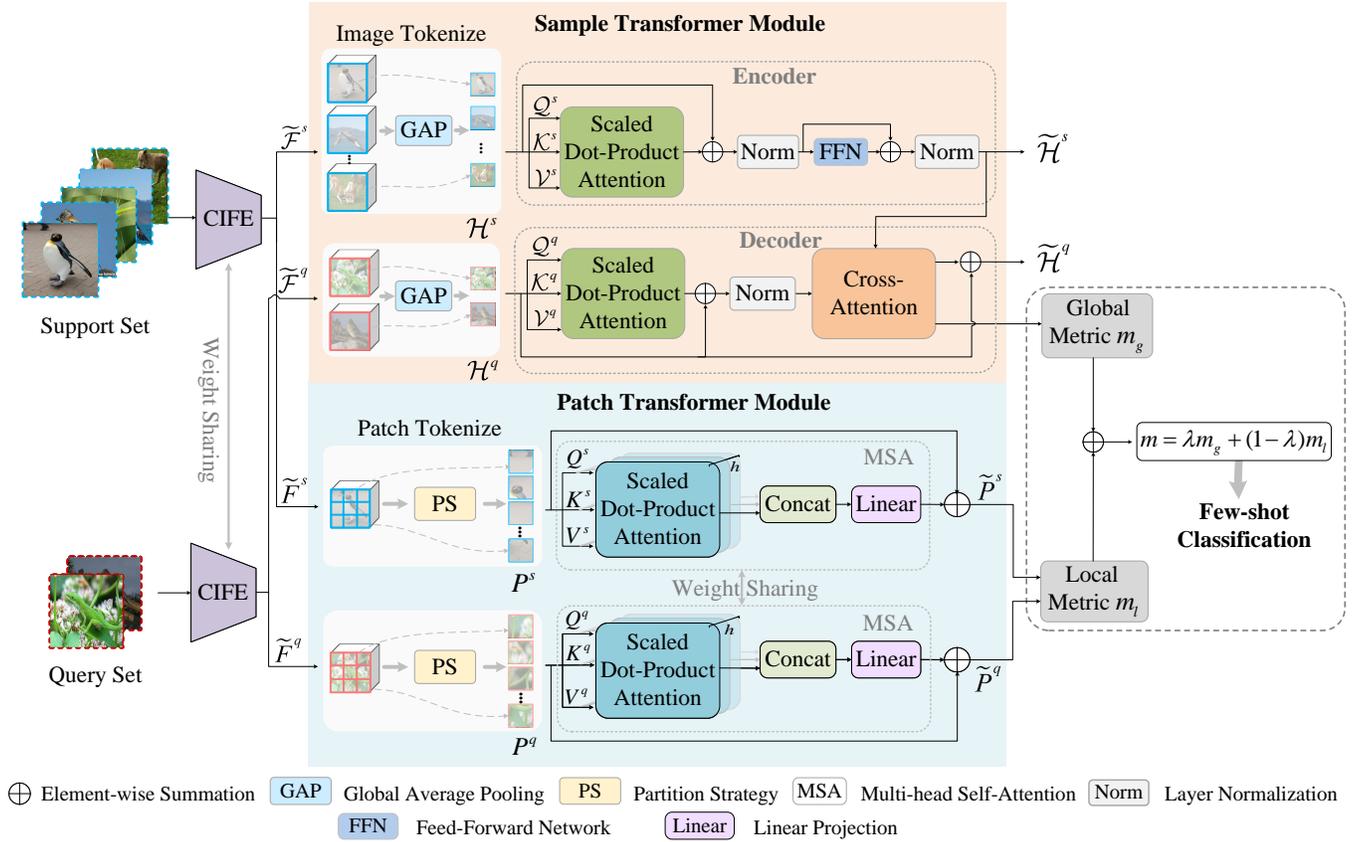
Fig. 2. An overview of the proposed QSFormer framework, which mainly consists of Cross-scale Interactive Feature Extractor (CIFE), Sample Transformer Module, Patch Transformer Module, Metric Learning and Few-shot Classification. More details can be found in Section III.

correlation of patches in each image sample to conduct the local metric learning between query-support sample pairs.

- **Metric Learning and Few-shot Classification:** we acquire the final metric by combining global metric obtained via sampleFormer and local metric obtained via patchFormer together and final achieve few-shot classification.

Below, we introduce the details of these modules.

### A. Cross-scale Interactive Feature Extractor

We introduce a novel Cross-scale Interactive Feature Extractor (CIFE) as backbone module, which aims to obtain the ego-context CNN feature representations for support and query samples.

As shown in Figure 3, taking the support image set $\mathcal{X}^s = \{X_1^s, X_2^s, ..., X_{n_s}^s\}$ as inputs, we first use the pre-trained ResNet-12 to generate the initial multi-scale feature representations $\mathcal{F}_l^s \in \mathbb{R}^{n_s \times c_l \times h_l \times w_l}, l \in \{1, 2, 3, 4\}$, where $n_s$ represents the number of support samples in each episode and $c_l$, $h_l$ and $w_l$ denote the channel, height and width of support feature map in the $l$-th level respectively. Then, we employ a Transformer architecture [16] consisting of multi-head self-attention (MSA), layer normalization (LN), feed-forward network (FFN) and residual connection to achieve the interaction of multi-scale features. Finally, we can obtain the
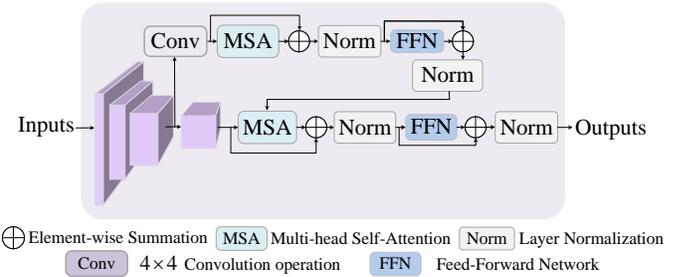


Fig. 3. Illustration of Cross-scale Interactive Feature Extractor (CIFE) for feature extraction.

spatial enhanced feature representations for support samples as $\widetilde{\mathcal{F}}^s = \{\widetilde{F}_1^s, \widetilde{F}_2^s, \cdots, \widetilde{F}_{n_s}^s\} \in \mathbb{R}^{n_s \times c \times h \times w}$. Similarly, we obtain the spatial enhanced features for query samples as $\widetilde{\mathcal{F}}^q = \{\widetilde{F}_1^q, \widetilde{F}_2^q, \cdots, \widetilde{F}_{n_q}^q\} \in \mathbb{R}^{n_q \times c \times h \times w}$. The parameters of CIFE are shared for support and query branches. In practice, we empirically set $c = 640$ and $h = w = 5$.

### B. Sample Transformer Module

To achieve both image sample representation and metric learning of samples in a unified module, we design a novel query-support sample Transformer module, named sampleFormer. The proposed sampleFormer mainly consists of Encoder and Decoder, as shown in Figure 2.

**Encoder.** The purpose of the Encoder is to mine the relationships of samples in support set to obtain better support feature representations. To this end, based on the aforementioned support features $\widetilde{\mathcal{F}}^s \in \mathbb{R}^{n_s \times c \times h \times w}$, we first introduce *image tokenize*, which utilizes a global average pooling and reshape operation to gain the token sequence $\mathcal{H}^s = \{H_1^s, H_2^s, \cdots, H_{n_s}^s\} \in \mathbb{R}^{n_s \times c}$ of support samples, where each token $H_i^s$ denotes a support image sample. As shown in Figure 2, we can see that the main component of encoder is *attention mechanism*, whose inputs are Query $\mathcal{Q}^s \in \mathbb{R}^{n_s \times c}$, Key $\mathcal{K}^s \in \mathbb{R}^{n_s \times c}$, and Value $\mathcal{V}^s \in \mathbb{R}^{n_s \times c}$ obtained by conducting three linear projections on $\mathcal{H}^s$ respectively. Next, it employs dot-product operation to obtain a correlation/affinity matrix $Attn_{s \to s}(\mathcal{Q}^s, \mathcal{K}^s)$ of different support samples as

$$Attn_{s \to s}(\mathcal{Q}^s, \mathcal{K}^s) = Softmax(\frac{\mathcal{Q}^s(\mathcal{K}^s)^T}{\sqrt{c}}) \qquad (1)$$

where $c$ denotes the dimension of support features. It learns the representations for support samples by conducting the message passing operation as

$$\widehat{\mathcal{H}}^s = LN(\mathcal{H}^s + Attn_{s \to s}(\mathcal{Q}^s, \mathcal{K}^s)\mathcal{V}^s) \qquad (2)$$

where $LN(\cdot)$ refers to layer normalization. Besides, we add Feed-Forward Network (FFN) [8] and residual operation to obtain the final support sample representations as,

$$\widetilde{\mathcal{H}}^s = LN(\widehat{\mathcal{H}}^s + FFN(\widehat{\mathcal{H}}^s)) \qquad (3)$$

where $\widetilde{\mathcal{H}}^s = \{\widetilde{H}_1^s, \widetilde{H}_2^s \cdots, \widetilde{H}_{n_s}^s\} \in \mathbb{R}^{n_s \times c}$. $n_s$ denotes the number of support samples and $c$ is the feature dimension. FFN consists of two fully-connection layers.

**Decoder.** The Decoder aims to explore the dependence of samples in query set to learn the representations for query samples and also mines the intrinsic metrics of samples in query and support sets. To be specific, it takes the aforementioned encoded support features $\widetilde{\mathcal{H}}^s \in \mathbb{R}^{n_s \times c}$ and query feature embeddings $\widetilde{\mathcal{F}}^q \in \mathbb{R}^{n_q \times c \times h \times w}$ as its inputs. The *image tokenize* is applied on $\widetilde{\mathcal{F}}^q$ to obtain the initial query token sequence $\mathcal{H}^q = \{H_1^q, H_2^q, \cdots, H_{n_q}^q\} \in \mathbb{R}^{n_q \times c}$, where each token $H_j^q$ denotes a query image sample. Similar to the Encoder branch, we first leverage self-attention message passing mechanism to model the relationships among query samples and learn representations for query samples as

$$Attn_{q \to q}(\mathcal{Q}^q, \mathcal{K}^q) = Softmax(\frac{\mathcal{Q}^q(\mathcal{K}^q)^T}{\sqrt{c}}) \qquad (4)$$

$$\widehat{\mathcal{H}}^q = LN(\mathcal{H}^q + Attn_{q \to q}(\mathcal{Q}^q, \mathcal{K}^q)\mathcal{V}^q) \qquad (5)$$

where $LN(\cdot)$ denotes layer normalization.

Afterward, based on the support features $\widetilde{\mathcal{H}}^s$ and query features $\widehat{\mathcal{H}}^q$, we employ a **cross-attention** mechanism to explore the relationships between support and query samples for query sample representations. Specifically, it first computes the cross-affinities between support and query samples as follows

$$Attn_{q \to s}(\mathcal{Q}^q, \mathcal{K}^s) = Softmax(\mathcal{Q}^q(\mathcal{K}^s)^T) \qquad (6)$$

Then, it learns query sample representations by aggregating the information from support samples as follows

$$\widetilde{\mathcal{H}}^q = \widehat{\mathcal{H}}^q + LN(Attn_{q \to s}(\mathcal{Q}^q, \mathcal{K}^s)\mathcal{V}^s) \qquad (7)$$

where $\widetilde{\mathcal{H}}^q \in \mathbb{R}^{n_q \times c}$ and $LN(\cdot)$ denotes layer normalization. $\mathcal{Q}^q \in \mathbb{R}^{n_q \times c}$ is computed by conducting a linear projection on $\widehat{\mathcal{H}}^q$. $\mathcal{K}^s \in \mathbb{R}^{n_s \times c}$ and $\mathcal{V}^s \in \mathbb{R}^{n_s \times c}$ are obtained by conducting two different linear projections on $\widetilde{\mathcal{H}}^s$, respectively.

**Remark.** The above cross-affinities $Attn_{q \to s}(\mathcal{Q}^q, \mathcal{K}^s)$ naturally reflect the similarities/affinities between support and query samples. In our work, we regard them as global metric $m_g$ for all support and query samples, i.e.,

$$m_g(\mathcal{X}^s, \mathcal{X}^q) = Attn_{q \to s}(\mathcal{Q}^q, \mathcal{K}^s) \qquad (8)$$

where $m_g(\mathcal{X}^s, \mathcal{X}^q)$ contains the similarities for all query-support sample pairs in each episode. For convenience, in the following, we also use $m_g(X^s, X^q)$ to denote the metric between image $X^s$ and $X^q$, where $X^s \in \mathcal{X}^s, X^q \in \mathcal{X}^q$. We can utilize $m_g(X^s, X^q)$ for query sample classification, as discussed in the following Section Metric Learning and Few-shot Classification. Therefore, we can note that both query/support **sample representation** and **metric learning** in few-shot learning task are conducted simultaneously in our sampleFormer architecture. This is one main aspect of the proposed sampleFormer module.

### C. Patch Transformer Module

As a complementary to the above sampleFormer branch, we also develop a query-support Patch Transformer Module (patchFormer) to capture the more visual content of each image sample for local metric. As shown in Figure 2, patchFormer mainly consists of multi-head self-attention (MSA) and residual connection. Here, we omit Feed-Forward Network used in regular Transformer [8] for simplicity consideration. The parameters of MSA are shared on both support and query branches.

Concretely, for each input support sample $X^s$ and query sample $X^q$, we first obtain their feature embedding $\widetilde{F}^s \in \mathbb{R}^{c \times h \times w}$ and $\widetilde{F}^q \in \mathbb{R}^{c \times h \times w}$ by using the above CIFE, followed by the *patch tokenize* [8] to obtain the initial patch token sequence for each support and query image, i.e., $P^s = \{p_1^s, p_2^s, \cdots, p_{hw}^s\} \in \mathbb{R}^{hw \times c}$ and $P^q = \{p_1^q, p_2^q, \cdots, p_{hw}^q\} \in \mathbb{R}^{hw \times c}$. Then, we employ multi-head self-attention (MSA) [16] with shared weights and residual operation to transform the support and query image patch features as

$$\begin{aligned} \widetilde{P}^s &= LN(P^s + MSA(P^s)) \\ \widetilde{P}^q &= LN(P^q + MSA(P^q)) \end{aligned} \qquad (9)$$

where $LN(\cdot)$ denotes layer normalization.

Based on the above patch representations $\widetilde{P}^s = \{\widetilde{p}_1^s, \widetilde{p}_2^s, \cdots, \widetilde{p}_{hw}^s\}$ and $\widetilde{P}^q = \{\widetilde{p}_1^q, \widetilde{p}_2^q, \cdots, \widetilde{p}_{hw}^q\}$, we then adopt the Earth Mover's Distance (EMD) [3], [32] to compute their structural similarity. It first computes the distance between all patch pairs $(\widetilde{p}_i^s, \widetilde{p}_j^q)$ and then acquires the optimal matching between patches of two images that have the minimum distance cost. Finally, it returns the image-level metric

by aggregating the metrics of all matched patch pairs. In this paper, we denote this metric as local metric between support sample $X^s$ and query sample $X^q$, i.e.,

$$m_l(X^s, X^q) = EMD(\widetilde{P}^q, \widetilde{P}^s) \tag{10}$$

### D. Metric Learning and Few-Shot Classification

Given the support samples $(X^s, Y^s) \in \mathcal{X}^s$ with known labels and input query sample $X^q \in \mathcal{X}^q$, few-shot classification aims to determine the label of the query sample. To achieve this task, we first obtain the sample-based global metric $m_g(X^s, X^q)$ via Equ. (8) and patch-based local metric $m_l(X^s, X^q)$ via Equ. (10) respectively and combine them together to obtain the final metric/similarity between $X^s$ and $X^q$ as

$$m(X^s, X^q) = \lambda m_g(X^s, X^q) + (1 - \lambda)m_l(X^s, X^q) \tag{11}$$

where $\lambda \in (0, 1)$ is a tradeoff parameter.

Then, we can conduct few-shot classification by using the nearest neighbor classification strategy, i.e., the label of query $X^q$ is determined by the label $Y^{s^*}$ of the support sample $X^{s^*}$ that is most similar with query $X^q$, as used in previous works [3], [4].

**Loss Function.** In the training phase, we employ two loss functions for the proposed QSFormer. First, for the sample-Former module, we specifically introduce a contrastive loss as suggested in work [33], [34], which encourages the positive query-support sample pairs with same label (i.e., $Y^s = Y^q$) to be closing and the negative query-support sample pairs with different labels (i.e., $Y^s \neq Y^q$) are far away in each episode. This loss function can be written as follows,

$$L_{cl} = -log\frac{\sum\limits_{Y^s=Y^q} e^{m_g(X^s,X^q)}}{\sum\limits_{Y^s=Y^q} e^{m_g(X^s,X^q)} + \sum\limits_{Y^s \neq Y^q} e^{m_g(X^s,X^q)}} \tag{12}$$

where $m_g(X^s, X^q)$ is the global metric between query $X^q$ and support sample $X^s$. The whole network is trained in an end-to-end way by minimizing the Cross-Entropy (CE) loss function $\mathcal{L}_{ce}$ [3]. Thus, the total loss function can be formulated as

$$\mathcal{L}_{total} = \alpha\mathcal{L}_{ce}(\hat{Y}^q, Y^q) + (1 - \alpha)\mathcal{L}_{cl} \tag{13}$$

where $\hat{Y}^q$ is the label prediction obtained by our method and $Y^q$ denotes the corresponding ground-truth label. $\alpha \in (0, 1)$ is the balanced hyper-parameter.

**Implementation Details.** To achieve a fair comparison, the ResNet-12 [3], [6] with fully connected layers removed is adopted as the backbone module. It is firstly pre-trained from scratch and then use the episodic training based on meta-learning framework by following works [3], [7]. We empirically conduct the feature interaction of the last two levels in CIFE to obtain the enhanced sample features. We randomly sample 50/1000/5000 episodes from the training/validation/testing set on four public datasets. We compute the average accuracy and the corresponding 95% confidence interval to obtain the final performances of four datasets. Our proposed method is implemented by using Python on a server with a single 11G NVIDIA 2080Ti GPU. More hyper-parameter settings on four benchmarks for the proposed QSFormer are shown in Table VI.

## IV. EXPERIMENTS

### A. Datasets and Evaluation Metric

To verify our proposed QSFormer, we conduct extensive experiments on four publicly popular datasets for few-shot classification task, including **miniImageNet** [4], **tieredImageNet** [44], **Fewshot-CIFAR100** [36] and **Caltech-UCSD Birds-200-2011** [45]. We also conduct cross-domain experiments to evaluate the domain transfer ability of the proposed model. The recognition accuracy is adopted as the evaluation metric for our experiments. More details of datasets description are as follow.

**miniImageNet.** This dataset is a sub-dataset of ImageNet [46]. It contains a total of 100 classes with 600 samples in each class. As suggested in work [47], we divide these classes into training set, validation set and testing set, which respectively contains 64, 16 and 20 classes.

**tieredImageNet.** It contains 608 classes from 34 super-classes, with a total of 779,165 samples. Following [44], we split 34 super-classes into 20 super-classes (351 classes) for meta-training, 6 super-classes (97 classes) for meta-validation and 8 super-classes (160 classes) for meta-testing.

**FC100.** Fewshot-CIFAR100 is built upon the CIFAR100 dataset for few-shot classification task. It's named FC100 for short hereafter. It contains a total of 60,000 images from 100 classes. To reduce the information overlap, we group the 100 classes into 20 super-classes by following work [36]. Then, we divide these super-classes into training set, validation set and testing set, which contains 12, 4 and 4 super-classes respectively.

**CUB.** Caltech-UCSD Birds-200-2011 dataset is an extended vision of CUB-200 dataset. It's termed CUB for short hereafter. CUB is originally presented in fine-grained bird classification task. It contains the total of 11,788 images from 200 classes. As suggested by [27], we divide 200 classes into 100 classes for meta-training, 50 classes for meta-validation and 50 classes for meta-testing.

**miniImageNet → CUB.** By following [6], we train a model on miniImageNet dataset and evaluate on the CUB dataset to verify the transfer ability of model. In this experimental setting, specifically, we use all 100 classes of miniImageNet, with 600 samples per class for meta-training and use the meta-testing set (50 classes) of CUB dataset for meta-testing.

### B. Comparison with State-of-the-art Methods

As shown in Table I, we report our results and compare with other state-of-the-art (SOTA) approaches on miniImageNet [4] and tieredImageNet [44] datasets. From this Table, we can find that the proposed QSFormer beats many SOTA models on the miniImageNet dataset. For example, QSFormer exceeds the transformer-based HT [12] method by +11.14% and +11.46% in 1-shot and 5-shot tasks, respectively. For the attention mechanism based CAN [13], our model also outperforms it on the 1-shot/5-shot task by +1.39%/+0.52%. Compared with FETA [27] that is also developed based on ResNet12 and Transformer, the proposed QSFormer has better results.

From Table I, we can see that QSFormer achieves the best performance on the tieredImageNet dataset, i.e., 72.47±0.31

TABLE I
5-WAY RESULT COMPARISON OF OURS AND STATE-OF-THE-ART METHODS ON MINIIMAGENET AND TIEREDIMAGENET DATASETS. MOST RESULTS ARE FROM [3] OR THE ORIGINAL PAPERS. THE $1^{st}$, $2^{rd}$ AND $3^{rd}$ ARE RESPECTIVELY IN RED, BLUE AND GREEN. * DENOTES THIS METHOD IS REPRODUCED WITH OUR SETTINGS.

| Method | Backbone | miniImagenet | | tieredImagenet | |
|---|---|---|---|---|---|
| | | 1-shot | 5-shot | 1-shot | 5-shot |
| DHL [35] | Conv4 | $61.99 \pm -$ | $78.71 \pm -$ | $57.89 \pm -$ | $73.62 \pm -$ |
| cosine classifier [6]* | ResNet12 | $59.64 \pm 0.27$ | $75.80 \pm 0.21$ | $55.87 \pm 0.31$ | $80.92 \pm 0.23$ |
| TADAM [36] | ResNet12 | $58.50 \pm 0.30$ | $76.70 \pm 0.30$ | $-$ | $-$ |
| ECM [37] | ResNet12 | $59.00 \pm -$ | $77.46 \pm -$ | $63.99 \pm -$ | $81.97 \pm -$ |
| TPN [38] | ResNet12 | $59.46 \pm -$ | $75.65 \pm -$ | $59.91 \pm 0.94$ | $73.30 \pm 0.75$ |
| ProtoNet [5]* | ResNet12 | $63.03 \pm 0.29$ | $78.72 \pm 0.21$ | $68.68 \pm 0.34$ | $85.09 \pm 0.23$ |
| MTL [39] | ResNet12 | $61.20 \pm 1.80$ | $75.50 \pm 0.80$ | $-$ | $-$ |
| DC [40] | ResNet12 | $62.53 \pm 0.19$ | $79.77 \pm 0.19$ | $-$ | $-$ |
| MetaOptNet [41] | ResNet12 | $62.64 \pm 0.82$ | $78.63 \pm 0.46$ | $65.99 \pm 0.72$ | $81.56 \pm 0.53$ |
| MatchNet [4]* | ResNet12 | $61.24 \pm 0.29$ | $73.93 \pm 0.23$ | $71.01 \pm 0.33$ | $83.12 \pm 0.24$ |
| Meta-Baseline [7] | ResNet12 | $63.17 \pm 0.23$ | $79.26 \pm 0.17$ | $68.62 \pm 0.27$ | $83.74 \pm 0.18$ |
| CAN [13] | ResNet12 | $63.85 \pm 0.48$ | $79.44 \pm 0.34$ | $69.89 \pm 0.51$ | $84.23 \pm 0.37$ |
| PPA [42] | WRN-28-10 | $59.60 \pm 0.41$ | $73.74 \pm 0.19$ | $65.65 \pm 0.92$ | $83.40 \pm 0.65$ |
| wDAE-GNN [43] | WRN-28-10 | $61.07 \pm 0.15$ | $76.75 \pm 0.11$ | $68.18 \pm 0.16$ | $83.09 \pm 0.12$ |
| LEO [1] | WRN-28-10 | $61.76 \pm 0.08$ | $77.59 \pm 0.12$ | $66.33 \pm 0.05$ | $81.44 \pm 0.09$ |
| FEAT [27]* | ResNet12 | $64.75 \pm 0.28$ | $79.96 \pm 0.20$ | $71.34 \pm 0.33$ | $85.28 \pm 0.23$ |
| HT [12] | Transformer | $54.10 \pm -$ | $68.50 \pm -$ | $56.10 \pm -$ | $73.30 \pm -$ |
| DeepEMD [3]* | ResNet12 | $65.43 \pm 0.28$ | $79.28 \pm 0.20$ | $69.84 \pm 0.32$ | $84.06 \pm 0.23$ |
| DeepBDC [14]* | ResNet12 | $60.76 \pm 0.28$ | $78.25 \pm 0.20$ | $63.03 \pm 0.31$ | $81.57 \pm 0.22$ |
| QSFormer (Ours) | ResNet12 | $65.24 \pm 0.28$ | $79.96 \pm 0.20$ | $72.47 \pm 0.31$ | $85.43 \pm 0.22$ |

TABLE II
5-WAY RESULT COMPARISON OF OURS AND STATE-OF-THE-ART METHODS ON FEWSHOT-CIFAR100 DATASET. THE $1^{st}$, $2^{rd}$ AND $3^{rd}$ ARE RESPECTIVELY IN RED, BLUE AND GREEN. * DENOTES THIS METHOD IS REPRODUCED WITH OUR SETTINGS.

| Method | 1-shot | 5-shot |
|---|---|---|
| cosine classifier [6]* | $39.47 \pm 0.23$ | $56.29 \pm 0.25$ |
| FEAT [27]* | $42.28 \pm 0.26$ | $56.37 \pm 0.25$ |
| TADAM [36] | $40.10 \pm 0.40$ | $56.10 \pm 0.40$ |
| ProtoNet [5]* | $40.91 \pm 0.26$ | $56.66 \pm 0.25$ |
| MTL [39] | $45.10 \pm 1.8$ | $57.60 \pm 0.9$ |
| DC [40] | $42.04 \pm 0.17$ | $57.05 \pm 0.16$ |
| MetaOptNet [41] | $41.10 \pm 0.60$ | $55.50 \pm 0.60$ |
| MatchNet [4]* | $41.90 \pm 0.27$ | $54.41 \pm 0.25$ |
| TDE-FSL [48] | $44.61 \pm 0.96$ | $57.93 \pm 0.81$ |
| DeepEMD [3]* | $45.58 \pm 0.26$ | $62.08 \pm 0.25$ |
| DeepBDC [14]* | $43.57 \pm 0.25$ | $59.49 \pm 0.25$ |
| QSFormer (Ours) | $46.51 \pm 0.26$ | $61.58 \pm 0.25$ |

TABLE III
5-WAY RESULT COMPARISON OF OURS AND STATE-OF-THE-ART METHODS ON CALTECH-UCSD BIRDS-200-2011 DATASET. THE $1^{st}$, $2^{rd}$ AND $3^{rd}$ ARE RESPECTIVELY IN RED, BLUE AND GREEN. * DENOTES THIS METHOD IS REPRODUCED WITH OUR SETTINGS.

| Method | 1-shot | 5-shot |
|---|---|---|
| MELR [30] | $70.26 \pm 0.50$ | $85.01 \pm 0.32$ |
| IEPT [49] | $69.97 \pm 0.49$ | $84.33 \pm 0.33$ |
| MVT [50] | $-$ | $85.35 \pm 0.55$ |
| FEAT [27]* | $75.00 \pm 0.29$ | $86.24 \pm 0.19$ |
| cosine classifier [6]* | $62.09 \pm 0.29$ | $80.04 \pm 0.21$ |
| ProtoNet [5]* | $70.93 \pm 0.30$ | $85.55 \pm 0.19$ |
| MatchNet [4]* | $70.21 \pm 0.30$ | $82.69 \pm 0.22$ |
| RelationNet [15] | $66.20 \pm 0.99$ | $82.30 \pm 0.58$ |
| MAML [51] | $67.28 \pm 1.08$ | $83.47 \pm 0.59$ |
| DEML [52] | $66.95 \pm 1.06$ | $77.11 \pm 0.78$ |
| DeepEMD [3]* | $70.71 \pm 0.30$ | $86.13 \pm 0.19$ |
| DeepBDC [14]* | $65.45 \pm 0.29$ | $85.01 \pm 0.19$ |
| QSFormer (Ours) | $75.44 \pm 0.29$ | $86.30 \pm 0.19$ |

and 85.43±0.22 in 1-shot and 5-shot tasks. It exceeds the CAN [13] by +2.58 and +1.2 points in 1-shot and 5-shot tasks. Similar conclusions can also be drawn from the experimental results of Fewshot-CIFAR100 [36] and CUB [45] datasets, as illustrated in Table II and Table III. All in all, the proposed QS-Former attains SOTA performance on multiple FSL datasets, which fully demonstrates the effectiveness and advantages of our proposed QSFormer model.

## C. Ablation Study

To better understand the effectiveness of our proposed QSFormer, in this section, we conduct extensive ablation studies, including component analysis, similarity metric analysis, cross-domain analysis, etc.

**Component Analysis.** Our proposed QSFormer mainly contains three components: Cross-scale Interactive Feature Extractor (CIFE), Sample Transformer Module (sampleFormer) and Patch Transformer Module (patchFormer). The experimental results of ablation study are shown in Table V. We

TABLE IV
CROSS-DOMAIN EXPERIMENTS ($miniImagenet \rightarrow CUB$). * DENOTES THIS METHOD IS REPRODUCED WITH OUR SETTINGS. THE RED REPRESENTS THE BEST RESULTS AND BLUE DENOTES THE SECOND-BEST RESULTS.

| Methods | 1-shot | 5-shot |
|---|---|---|
| ProtoNet [5] | $50.01 \pm 0.82$ | $72.02 \pm 0.67$ |
| MatchNet [4] | $51.65 \pm 0.84$ | $69.14 \pm 0.72$ |
| cosine classifier [6] | $44.17 \pm 0.78$ | $69.01 \pm 0.74$ |
| Baseline [6] | $-$ | $65.57 \pm 0.70$ |
| Baseline++ [6] | $-$ | $62.04 \pm 0.76$ |
| FEAT [27]* | $52.67 \pm 0.29$ | $72.65 \pm 0.25$ |
| DeepEMD [3] | $54.24 \pm 0.86$ | $78.86 \pm 0.65$ |
| DeepBDC [14]* | $50.28 \pm 0.27$ | $76.49 \pm 0.23$ |
| QSFormer (Ours) | $55.04 \pm 0.29$ | $77.12 \pm 0.24$ |

reproduce cosine classifier method [6] consisting of CNN network and cosine distance as the Baseline network for comparison. From Table V, we can observe: (1) By comparing #1 with #2, the performance of Baseline network can be significantly improved with the help of CIFE, which demonstrates the

TABLE V
ABLATION STUDY FOR THE DIFFERENT COMPONENTS OF THE PROPOSED QSFORMER. THE BEST RESULTS ARE HIGHLIGHTED IN **BOLD**.

| # | Different Components | | | | Datasets | | | |
|---|---|---|---|---|---|---|---|---|
| | Baseline | CIFE | sampleFormer | patchFormer | miniImageNet | tieredImageNet | FC100 | CUB |
| 1 | ✓ | | | | 59.64 ± 0.27 | 55.87 ± 0.31 | 39.47 ± 0.23 | 62.09 ± 0.29 |
| 2 | ✓ | ✓ | | | 61.15 ± 0.28 | 70.73 ± 0.32 | 41.54 ± 0.25 | 65.95 ± 0.30 |
| 3 | ✓ | ✓ | ✓ | | 63.97 ± 0.28 | 71.64 ± 0.32 | 45.46 ± 0.26 | 72.93 ± 0.29 |
| 4 | ✓ | ✓ | ✓ | ✓ | **65.24 ± 0.28** | **72.47 ± 0.31** | **46.51 ± 0.26** | **75.44 ± 0.29** |

TABLE VI
HYPERPARAMETER SETTINGS OF OUR PROPOSED QSFORMER.

| Hyper-parameters | Datasets | | | | |
|---|---|---|---|---|---|
| | miniImageNet | tieredImageNet | FC100 | CUB | miniImageNet → CUB |
| Optimizer | SGD | SGD | SGD | SGD | SGD |
| Initial LR | 5e-4 | 5e-4 | 1e-4 | 5e-4 | 5e-4 |
| Steps of LR decay | 10 | 10 | 10 | 10 | 10 |
| Coefficient of LR decay | 0.9 | 0.5 | 0.9 | 0.95 | 0.9 |
| N | 3 | 3 | 4 | 2 | 3 |
| Number of Head | 10,8 | 8,8 | 8,1 | 8,1 | 10,8 |
| dropout rates | 0.5,0.5,0.5,0.1 | 0.5,0.5,0.5,0.1 | 0.5,0.5,0.5,0.1 | 0.1,0.5,0.5,0.1 | 0.5,0.5,0.5,0.1 |
| $\alpha$ | 0.7 | 0.5 | 0.5 | 0.05 | 0.7 |
| $\lambda$ | 0.1 | 0.1 | 0.4 | 0.3 | 0.1 |
| Epochs | 100 | 100 | 50 | 150 | 100 |

TABLE VII
PERFORMANCE COMPARISON OF THE CLASSICAL METHODS BASED ON DIFFERENT METRIC LEARNING. * DENOTES THE COMPARISON METHODS IS
REPRODUCED WITH OUR SETTING. THE **BOLD BLACK** REPRESENTS THE BEST RESULTS.

| Methods | Metric | miniImageNet | tieredImageNet | FC100 | CUB |
|---|---|---|---|---|---|
| cosine classifier [6]* | Cosine | 59.64 ± 0.27 | 55.87 ± 0.31 | 39.47 ± 0.23 | 62.09 ± 0.29 |
| MatchNet [4]* | Cosine | 61.24 ± 0.29 | 71.01 ± 0.33 | 41.90 ± 0.27 | 70.20 ± 0.30 |
| ProtoNet [5]* | Euclidean | 63.03 ± 0.29 | 68.68 ± 0.34 | 40.91 ± 0.26 | 70.93 ± 0.30 |
| DeepEMD [3]* | EMD | **65.43 ± 0.28** | 69.84 ± 0.32 | 45.58 ± 0.26 | 70.71 ± 0.30 |
| QSFormer | Ours | 65.24 ± 0.28 | **72.47 ± 0.31** | **46.51 ± 0.26** | **75.44 ± 0.29** |

effectiveness of CIFE. (2) By comparing #2 with #3, we can find that sampleFormer significantly improves the performance of model based on #2, which indicates the effectiveness of sampleFormer module. (3) By adding patchFormer into #3, we further improve the performance of whole network, which shows the effectiveness of patchFormer module. All these experiments fully validate the effectiveness of each component in our proposed QSFormer framework.

**Similarity Metric Analysis.** To verify the effectiveness of the proposed QSFormer on metric learning, we visualize the similarity distribution of Baseline and QSFormer on the more challenging 5-way 1-shot task, as shown in Figure 4. For 5-way 1-shot task, each query sample generates the similarity results of one positive query-support sample pair (i.e., "Q-S pos") and four negative query-support sample pairs (i.e., "Q-S neg") during the metric learning process. To facilitate the comparison of the similarity results of "Q-S pos" and "Q-S neg", we average the similarity values of four "Q-S neg" corresponding to each query sample. For this experiment, we perform 10 episodes, where each episode random selects $15 \times 5 = 75$ query samples for classification, i.e., we can get the $75 \times 10 = 750$ similarity values of "Q-S pos" and "Q-S neg", respectively. Subsequently, we count the number of "Q-S pos" and "Q-S neg" within a certain range according to the normalized similarity values and thus produce the similarity distribution as shown in Figure 4. We can observe that: (1) the similarity values of "Q-S pos" obtained by the

Baseline method are generally below 0.5, while "Q-S neg" are above 0.25. (2) In our proposed QSFormer, the similarity values of "Q-S pos" are mostly above 0.5, while "Q-S neg" are mostly below 0.25. Therefore, our proposed QSFormer can separate positive and negative query-support sample pairs more accurately.

In addition, we also compare our QSFormer with other metric learning algorithms, including cosine classifier [6], MatchNet [4], ProtoNet [5] and DeepEMD [3]. These compared methods are reproduced with the same settings and training schemes as ours for a more fair comparison. As shown in Table VII, we can observe that our proposed method obtains the best performance on four publicly popular datasets, which fully demonstrates the effectiveness and superiority of our proposed QSFormer. These experiments fully demonstrate the effectiveness of our proposed QSFormer for metric learning.

**Cross-domain Analysis.** To validate the transferable ability of our proposed QSFormer, we conduct a cross-domain experiment by following [3], [6]. The training and testing are implemented on miniImagenet dataset and CUB dataset, respectively. As shown in Table IV, our proposed QSFormer achieves the best performance on the 1-shot setting (55.04 ± 0.29) and the second-best results on the 5-shot, i.e., 77.12 ± 0.24. These results demonstrate that the proposed QSFormer learns the discriminative information across domains, and adaptively explores the correspondence of query-support samples.
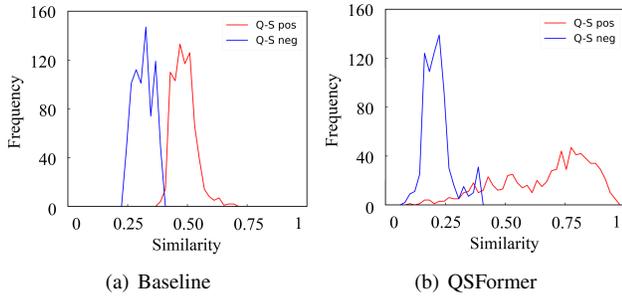
Fig. 4. Comparison of similarity distribution between Baseline and our QSFormer. The similarities of "Q-S pos" become larger while the similarities of "Q-S neg" become smaller, which indicates they are more easily separated.
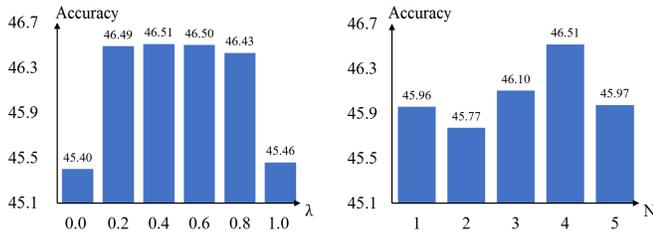


Fig. 5. Ablation study of two parameters (i.e., $\lambda$ and $N$).

**Parameter Analysis.** There are two important parameters in our model, including the balanced parameter $\lambda$ in Equ. (11) for local and global metric, and the number of sampleFormer layers $N$. In this section, we conduct experiments on the FC100 dataset on 5-way 1-shot task to check their influence. As shown in Figure 5, we can observe that the performance is relatively stable when we slightly adjust the balanced parameter $\lambda$ in the range of (0.2, 0.6). For the number $N$ of sampleFormer layers, we can find that our performance is increasing continuously when the $N$ is changing from 2 to 4. Therefore, we set $\lambda = 0.4$ and $N = 4$ for our experiments.

## V. CONCLUSION

In this paper, we propose a novel unified Query-Support Transformer (QSFormer) to deeply exploit the sample relationships in query and support sets for few-shot classification task. QSFormer mainly contains sample Transformer (sampleFormer) module and patch Transformer (patchFormer) module. sampleFormer is designed to meet the problem setting of few-shot classification, i.e., it couples the sample representation and metric learning between query and support sets together via a single Transformer architecture. Meanwhile, as a complementary, patchFormer is also adopted to model the local structural metric between query and support samples. A new CNN feature extractor (CIFE) is also proposed to provide an effective CNN backbone for our approach. Extensive experiments demonstrate the effectiveness and superiority of our proposed QSFormer approach.

## REFERENCES

[1] A. A. Rusu, D. Rao, J. Sygnowski, O. Vinyals, R. Pascanu, S. Osindero, and R. Hadsell, "Meta-learning with latent embedding optimization," in *Proceedings of the IEEE/CVF International Conference on Learning Representations*, 2018, pp. 6907–6917.

[2] M. A. Jamal and G.-J. Qi, "Task agnostic meta-learning for few-shot learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 719–11 727.

[3] C. Zhang, Y. Cai, G. Lin, and C. Shen, "Deepemd: Few-shot image classification with differentiable earth mover's distance and structured classifiers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 203–12 213.

[4] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, "Matching networks for one shot learning," *Advances in Neural Information Processing Systems*, vol. 29, 2016.

[5] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," *Advances in Neural Information Processing Systems*, vol. 30, pp. 4080–4090, 2017.

[6] W.-Y. Chen, Y.-C. Liu, Z. Kira, Y.-C. F. Wang, and J.-B. Huang, "A closer look at few-shot classification," in *Proceedings of the IEEE/CVF International Conference on Learning Representations*, 2019.

[7] Y. Chen, Z. Liu, H. Xu, T. Darrell, and X. Wang, "Meta-baseline: exploring simple meta-learning for few-shot learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9062–9071.

[8] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[9] Z. Peng, W. Huang, S. Gu, L. Xie, Y. Wang, J. Jiao, and Q. Ye, "Conformer: Local features coupling global representations for visual recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 367–376.

[10] Y. He, W. Liang, D. Zhao, H.-Y. Zhou, W. Ge, Y. Yu, and W. Zhang, "Attribute surrogates learning and spectral tokens pooling in transformers for few-shot learning," *arXiv preprint arXiv:2203.09064*, 2022.

[11] B. Dong, P. Zhou, S. Yan, and W. Zuo, "Self-promoted supervision for few-shot transformer," *arXiv preprint arXiv:2203.07057*, 2022.

[12] A. Zhmoginov, M. Sandler, and M. Vladymyrov, "Hypertransformer: Model generation for supervised and semi-supervised few-shot learning," *arXiv preprint arXiv:2201.04182*, 2022.

[13] R. Hou, H. Chang, B. Ma, S. Shan, and X. Chen, "Cross attention network for few-shot classification," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[14] J. Xie, F. Long, J. Lv, Q. Wang, and P. Li, "Joint distribution matters: Deep brownian distance covariance for few-shot classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 7972–7981.

[15] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1199–1208.

[16] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[17] N. Wang, W. Zhou, J. Wang, and H. Li, "Transformer meets tracker: Exploiting temporal context for robust visual tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 1571–1580.

[18] E. Yu, Z. Li, S. Han, and H. Wang, "Relationtrack: Relation-aware multiple object tracking with decoupled representation," *IEEE Transactions on Multimedia*, pp. 1–12, 2022.

[19] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Proceedings of European Conference on Computer Vision*, 2020, pp. 213–229.

[20] T. Guan, J. Wang, S. Lan, R. Chandra, Z. Wu, L. Davis, and D. Manocha, "M3detr: Multi-representation, multi-scale, mutual-relation 3d object detection with transformers," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 772–782.

[21] S. Liao and L. Shao, "Transmatcher: Deep image matching through transformers for generalizable person re-identification," *Advances in Neural Information Processing Systems*, vol. 34, pp. 1992–2003, 2021.

[22] M. Jia, X. Cheng, S. Lu, and J. Zhang, "Learning disentangled representation implicitly via transformer for occluded person re-identification," *IEEE Transactions on Multimedia*, pp. 1–11, 2022.

[23] S. Liu, L. Zhang, X. Yang, H. Su, and J. Zhu, "Query2label: A simple transformer way to multi-label classification," *arXiv preprint arXiv:2107.10834*, 2021.

[24] Z.-M. Chen, Q. Cui, B. Zhao, R. Song, X. Zhang, and O. Yoshie, "Sst: Spatial and semantic transformers for multi-label image recognition," *IEEE Transactions on Image Processing*, vol. 31, pp. 2570–2583, 2022.

[25] J. M. J. Valanarasu, P. Oza, I. Hacihaliloglu, and V. M. Patel, "Medical transformer: Gated axial-attention for medical image segmentation," in *Proceedings of International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2021, pp. 36–46.

[26] O. Petit, N. Thome, C. Rambour, L. Themyr, T. Collins, and L. Soler, "U-net transformer: Self and cross attention for medical image segmentation," in *Proceedings of International Workshop on Machine Learning in Medical Imaging*, 2021, pp. 267–276.

[27] H.-J. Ye, H. Hu, D.-C. Zhan, and F. Sha, "Few-shot learning via embedding adaptation with set-to-set functions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8808–8817.

[28] L. Liu, W. Hamilton, G. Long, J. Jiang, and H. Larochelle, "A universal representation transformer layer for few-shot image classification," in *Proceedings of the IEEE/CVF International Conference on Learning Representations*, 2021.

[29] B. Jiang, K. Zhao, and J. Tang, "Rgtransformer: Region-graph transformer for image representation and few-shot classification," *IEEE Signal Processing Letters*, vol. 29, pp. 792–796, 2022.

[30] N. Fei, Z. Lu, T. Xiang, and S. Huang, "Melr: Meta-learning via modeling episode-level relationships for few-shot learning," in *Proceedings of the IEEE/CVF International Conference on Learning Representations*, 2021.

[31] J. Wang, B. Song, D. Wang, and H. Qin, "Two-stream network with phase map for few-shot classification," *Neurocomputing*, vol. 472, pp. 45–53, 2022.

[32] F. L. Hitchcock, "The distribution of a product from several sources to numerous localities," *Journal of Mathematics and Physics*, vol. 20, no. 1-4, pp. 224–230, 1941.

[33] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[34] C. Liu, Y. Fu, C. Xu, S. Yang, J. Li, C. Wang, and L. Zhang, "Learning a few-shot embedding model with contrastive learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 10, 2021, pp. 8635–8643.

[35] X. Zhang, Y. Zhang, Z. Zhang, and J. Liu, "Discriminative learning of imaginary data for few-shot classification," *Neurocomputing*, vol. 467, pp. 406–417, 2022.

[36] B. Oreshkin, P. Rodríguez López, and A. Lacoste, "Tadam: Task dependent adaptive metric for improved few-shot learning," *Advances in Neural Information Processing Systems*, vol. 31, pp. 719–729, 2018.

[37] A. Ravichandran, R. Bhotika, and S. Soatto, "Few-shot learning with embedded class models and shot-free meta training," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 331–339.

[38] Y. Liu, J. Lee, M. Park, S. Kim, E. Yang, S. J. Hwang, and Y. Yang, "Learning to propagate labels: Transductive propagation network for few-shot learning," in *Proceedings of the IEEE/CVF International Conference on Learning Representations*, 2019.

[39] Q. Sun, Y. Liu, T.-S. Chua, and B. Schiele, "Meta-transfer learning for few-shot learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 403–412.

[40] Y. Lifchitz, Y. Avrithis, S. Picard, and A. Bursuc, "Dense classification and implanting for few-shot learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9258–9267.

[41] K. Lee, S. Maji, A. Ravichandran, and S. Soatto, "Meta-learning with differentiable convex optimization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 657–10 665.

[42] S. Qiao, C. Liu, W. Shen, and A. L. Yuille, "Few-shot image recognition by predicting parameters from activations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7229–7238.

[43] S. Gidaris and N. Komodakis, "Generating classification weights with gnn denoising autoencoders for few-shot learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 21–30.

[44] 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 *Proceedings of the IEEE/CVF International Conference on Learning Representations*, 2018.

[45] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, "The caltech-ucsd birds-200-2011 dataset," 2011.

[46] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

[47] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," 2017.

[48] L. Xing, S. Shao, W. Liu, A. Han, X. Pan, and B.-D. Liu, "Learning task-specific discriminative embeddings for few-shot image classification," *Neurocomputing*, vol. 488, pp. 1–13, 2022.

[49] M. Zhang, J. Zhang, Z. Lu, T. Xiang, M. Ding, and S. Huang, "Iept: Instance-level and episode-level pretext tasks for few-shot learning," in *Proceedings of the IEEE/CVF International Conference on Learning Representations*, 2021.

[50] S.-J. Park, S. Han, J.-W. Baek, I. Kim, J. Song, H. B. Lee, J.-J. Han, and S. J. Hwang, "Meta variance transfer: Learning to augment from the others," in *Proceedings of the IEEE/CVF International Conference on Machine Learning*, 2020, pp. 7510–7520.

[51] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proceedings of the IEEE/CVF International Conference on Machine Learning*, 2017, pp. 1126–1135.

[52] F. Zhou, B. Wu, and Z. Li, "Deep meta-learning: Learning to learn in the concept space," *arXiv preprint arXiv:1802.03596*, 2018.