

# tSF: Transformer-based Semantic Filter for Few-Shot Learning

Jinxiang Lai, Siqian Yang, Wenlong Liu, Yi Zeng, Zhongyi Huang, Wenlong Wu, Jun Liu, Bin-Bin Gao, and Chengjie Wang\*

Youtu Lab, Tencent

{jinxianglai, seasonsyang, sylviazeng, ezrealwu, jasoncjiang}@tencent.com, {nicehuster, huangzhny, junsenselee, csgaobb}@gmail.com

**Abstract.** Few-Shot Learning (FSL) alleviates the data shortage challenge via embedding discriminative target-aware features among plenty seen (base) and few unseen (novel) labeled samples. Most feature embedding modules in recent FSL methods are specially designed for corresponding learning tasks (e.g., classification, segmentation, and object detection), which limits the *utility* of embedding features. To this end, we propose a light and universal module named transformer-based Semantic Filter (tSF), which can be applied for different FSL tasks. The proposed tSF redesigns the inputs of a transformer-based structure by a semantic filter, which not only embeds the knowledge from whole base set to novel set but also filters semantic features for target category. Furthermore, the parameters of tSF is equal to half of a standard transformer block (less than  $1M$ ). In the experiments, our tSF is able to boost the performances in different classic few-shot learning tasks (about 2% improvement), especially outperforms the state-of-the-arts on multiple benchmark datasets in few-shot classification task.

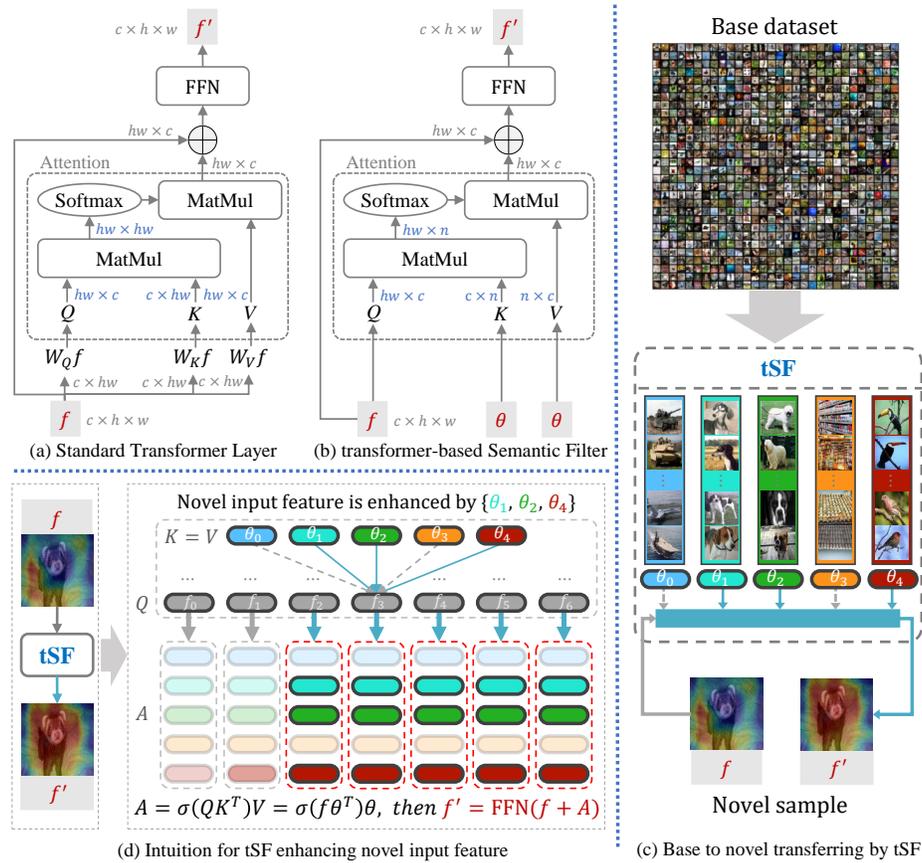
## 1 Introduction

Few-Shot Learning (FSL) aims to recognize unseen objects with plenty known data (base) and few labeled unknown samples (novel). Due to the shortage of novel data, FSL tasks suffer from weak representation problem. Hence, researchers [17,60,51,41,79,80,25,67] manage to design a embedding network to make extracted features robust and fine-grained enough in unseen instances recognition. To deal with different FSL tasks (e.g., classification, segmentation, and object detection), researchers propose different feature embedding modules, e.g., FEAT [71], CTX [6] for classification, HSNNet [31] for segmentation, and FSCE [48] for detection, respectively.

Nevertheless, these methods are limited by the purposes of different tasks. In the classification task, methods put emphasis on locating the representative prototype for each class. Methods in the detection task aim to distinguish the similar objects and correct the bounding box, while in the segmentation task, methods manage to generate a precise mask. In consequence, classification task requires more robust features, while detection task and segmentation task need more fine-grained features. To satisfy

---

\* Corresponding Author



**Fig. 1.** (a) Standard Transformer layer. (b) transformer-based Semantic Filter (tSF), where  $\theta$  is learnable semantic filter. (c) Base to novel transferring by tSF. After training, the semantic info of base dataset are embedded into  $\theta$ , e.g. there are  $n = 5$  semantic groups and  $\theta_1$  represents dog-like group. Then, given a novel input sample, tSF enhances its regions which are semantic similar to  $\theta$ . (d) Intuition for tSF enhancing novel input feature.

the demands of different tasks concurrently, a *target-aware and image-aware feature embedding method* is inevitable.

Throughout the recent investigations, transformer-base structure [54] brings an important significance at computer vision field, which works on almost all common tasks due to its sensibility on both big dataset (macro) and a single image (micro). Specifically, transformer is able to store the information of whole dataset modeling spatio-temporal correlations among instances. The property exactly satisfies the purpose of the feature embedding operation. Besides, observing from Transformer [54], Feat [72], SuperGlue [43], CTX [6] and DETR [4], we notice that the transformer layer could perform different learning behaviors with different input forms of  $\{Q, K, V\}$ . In common, a traditional transformer structure needs big training dataset to achieve high perfor-

mances. However, in few-shot learning field, it may fall into overfitting problem due to the shortage of data without a carefully designed framework.

To this end, we propose a light and general feature embedding module, named transformer based Semantic Filter (tSF), as illustrated in Fig. 1(b). We redesign the inputs  $Q, K, V$  of traditional transformer as  $f, \theta, \theta$ , where  $f$  is the extracted feature and  $\theta$  is a learnable weight, named semantic filter. The tSF uses the correlation matrix between  $(f, \theta)$  to re-weight  $\theta$ . The average of the re-weighted  $\theta$  is involved the dataset-attention response, which enhances the feature  $f$  from the views of both global dataset and local image. In this way, tSF can be trained without big data, while keeping the macro and micro information at the same time. Intuitively, Fig. 1(c) and Fig. 1(d) show that the proposed tSF is able to enhance the foreground regions of input novel feature via the embedded semantic info of base dataset. Besides, to further show the efficient of tSF, we insert the tSF into a strong few-shot classification framework, named PatchProto. The experimental results show that tSF helps PatchProto achieve SOTA performance with about 2% improvement. In addition, the parameter size of tSF is less than  $1M$ , half of that in traditional transformer. Moreover, to prove that tSF can suit different FSL frameworks, we conduct massive experiments on different few-shot learning tasks, and the results show that tSF can make 2% – 3% improvements on classification, detection and segmentation tasks.

In summary, our contributions are listed as follows:

- An effective and light module named transformer-based Semantic Filter (tSF) is proposed, which is helpful to learn a generalized-well embedding for novel targets (unseen in model training phase). The tSF leverages dataset-attention mechanism to realize information interaction between single input instance and whole dataset.
- A strong baseline framework called PatchProto is introduced for few-shot classification. Experimental results show our PatchProto+tSF approach outperforms the state-of-the-arts on multiple benchmark datasets such as miniImageNet and tieredImageNet.
- As a universal module, tSF is applied in different few-shot learning tasks, including classification, semantic segmentation and object detection. The massive experiments demonstrate that tSF is able to boost their performances.

## 2 Related Work

**Few-Shot Classification** FSL algorithms first pre-train a base classifier with abundant samples (seen images), then learn to recognize novel classes (unseen images) with a few labeled samples. According to recent investigations, there are four representative directions: optimization-based, parameter-generating based, metric-learning based, and embedding-based methods. *Optimization-based methods* are able to perform rapid adaption with a few training samples for new classes by learning a good optimizer [39,2] or learning a well-initialized model [12,35,42]. *Parameter-generating methods* [3,32,36,33,14] focus on learning a parameter generating network. *Metric-learning based methods* learn to compare to tackle the few-shot classification problem. The main idea is classifying a new input image by computing the similarity compared with labeled instances [15,55,17,67,46,50,63,73]. These methods design carefully to embedding network to match their corresponding distance metrics. *Embedding-based*

*methods* [60,51,41,79,80,25] firstly focus on learning a generalize-well embedding with supervised or self-supervised learning tasks, and then freeze this embedding and further train a linear classifier or design a metric classifier on novel classes.

**Auxiliary Task for Few-shot Classification** Some recent works gain a performance improvement by training few-shot models with supervised and self-supervised auxiliary tasks. The supervised task for FSL simply performs global classification on the base dataset as in CAN [17]. The effectiveness of self-supervised learning for FSL has been demonstrated, such as contrastive learning in either unsupervised pre-training [30] or episodic training [6,25], and auxiliary rotation prediction task [13,47,41].

**Few-Shot Semantic Segmentation** Early few-shot semantic segmentation methods apply a dual-branch architecture [44,7,38], one segmenting query-images with the prototypes learned by the other branch. In recently, the dual-branch architecture is unified into a single-branch, using the same embedding for support and query images [77,45,57,69,27]. These methods aim to leverage better guidance for the segmentation of query-images [77,34,56,74], via learning better class-specific representations [57,26,27,69,45] or iteratively refining [75].

**Few-Shot Object Detection** Existing few-shot object detection approaches can be divided into two paradigms: meta-learning based [20,64,10,18] and transfer learning based [59,62,48,11,37]. The majority of meta-learning approaches adopt *feature reweighting* or its variants to aggregate query and support features, which predict detections conditioned on support sets. Differently, the transfer learning based approaches firstly train the detectors on base set, then fine-tune the task-head layer on novel set, which achieve competitive results comparing to meta-learning approaches.

**Transformer** Transformer is an attention-based network architecture that is widely applied in natural language processing area [54,5]. Due to its power in learning representation, it has been introduced in many computer vision tasks, such as image classification [8,58,53], detection [76,4,82], segmentation [78,22,66], image matching [49,43] and few-shot learning [6,70,28].

To best of our knowledge, there is no general feature embedding method, which can be applied on multiple few-shot learning tasks.

### 3 Transformer-based Semantic Filter (tSF)

#### 3.1 Related Transformer

**Transformer-based Self-Attention** The architecture of a standard Transformer [54] layer is presented in Fig. 1(a), which consists of Attention and Feed-Forward Network (FFN). Given input feature  $f \in \mathbb{R}^{c \times h \times w}$ , the Transformer layer outputs  $f' \in \mathbb{R}^{c \times h \times w}$ . The key operation *Attention* is expressed as:

$$Attention(Q, K, V) = \sigma(QK^T)V, \quad (1)$$

where,  $\{Q, K, V\}$  are known as query, key and value respectively,  $\sigma$  is softmax function. The forms of  $\{Q, K, V\}$  in *transformer-based self-attention* are:

$$Q = W_Q f, \quad K = W_K f, \quad V = W_V f, \quad (2)$$

where,  $\{Q, K, V\} \in \mathbb{R}^{hw \times c}$  and some feature-reshaping operations are omitted for simplicity,  $W_Q, W_K, W_V$  are learnable weights (e.g. convolution layers). For few-shot classification, based on self-attention mechanism, Feat [72] used the standard transformer layer as a set-to-set function to perform embedding adaptation, formally:

$$Q = W_Q f_{set}, \quad K = W_K f_{set}, \quad V = W_V f_{set}, \quad (3)$$

where  $f_{set}$  is a set of features of all the instances in the support set. Both the standard transformer [54] and Feat [72] are based on self-attention mechanism which learns to model the relationship between the feature-nodes insight the input features. Differently, as illustrated in Fig. 1(b) and Eq. 6, our Transformer-based Semantic Filter (tSF) applies *Transformer-based Dataset-Attention*, which makes information interaction between single input sample and whole dataset.

**Transformer-based Cross-Attention** The standard Transformer performs self-attention behavior, and SuperGlue [43] found that Transformer can be used to make cross-attention between pair-features. Given input pair-features  $(f_{ref}, f)$ , it can obtain a cross-correlation matrix between  $(f_{ref}, f)$  which is used to re-weight  $f$  and achieves the cross-attention implementation. Specifically,  $\{Q, K, V\}$  forms in *transformer-based cross-attention* are:

$$Q = W_Q f_{ref}, \quad K = W_K f, \quad V = W_V f. \quad (4)$$

For few-shot classification, based on cross-attention mechanism, CTX [6] used the transformer layer to generate query-aligned prototype for support class.

**Transformer-based Decoder** The forms of  $\{Q, K, V\}$  in transformer-based decoder of DETR [4] are:

$$Q = \varphi, \quad K = V = f, \quad (5)$$

where  $\varphi \in \mathbb{R}^{u \times c}$  is learnable weights which is called as object queries. Given input feature  $f \in \mathbb{R}^{c \times h \times w}$ , the DETR [4] decoder outputs  $\varphi' \in \mathbb{R}^{u \times c}$  which is used to locate the objects. The dimension  $u$  of object queries  $\varphi$  represents the maximum number of objects insight a image.

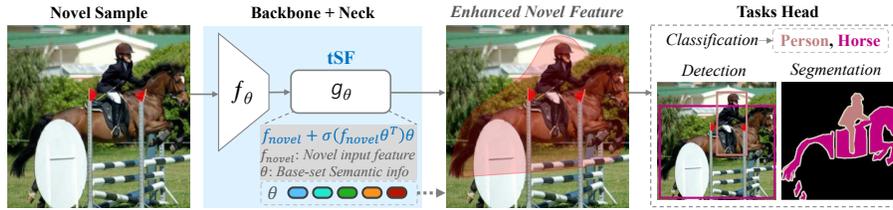
### 3.2 tSF Methodology

In few-shot learning task, obtaining a generalized-well embedding for novel categories is one of the key problem. To this end, we plan to model the whole dataset information and then transfer the knowledge from base set to novel set. Benefiting from the property of modeling spatio-temporal correlations among instances, transformer is able to learn the whole dataset information. Besides, observing from Transformer [54], Feat [72], SuperGlue [43], CTX [6] and DETR [4], we notice that the transformer layer performs different learning behavior with different input forms of  $\{Q, K, V\}$ . Therefore, in order to transfer the knowledge from whole base set to novel set, we propose a transformer-based Semantic Filter (tSF) as illustrated in Fig. 1(b), where the forms of  $\{Q, K, V\}$  are designed as:

$$Q = f, \quad K = V = \theta, \quad (6)$$

where,  $Q \in \mathbb{R}^{hw \times c}$  after feature-reshaping, and  $\theta \in \mathbb{R}^{n \times c}$  is learnable weights which is called as semantic filter. Formally, tSF is expressed as:

$$f' = tSF(f) = FFN(f + \sigma(f\theta^T)\theta). \quad (7)$$



**Fig. 2.** The tSF for few-shot learning tasks such as classification, semantic segmentation and object detection.

The input and output features of tSF are  $\{f, f'\} \in \mathbb{R}^{c \times h \times w}$  respectively, which is consistent with the standard Transformer. The tSF can be utilized as a feature-to-feature function, e.g. stacking tSF as a neck after the backbone architecture for few-shot learning as illustrated in Fig. 2 and Fig. 3(a).

Formally, let's define  $C_{base}$  and  $C_{novel}$  as categories of base and novel respectively. Although  $C_{base} \cap C_{novel} = \emptyset$ , we assume  $C_{base}$  consists of sub-sets  $C_{base}^{sim}$  and  $C_{base}^{diff}$  which are semantically similar and different from  $C_{novel}$  respectively. Then, tSF transfers knowledge from base to novel:

$$\begin{aligned} \text{Base Training} : f'_{base} &= FFN(f_{base} + \sigma(f_{base}\theta^T)\theta), \\ \text{Novel Testing} : f'_{novel} &= FFN(f_{novel} + \sigma(f_{novel}\theta^T)\theta). \end{aligned} \quad (8)$$

After training on base, semantic info of  $C_{base}$  are embedded into  $\theta$  of which dimension  $n$  are interpreted as projected semantic groups, i.e.  $\theta$  also consists of sub-sets  $\theta^{sim}$  and  $\theta^{diff}$  which are similar and different from  $C_{novel}$  respectively. Given a novel image, tSF enhances its regions which are similar to  $\theta^{sim}$  while  $\theta^{diff}$  doesn't. For example, if 'dog' in base, tSF enhances novel 'wolf' due to their high similarity relation. In Fig. 1(d),  $\theta^{sim}$  is  $\{\theta_1, \theta_2, \theta_4\}$ ,  $\theta^{diff}$  is  $\{\theta_0, \theta_3\}$ .

Intuitively, as illustrated in Fig. 1(c), with the model training on base set, the semantic filter  $\theta$  learns the whole dataset information. In the model testing on novel set, according to Eq. 7 and Fig. 1(d), the tSF uses the correlation matrix between  $(f, \theta)$  to re-weight  $\theta$ . The average of the re-weighted  $\theta$  is the dataset-attention response  $A$ , which is semantically similar to  $f$  and can be used to enhance the target object as  $f' = FFN(f + A)$ .  $f$  and  $\theta$  contains the info of one sample and whole dataset respectively, and the tSF can collect the target information (i.e.  $A$  which is semantically similar to  $f$ ) from  $\theta$  to enhance  $f$ . Therefore the foreground region of the input novel sample feature can be enhanced by the dataset-attention response.

### 3.3 Discussions

**Visualizations** Under the few-shot classification framework of PatchProto+tSF introduced in Sec. 4.2, we give the visualizations of feature response map and t-SNE for tSF as shown in Fig. 3(c) and Fig. 3(b) respectively. In detail,  $f' = g_\theta(f)$  and  $g_\theta$  is the proposed tSF, the correlation matrix  $R = \sigma(f_{novel}\theta^T) \in \mathbb{R}^{hw \times n}$  and  $R \cdot \theta_i \in \mathbb{R}^{hw}$  represents the correlation vector between  $f_{novel}$  and  $\theta_i \in \mathbb{R}^c$  ( $i^{th}$  position of  $\theta \in \mathbb{R}^{n \times c}$ ). In Fig. 3(c), comparing to  $f$ ,  $f'$  obtains more accurate and complete response map

focusing on the foreground region, which indicates that tSF is able to transfer semantic knowledge from base set to novel set. The visualizations of the correlation vector  $R.\theta_i$  show that  $\theta_i$  learns semantic information from base set, specifically, these targets foreground are enhanced mainly contributed from  $\{\theta_1, \theta_2, \theta_4\}$ . In addition, the t-SNE visualizations in Fig. 3(b) show that  $f'$  obtains more clear category boundaries than  $f$ , which demonstrates that tSF is able to produce more discriminative embedding features for novel categories.

**Properties** The properties of tSF are as follows: (i) Generalization ability: Based on dataset-attention mechanism, the tSF models the whole dataset information and then transfer the knowledge from base set to novel set. The tSF makes information interaction between input sample and whole dataset, while self-attention based transformer interacts info insight input sample itself which leads to overfitting problem due to insufficient information interaction. (ii) Representation ability: The low dimension semantic filter  $\theta \in \mathbb{R}^{n \times c}$  learns high-level semantic information from whole dataset. (iii) Efficiency: The computational complexity of tSF is less than self-attention based transformer. The complexity of transformer in calculating *Attention* by Eq. 1 is  $O(\text{transformer}) = (h \times w)^3 \times c$ , while our tSF is only  $O(\text{tSF}) = (h \times w)^2 \times n \times c$ , i.e.  $\frac{O(\text{transformer})}{O(\text{tSF})} = \frac{h \times w}{n} \gg 1$ .

**Comparisons** Comparing our tSF with Transformer and DETR decoder, in model testing on novel set, the input novel feature  $f_{\text{novel}}$  is enhanced by different information. Our tSF enhances  $f_{\text{novel}}$  by base dataset info (i.e. the semantic filter  $\theta$  learned on base set) as defined in Eq. 8, which fulfils base to novel transferring. Differently, the standard Transformer enhances  $f_{\text{novel}}$  by itself instance info:

$$\text{Transformer} : f'_{\text{novel}} = \text{FFN} \left( f_{\text{novel}} + \sigma \left( f_{\text{novel}} f_{\text{novel}}^T \right) f_{\text{novel}} \right). \quad (9)$$

Besides, the DETR decoder also enhances  $f_{\text{novel}}$  by itself instance info:

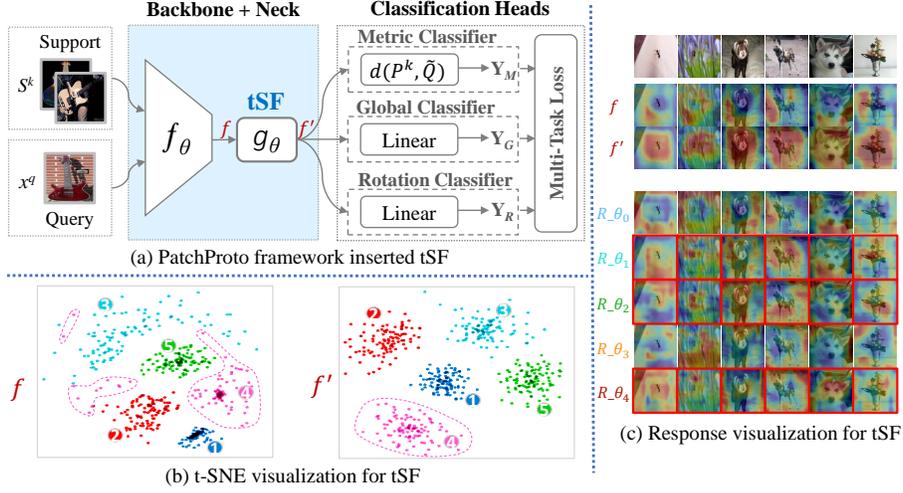
$$\text{DETR} : f'_{\text{novel}} = \text{FFN} \left( f_{\text{novel}} + \sigma \left( \theta f_{\text{novel}}^T \right) f_{\text{novel}} \right). \quad (10)$$

As illustrated in Tab. 2, the experimental comparisons show that our tSF obtains obvious performance gains, while Transformer and DETR show performance degradation due to overfitting problem.

## 4 tSF for Few-Shot Classification

### 4.1 Problem Definition

$N$ -way  $M$ -shot task to learn a classifier for  $N$  unseen classes with  $M$  labeled samples. Formally, we have three mutually disjoint datasets: a base set  $X_{\text{base}}$  for training, a validation set  $X_{\text{val}}$ , and a novel set  $X_{\text{novel}}$  for testing. Following [55,50,17,67], the episodic training strategy is adopted to mimic the few-shot learning setting, which has shown that it can effectively train a meta-learner (i.e., a few-shot classification model). Each episode contains  $N$  classes with  $M$  samples per class as the support set  $\mathcal{S} = \{(x_i^s, y_i^s)\}_{i=1}^{m_s}$  ( $m_s = N \times M$ ), and a fraction of the rest samples as the query set  $\mathcal{Q} = \{(x_i^q, y_i^q)\}_{i=1}^{m_q}$ . And the support subset of the  $k^{\text{th}}$  class is denoted as  $\mathcal{S}^k$ .



**Fig. 3.** (a) The PatchProto framework inserted tSF for few-shot classification. (b) The t-SNE visualization comparison for PatchProto+tSF, where  $f' = g_\theta(f)$  and  $g_\theta$  is the proposed tSF. (c) The visualizations of response map with the input of novel sample, where  $R_{\theta_i}$  is the correlation vector between  $(f, \theta_i)$ , and the dimension  $n$  of  $\theta \in \mathbb{R}^{n \times c}$  is set to 5.

## 4.2 PatchProto Framework with tSF

As illustrated in Fig. 3, the proposed PatchProto network consists of five components: feature extractor backbone  $f_\theta$ , transformer-based Semantic Filter (tSF)  $g_\theta$ , Metric classifier  $f_M$  for few-shot classification, and Global  $f_G$  and Rotation  $f_R$  classifiers for auxiliary tasks which are only used to assist model training.

The input image  $x^q$  in the query set  $\mathcal{Q} = \{(x_i^q, y_i^q)\}_{i=1}^{m_q}$  is rotated with  $[0^\circ, 90^\circ, 180^\circ, 270^\circ]$  and outputs a rotated  $\tilde{\mathcal{Q}} = \{(\tilde{x}_i^q, \tilde{y}_i^q)\}_{i=1}^{m_q \times 4}$ . Each support subset  $\mathcal{S}^k$  and a rotated query image  $\tilde{x}^q$  are fed through the feature extractor backbone  $f_\theta$  and tSF  $g_\theta$ , and produces the class feature  $P^k = \frac{1}{|\mathcal{S}^k|} \sum_{x_i^s \in \mathcal{S}^k} g_\theta(f_\theta(x_i^s))$  and query feature  $\tilde{Q} = g_\theta(f_\theta(\tilde{x}^q)) \in \mathbb{R}^{c \times h \times w}$ , respectively. Then the Metric classifier  $f_M$  makes classification via measuring the similarity between each pair-features  $(P^k, \tilde{Q})$ . Finally, PatchProto network is trained by optimizing the multi-task loss contributing from meta loss and auxiliary loss. In inductive inference phase, with the embedding learned in train set, the Metric classifier predicts the query  $x^q$  as  $Y_M$  based on cosine similarity measurement.

**Objective functions** *Meta loss:* As a metric-based meta learner, the Metric classifier predicts the query into  $N$  support categories by measuring cosine similarity. Following [17], we adopt the patch-wise classification mechanism to generate precise embeddings. Specifically, each local feature  $\tilde{Q}_m$  at  $m^{\text{th}}$  spatial position of  $\tilde{Q}$ , is predicted into  $N$  categories. Formally, the probability of recognizing  $\tilde{Q}_m$  as  $k^{\text{th}}$  category is:

$$\hat{Y}(y = k | \tilde{Q}_m) = \sigma \left( d \left( \tilde{Q}_m, \text{GAP} (P^k) \right) \right), \quad (11)$$

where *GAP* denotes global average pooling,  $d$  is cosine distance. Then the metric classification loss with the few-shot label  $\tilde{y}^q$  is:

$$\mathcal{L}_M = - \sum_{i=1}^{m_q} \sum_{m=1}^{h \times w} \log \hat{Y}(y = \tilde{y}_i^q | (\tilde{Q}_m)_i). \quad (12)$$

*Auxiliary loss:* The Global classifier predicts the query into all  $C$  categories of train set, thus its loss is:

$$\mathcal{L}_G = PCE(\tilde{Q}, C^q) = - \sum_{i=1}^{m_q} \sum_{n=1}^{h \times w} C_i^q \log \left( \sigma(W_l(\tilde{Q}_m)_i) \right). \quad (13)$$

where,  $W_l$  is a linear layer,  $C_i^q$  is the global category of  $\tilde{x}_i^q$  with all  $C$  categories, and *PCE* denotes the patch-wise cross-entropy function. Similarly, the loss of Rotation classifier is derived by  $\mathcal{L}_R = PCE(\tilde{Q}, B^q)$ , where  $B_i^q$  is the rotation category of  $\tilde{x}_i^q$  with four categories.

*Multi-task loss:* Therefore, inspired by [19], the overall classification loss is:

$$\mathcal{L} = \frac{1}{2} \mathcal{L}_M + \sum_{j=G,R} \left( (\lambda + w_j) \mathcal{L}_j + \log \frac{1}{(\lambda + w_j)} \right), \quad (14)$$

where,  $w = \frac{1}{2\alpha^2}$ ,  $\alpha$  is learnable variable,  $\lambda$  is a hyper-parameter to balance the effects of the losses of few-shot task and auxiliary tasks. The influence of  $\lambda$  is studied in Tab. 3.

## 5 tSF for Few-Shot Segmentation and Detection

As shown in Fig. 2, the proposed tSF is stacked after the backbone architecture (i.e. Backbone + tSF) for few-shot learning tasks, such as classification, segmentation and detection. To verify the effectiveness and the universality of our tSF module, we insert the tSF into the current state-of-the-art few-shot segmentation and detection methods. And the details are introduced in the next.

**RePRI+tSF for Segmentation** RePRI (Region Proportion Regularized Inference) [29] approach leverages the statistics of unlabeled pixels for the input image. It optimizes three complementary loss terms, including cross-entropy on labeled support pixels, Shannon entropy on unlabeled query pixels and a global KL-divergence regularizer on predicted foreground. RePRI achieves state-of-the-art on few-shot segmentation benchmark PASCAL-5<sup>i</sup> built from PASCAL VOC [9]. Based on the RePRI framework, we simply stack our tSF behind its backbone which obtains the RePRI+tSF approach.

**DeFRCN+tSF for Detection** DeFRCN (Decoupled Faster R-CNN) [37] is a simple yet effective fine-tuned approach for few-shot object detection, which proposes Gradient Decoupled Layer and Prototypical Calibration Block to alleviate the contradictions of Faster R-CNN in FSL scenario. Due to its simplicity and effective, DeFRCN achieves state-of-the-art on PASCAL VOC [9] and COCO [23]. To verify the effectiveness of tSF module in few shot object detection task, we use DeFRCN as baseline, and insert the tSF into the ResNet-101 backbone to obtain the DeFRCN+tSF approach, specifically, tSF module is followed in the 5th residual block. For the hyperparameter settings, we use the default parameter as same with DeFRCN.

**Table 1.** Comparison with existing methods on 5-way classification task on benchmark miniImageNet and tieredImageNet datasets.

model	Backbone	miniImageNet		tieredImageNet	
		1-shot	5-shot	1-shot	5-shot
MatchingNet [55]	Conv4	43.44 ± 0.77	60.60 ± 0.71	-	-
ProtoNet [46]	Conv4	49.42 ± 0.78	68.20 ± 0.66	53.31 ± 0.89	72.69 ± 0.74
RelationNet [50]	Conv4	50.44 ± 0.82	65.32 ± 0.70	54.48 ± 0.93	71.32 ± 0.78
<b>PatchProto</b>	Conv4	54.71 ± 0.46	70.67 ± 0.38	56.90 ± 0.51	71.47 ± 0.42
<b>PatchProto+tSF</b>	Conv4	<b>57.39 ± 0.47</b>	<b>73.34 ± 0.37</b>	<b>59.79 ± 0.51</b>	<b>74.56 ± 0.41</b>
CAN [17]	ResNet-12	63.85 ± 0.48	79.44 ± 0.34	69.89 ± 0.51	84.23 ± 0.37
MetaOpt+ArL [16]	ResNet-12	65.21 ± 0.58	80.41 ± 0.49	-	-
DeepEMD [73]	ResNet-12	65.91 ± 0.82	82.41 ± 0.56	71.16 ± 0.87	86.03 ± 0.58
IENet [41]	ResNet-12	66.82 ± 0.80	<b>84.35 ± 0.51</b>	71.87 ± 0.89	<b>86.82 ± 0.58</b>
DANet [67]	ResNet-12	67.76 ± 0.46	82.71 ± 0.31	71.89 ± 0.52	85.96 ± 0.35
<b>PatchProto</b>	ResNet-12	68.46 ± 0.47	82.65 ± 0.31	70.50 ± 0.50	83.60 ± 0.37
<b>PatchProto+tSF</b>	ResNet-12	<b>69.74 ± 0.47</b>	83.91 ± 0.30	<b>71.98 ± 0.50</b>	85.49 ± 0.35
wDAE-GNN [14]	WRN-28	61.07 ± 0.15	76.75 ± 0.11	68.18 ± 0.16	83.09 ± 0.12
LEO [42]	WRN-28	61.76 ± 0.08	77.59 ± 0.12	66.33 ± 0.05	81.44 ± 0.09
PSST [79]	WRN-28	64.16 ± 0.44	80.64 ± 0.32	-	-
FEAT [72]	WRN-28	65.10 ± 0.20	81.11 ± 0.14	70.41 ± 0.23	84.38 ± 0.16
CA [1]	WRN-28	65.92 ± 0.60	82.85 ± 0.55	74.40 ± 0.68	86.61 ± 0.59
<b>PatchProto</b>	WRN-28	69.34 ± 0.46	83.46 ± 0.30	73.40 ± 0.50	86.85 ± 0.35
<b>PatchProto+tSF</b>	WRN-28	<b>70.23 ± 0.46</b>	<b>84.55 ± 0.29</b>	<b>74.87 ± 0.49</b>	<b>88.05 ± 0.32</b>

## 6 Experiments

To prove the effectiveness and universality of the proposed tSF, massive experiments are conducted on different few-shot learning tasks, including classification, semantic segmentation and object detection. Overall, the results show that tSF can make 2%–3% improvements on these tasks.

### 6.1 Few-Shot Classification

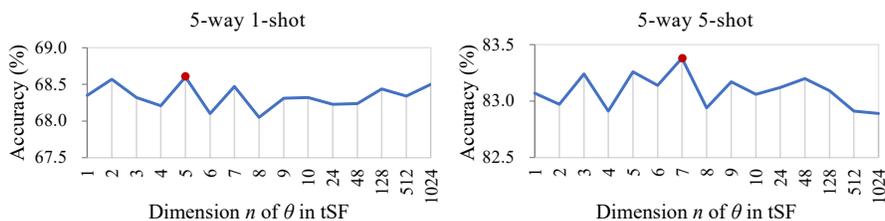
**Datasets** *miniImageNet* dataset is a subset of ImageNet [21], which consists of 100 classes. We split the 100 classes following the setting in [50,17,67], i.e. 64, 16 and 20 classes for training, validation and test respectively. *tieredImageNet* dataset [40] is also a subcollection of ImageNet [21]. It contains 608 classes, which are separated into 351 classes for training, 97 for validation and 160 for testing.

**Evaluation and Implementation details** We conduct experiments under 5-way 1-shot and 5-shot settings. We report the *average accuracy* and *95% confidence interval* over 2000 episodes sampled from the test set. Horizontal flipping, random cropping, random erasing [81] and color jittering are employed for data augmentation in training. According to the ablation study results in Tab. 3, the hyperparameter  $\lambda$  in Eq. 14 is set to 0.5 and 1.5 for ResNet-12 and WRN-28 respectively. The detailed info of optimizer, learning-rate and training-epochs are referred to our public source code.

**Comparison with State-of-the-arts** Tab. 1 compares our methods with existing few-shot classification algorithms on *miniImageNet* and *tieredImageNet*, which shows that

**Table 2.** The results on 5-way miniImageNet classification about the structure (refers to Fig.1 (a) and Fig.1 (b)) influence of tSF, which utilize PatchProto+tSF framework under ResNet-12 backbone. The dimension  $n$  of  $\theta$  in tSF is set to 5. The Metric and Global loss weights are set to 0.5 and 1.0 respectively, and the Rotation classifier is not applied.

Neck	$Q, K, V$	Attention Heads	Param	miniImageNet	
				1-shot	5-shot
None	-	-	7.75M	$67.47 \pm 0.47$	$81.85 \pm 0.32$
Transformer	$\bar{Q}=\bar{K}=\bar{V}=f$	1	8.75M	$63.25 \pm 0.45$	$79.44 \pm 0.33$
Transformer	$Q=W_Q f, K=W_K f, V=W_V f$	1	9.75M	$62.96 \pm 0.47$	$78.92 \pm 0.33$
Transformer	$Q=K=V=f$	4	8.75M	$62.68 \pm 0.47$	$78.98 \pm 0.34$
Transformer	$Q=W_Q f, K=W_K f, V=W_V f$	4	9.75M	$62.70 \pm 0.47$	$78.33 \pm 0.34$
DETR	$Q=\theta, K=V=f$	1	8.75M	$63.55 \pm 0.45$	$79.65 \pm 0.33$
tSF-V	$Q=\bar{K}=f, V=\theta$	1	9.00M	$61.84 \pm 0.48$	$76.11 \pm 0.36$
tSF-K	$Q=V=f, K=\theta$	1	9.00M	$64.73 \pm 0.45$	$80.43 \pm 0.33$
tSF	$Q=f, K=V=\theta$	1	8.75M	$68.37 \pm 0.46$	$83.08 \pm 0.31$
tSF	$Q=W_Q f, K=W_K \theta, V=W_V \theta$	1	9.75M	$68.27 \pm 0.47$	$83.01 \pm 0.31$
tSF	$Q=f, K=V=\theta$	4	8.75M	<b><math>68.60 \pm 0.47</math></b>	<b><math>83.26 \pm 0.31</math></b>
tSF	$Q=W_Q f, K=W_K \theta, V=W_V \theta$	4	9.75M	$68.49 \pm 0.47$	$82.95 \pm 0.31$
tSF	$Q=f, K=V=\theta$	8	8.75M	$68.46 \pm 0.46$	$83.14 \pm 0.31$
tSF	$Q=W_Q f, K=W_K \theta, V=W_V \theta$	8	9.75M	$68.42 \pm 0.47$	$83.12 \pm 0.31$



**Fig. 4.** The results on miniImageNet classification about the influence of dimension  $n$  of  $\theta \in \mathbb{R}^{n \times c}$  in tSF, which utilize PatchProto+tSF framework under ResNet-12 backbone without Rotation classifier.

the proposed PatchProto and PatchProto+tSF methods outperform the existing SOTAs under different backbones. And PatchProto+tSF shows obviously accuracy improvements under different backbones on 1-shot and 5-shot tasks than the strong baseline PatchProto, which demonstrates the effectiveness of the proposed tSF. The proposed PatchProto+tSF performs better than the optimization-based MetaOpt+ArL [16] and parameter-generating method wDAE-GCNN [14], with an improvement up to 4.53% and 9.16% respectively. Comparing to the competitive metric-based DeepEMD [73], PatchProto+tSF achieves 3.83% higher accuracy. Some metric-based methods [67,17] employing cross attention mechanism, and our PatchProto+tSF still surpasses the best DANet [67] with a performance improvement up to 1.98% on 1-shot. Overall our PatchProto+tSF obtains a new SOTA performance on both 1-shot and 5-shot classification tasks on miniImageNet and tieredImageNet, which demonstrates the strength of our framework and the effectiveness of the proposed tSF.

**Ablation Study** *Structure Influence of tSF:* As shown in Tab. 2, by comparing our tSF to the baseline without neck in first row, it shows consistent improvements on 1-

**Table 3.** The results on 5-way miniImageNet classification about the influence of multi-task loss employed in PatchProto+tSF under ResNet-12 and WRN-28 backbones. As introduced in Eq. 14,  $\lambda$  is the hyper-parameter, and  $w_G, w_R$  are learnable weights.

$\lambda$	Loss weights			ResNet-12		WRN-28	
	Metric	Global	Rotation	1-shot	5-shot	1-shot	5-shot
-	0.5	-	-	62.76 $\pm$ 0.49	80.07 $\pm$ 0.34	61.71 $\pm$ 0.50	77.53 $\pm$ 0.36
-	0.5	-	1.0	65.57 $\pm$ 0.49	80.81 $\pm$ 0.34	63.97 $\pm$ 0.50	79.25 $\pm$ 0.37
-	0.5	1.0	-	68.60 $\pm$ 0.47	82.98 $\pm$ 0.31	67.73 $\pm$ 0.47	82.59 $\pm$ 0.31
-	0.5	1.0	1.0	<b>69.41 <math>\pm</math> 0.46</b>	<b>83.82 <math>\pm</math> 0.30</b>	<b>69.64 <math>\pm</math> 0.47</b>	<b>84.01 <math>\pm</math> 0.30</b>
0.0	0.5	$\lambda+w_G$	$\lambda+w_R$	68.19 $\pm$ 0.47	83.00 $\pm$ 0.31	68.55 $\pm$ 0.48	83.30 $\pm$ 0.31
0.5	0.5	$\lambda+w_G$	$\lambda+w_R$	<b>69.74 <math>\pm</math> 0.47</b>	<b>83.91 <math>\pm</math> 0.30</b>	69.20 $\pm$ 0.46	84.03 $\pm$ 0.30
1.0	0.5	$\lambda+w_G$	$\lambda+w_R$	69.44 $\pm$ 0.46	83.90 $\pm$ 0.31	70.05 $\pm$ 0.46	84.17 $\pm$ 0.29
1.5	0.5	$\lambda+w_G$	$\lambda+w_R$	69.30 $\pm$ 0.45	83.82 $\pm$ 0.30	<b>70.23 <math>\pm</math> 0.46</b>	<b>84.55 <math>\pm</math> 0.29</b>
2.0	0.5	$\lambda+w_G$	$\lambda+w_R$	69.50 $\pm$ 0.45	83.86 $\pm$ 0.30	70.02 $\pm$ 0.45	83.61 $\pm$ 0.30

**Table 4.** The results on 1-way PASCAL-5<sup>i</sup> segmentation using mean-IoU. The dimension  $n$  of  $\theta$  in tSF is set to 5.

Method	Backbone	1 shot					5 shot				
		Fold-0	Fold-1	Fold-2	Fold-3	Mean	Fold-0	Fold-1	Fold-2	Fold-3	Mean
PANet [57]	VGG-16	42.3	58.0	51.1	41.2	48.1	51.8	64.6	59.8	46.5	55.7
RPMM [27]		47.1	<b>65.8</b>	50.6	<b>48.5</b>	53.0	50.0	66.5	51.9	47.6	54.0
RePRI [29]	VGG-16	49.7	63.4	58.2	42.8	53.5	54.5	67.2	63.7	48.8	58.6
<b>RePRI+tSF</b>		<b>53.0</b>	65.3	<b>58.3</b>	44.2	<b>55.2</b>	<b>57.0</b>	<b>67.9</b>	<b>63.9</b>	<b>50.8</b>	<b>59.9</b>
CANet [75]	ResNet-50	52.5	65.9	51.3	<b>51.9</b>	55.4	55.5	67.8	51.9	53.2	57.1
PGNet [74]		56.0	66.9	50.6	50.4	56.0	57.7	68.7	52.9	54.6	58.5
RPMM [69]		55.2	66.9	52.6	50.7	56.3	56.3	67.3	54.5	51.0	57.3
PPNet [27]		47.8	58.8	53.8	45.6	51.5	58.4	67.8	64.9	56.7	62.0
RePRI [29]		60.8	67.8	60.9	47.5	59.3	66.0	70.9	65.9	56.4	64.8
<b>RePRI+tSF</b>	<b>62.4</b>	<b>68.6</b>	<b>61.4</b>	49.4	<b>60.5</b>	<b>66.4</b>	<b>71.1</b>	<b>66.4</b>	<b>58.3</b>	<b>65.6</b>	

shot and 5-shot tasks, because tSF can transfer the knowledge from base set to novel set and generates more discriminative representations via focusing on the foreground regions. Comparing to the baseline without neck, the self-attention based transformer shows performance degradation due to overfitting on base dataset. Instead of behaving self-attention as standard transformer do, the proposed dataset-attention based tSF is able to prevent overfitting and generalize well on novel task, which is illustrated by the large accuracy improvement of tSF.

*Dimension Influence of  $\theta$  in tSF:* As shown in Fig.4, with a wide range [1, 1024] dimension  $n$  of  $\theta \in \mathbb{R}^{n \times c}$  in tSF, the accuracy differences are within 0.5% on 1-shot and 5-shot tasks, i.e. our PatchProto+tSF framework is not sensitive to the dimension of  $\theta$ . Considering the accuracy performance and computational complexity, we recommend to set the dimension  $n = 5$ . The  $n$  is interpreted as number of semantic groups. As  $n$  going larger, semantic groups become more fine-grained. The setting of  $n = 1$  represents one foreground group and is still able to obtain impressive performance.

*Influence of multi-task loss:* As illustrated in Tab. 3, the proposed PatchProto+tSF obtains its best results as setting  $\lambda$  to 0.5 and 1.5 under ResNet-12 and WRN-28 backbones

**Table 5.** The results on 1-way COCO-20<sup>i</sup> segmentation using mean-IoU. The dimension  $n$  of  $\theta$  in tSF is set to 5.

Method	Backbone	1 shot					5 shot				
		Fold-0	Fold-1	Fold-2	Fold-3	Mean	Fold-0	Fold-1	Fold-2	Fold-3	Mean
PPNet [27]		34.5	25.4	24.3	18.6	25.7	<b>48.3</b>	30.9	35.7	30.2	36.2
RPMM [69]	ResNet-50	29.5	36.8	29.0	27.0	30.6	33.8	42.0	33.0	33.3	35.5
PFENet [52]		36.5	38.6	34.5	33.8	35.8	36.5	43.3	37.8	38.4	39.0
<b>RePRI</b> [29]	ResNet-50	36.1	40.0	34.0	36.1	36.6	43.3	48.7	44.0	44.9	45.2
<b>RePRI+tSF</b>		<b>38.4</b>	<b>41.3</b>	<b>35.2</b>	<b>37.7</b>	<b>38.2</b>	45.0	<b>49.9</b>	<b>45.5</b>	<b>45.6</b>	<b>46.5</b>

**Table 6.** The results on VOC dataset. we evaluate the performance( $AP_{50}$ ) of DeFRCN under ResNet-101 with tSF module on three novel splits over multiple runs. The term  $w/G$  denotes whether using  $G$ -FSOD setting [59]. The dimension  $n$  of  $\theta$  in tSF is set to 5.

Method / Shots	$w/G$	Novel Set 1					Novel Set 2					Novel Set 3				
		1	2	3	5	10	1	2	3	5	10	1	2	3	5	10
MetaDet [61]	✗	18.9	20.6	30.2	36.8	49.6	21.8	23.1	27.8	31.7	43.0	20.6	23.9	29.4	43.9	44.1
TFA [59]	✗	39.8	36.1	44.7	55.7	56.0	23.5	26.9	34.1	35.1	39.1	30.8	34.8	42.8	49.5	49.8
<b>DeFRCN</b> [37]	✗	<b>53.6</b>	57.5	<b>61.5</b>	<b>64.1</b>	<b>60.8</b>	30.1	38.1	<b>47.0</b>	<b>53.3</b>	<b>47.9</b>	48.4	<b>50.9</b>	52.3	<b>54.9</b>	57.4
<b>DeFRCN+tSF</b>	✗	<b>53.6</b>	<b>58.1</b>	<b>61.5</b>	63.8	<b>60.8</b>	<b>31.5</b>	<b>39.3</b>	<b>47.0</b>	52.1	47.3	<b>48.5</b>	50.5	<b>52.8</b>	54.5	<b>58.0</b>
TFA [59]	✓	25.3	36.4	42.1	47.9	52.8	18.3	27.	30.9	34.1	39.5	17.9	27.2	34.3	40.8	45.6
FSDetView [65]	✓	24.2	35.3	42.2	49.1	57.4	21.6	24.6	31.9	37.0	45.7	21.2	30.0	37.2	43.8	49.6
<b>DeFRCN</b> [37]	✓	40.2	53.6	58.2	63.6	<b>66.5</b>	29.5	39.7	43.4	48.1	<b>52.8</b>	35.0	38.3	52.9	57.7	60.8
<b>DeFRCN+tSF</b>	✓	<b>43.6</b>	<b>57.4</b>	<b>61.2</b>	<b>65.1</b>	65.9	<b>31.0</b>	<b>40.3</b>	<b>45.3</b>	<b>49.6</b>	52.5	<b>39.3</b>	<b>51.4</b>	<b>54.8</b>	<b>59.8</b>	<b>62.1</b>

respectively. Comparing to the baseline (i.e. with Metric classification task only) in first row, our multi-task framework achieves large improvements on 1-shot and 5-shot tasks under different backbones. These results indicate that the auxiliary tasks (i.e. Global classification and Rotation classification) are useful for training a more robust embedding leading to an accuracy improvement, and the weights of the auxiliary tasks have a great influence on the overall few-shot classification performance.

## 6.2 Few-Shot Semantic Segmentation

**Datasets and Setting** *PASCAL-5<sup>i</sup> and COCO-20<sup>i</sup> Datasets*: PASCAL-5<sup>i</sup> is built from PASCAL VOC [9]. The 20 object categories are split into 4 folds. For each fold, 15 categories are utilized for training and the remaining 5 classes for testing. COCO-20<sup>i</sup> is built from MS-COCO [24]. COCO-20<sup>i</sup> dataset is divided into 4 folds with 60 base classes and 20 test classes in each fold.

*Evaluation Setting*: Following [27], the mean Intersection over Union (mIoU) is adopted for evaluation, and we report the average mIoU over 5 runs of 1000 tasks.

**Comparison with State-of-the-arts** Tab. 4 and Tab. 5 present our evaluation results on PASCAL-5<sup>i</sup> and COCO-20<sup>i</sup>. Comparing with existing few-shot semantic segmentation

methods, the RePRI+tSF approach achieves new state-of-the-art results. With the help of our tSF module, the RePRI+tSF obtains consistent performance improvement than RePRI, on 1-way 1-shot and 5-shot tasks under VGG-16 and ResNet-50 backbones.

### 6.3 Few-Shot Object Detection

**Datasets and Setting** *PASCAL VOC and COCO Datasets*: PASCAL VOC [9] are randomly sampled into 3 splits, and each contains 20 categories. For each split, there are 15 base and 5 novel categories. Each novel class has  $K = 1, 2, 3, 5, 10$  objects sampled from the train/val set of VOC2007 and VOC2012 for training, and the test set of VOC2007 for testing. COCO [23] use 60 categories disjoint with VOC as base set, and the remaining 20 categories are novel set with  $K = 1, 2, 3, 5, 10, 30$  shots. The total 5k images randomly sampled from the validation set are utilized for testing, while the rest for training.

*Evaluation Setting*: Following [59,65,20,68,37],

we conduct experiments on two evaluation protocols: few-shot object detection (*FSOD*) and generalized few-shot object detection (*G-FSOD*). The setting of FSOD only observes the performance of novel set. More comprehensively, that of G-FSOD considers both novel set and base set.

#### Comparison with State-of-the-arts

Tab. 6 and Tab. 7 present our evaluation results on VOC and COCO. Comparing with existing few-shot object detection methods, the DeFRCN+tSF approach achieves new state-of-the-art on VOC and COCO. On VOC three different data splits, under the FSOD and G-FSOD setting, the indicators of all

shots have been improved to a certain extent after adding tSF module. We achieve around 2.6*AP* improvement over the DeFRCN in all shots under G-FSOD setting. On COCO dataset, DeFRCN+tSF consistently outperforms DeFRCN in 1, 2, 3 and 5 shots.

**Table 7.** The results on COCO dataset. we report the performance (*mAP*) of DeFRCN under ResNet-101 with tSF module over multiple runs. The dimension  $n$  of  $\theta$  in tSF is set to 5.

Method / Shots	<i>w/G</i>	Shot Number					
		1	2	3	5	10	30
TFA [59]	✗	4.4	5.4	6.0	7.7	10.0	13.7
FSDetView [65]	✗	4.5	6.6	7.2	10.7	12.5	14.7
<b>DeFRCN [37]</b>	✗	9.3	12.9	<b>14.8</b>	16.1	<b>18.5</b>	<b>22.6</b>
<b>DeFRCN+tSF</b>	✗	<b>9.9</b>	<b>13.5</b>	<b>14.8</b>	<b>16.3</b>	18.3	22.5
TFA [59]	✓	1.9	3.9	5.1	7.0	9.1	12.1
FSDetView [65]	✓	3.2	4.9	6.7	8.1	10.7	15.9
<b>DeFRCN [37]</b>	✓	4.8	8.5	10.7	<b>13.6</b>	<b>16.8</b>	<b>21.2</b>
<b>DeFRCN+tSF</b>	✓	<b>5.0</b>	<b>8.7</b>	<b>10.9</b>	<b>13.6</b>	16.6	20.9

## 7 Conclusions

In this paper, we propose a transformer-based semantic filter (tSF) for few-shot learning problem. tSF leverages a well-designed transformer-based structure to encode the knowledge from base dataset and novel samples. In this way, a target-aware and image-aware feature can be generated via tSF. Moreover, tSF is a universal module, which can be applied into multiple few-shot learning tasks, e.g., classification, segmentation and detection. In addition, the parameter size of tSF is half of that of a standard transformer (less than  $1M$ ). The experimental results show that tSF is able to improve the performances about 2% in multiple classic few-shot learning tasks.

## References

1. Afrasiyabi, A., Lalonde, J.F., Gagné, C.: Associative alignment for few-shot image classification. In: ECCV (2020)
2. Andrychowicz, M., Denil, M., Gomez, S., Hoffman, M.W., Pfau, D., Schaul, T., Shillingford, B., De Freitas, N.: Learning to learn by gradient descent by gradient descent. In: NeurIPS (2016)
3. Bertinetto, L., Henriques, J.F., Valmadre, J., Torr, P., Vedaldi, A.: Learning feed-forward one-shot learners. In: NeurIPS (2016)
4. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: ECCV (2020)
5. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: NAACL-HLT (2020)
6. Doersch, C., Gupta, A., Zisserman, A.: Crosstransformers: spatially-aware few-shot transfer. In: NeurIPS (2020)
7. Dong, N., Xing, E.: Few-shot semantic segmentation with prototype learning. In: British Machine Vision Conference (BMVC) (2018)
8. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2021)
9. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International journal of computer vision* **88**(2), 303–338 (2010)
10. Fan, Q., Zhuo, W., Tang, C.K., Tai, Y.W.: Few-shot object detection with attention-rpn and multi-relation detector. In: CVPR (2020)
11. Fan, Z., Ma, Y., Li, Z., Sun, J.: Generalized few-shot object detection without forgetting. In: CVPR (2021)
12. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: ICML (2017)
13. Gidaris, S., Bursuc, A., Komodakis, N., Pérez, P., Cord, M.: Boosting few-shot visual learning with self-supervision. In: ICCV (2019)
14. Gidaris, S., Komodakis, N.: Generating classification weights with gnn denoising autoencoders for few-shot learning. In: CVPR (2019)
15. Gregory, K., Richard, Z., Ruslan, S.: Siamese neural networks for one-shot image recognition. In: ICML Workshops (2015)
16. Hongguang, Z., Piotr, K., Songlei, J., Hongdong, L., Philip, H. S., T.: Rethinking class relations: Absolute-relative supervised and unsupervised few-shot learning. In: CVPR (2021)
17. Hou, R., Chang, H., Bingpeng, M., Shan, S., Chen, X.: Cross attention network for few-shot classification. In: NeurIPS (2019)
18. Hu, H., Bai, S., Li, A., Cui, J., Wang, L.: Dense relation distillation with context-aware aggregation for few-shot object detection. In: CVPR (2021)
19. Jinxiang, L., Siqian, Y., Guannan, J., Xi, W., Yuxi, L., Zihui, J., Xiaochen, C., Jun, L., Binbin, G., Wei, Z., Yuan, X., Chengjie, W.: Rethinking the metric in few-shot learning: From an adaptive multi-distance perspective. In: ACMMM (2022)
20. Kang, B., Liu, Z., Wang, X., Yu, F., Feng, J., Darrell, T.: Few-shot object detection via feature reweighting. In: ICCV (2019)
21. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: NeurIPS (2012)
22. Liang, J., Homayounfar, N., Ma, W.C., Xiong, Y., Hu, R., Urtasun, R.: Polytransform: Deep polygon transformer for instance segmentation. In: CVPR (2020)

23. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014)
24. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014)
25. Liu, C., Fu, Y., Xu, C., Yang, S., Li, J., Wang, C., Zhang, L.: Learning a few-shot embedding model with contrastive learning. In: AAAI (2021)
26. Liu, W., Zhang, C., Lin, G., Liu, F.: CRNet: Cross-reference networks for few-shot segmentation. In: CVPR (2020)
27. Liu, Y., Zhang, X., Zhang, S., He, X.: Part-aware prototype network for few-shot semantic segmentation. In: ECCV (2020)
28. Lu, Z., He, S., Zhu, X., Zhang, L., Song, Y.Z., Xiang, T.: Simpler is better: Few-shot semantic segmentation with classifier weight transformer. In: ICCV (2021)
29. Malik, B., Hoel, K., Ziko, I.M., Pablo, P., Ismail, B.A., Jose, D.: Few-shot segmentation without meta-learning: A good transductive inference is all you need? In: CVPR (2021)
30. Medina, C., Devos, A., Grossglauser, M.: Self-supervised prototypical transfer learning for few-shot classification. arXiv:2006.11325 (2020)
31. Min, J., Kang, D., Cho, M.: Hypercorrelation squeeze for few-shot segmentation. In: ICCV (2021)
32. Munkhdalai, T., Yu, H.: Meta networks. In: ICML (2017)
33. Munkhdalai, T., Yuan, X., Mehri, S., Trischler, A.: Rapid adaptation with conditionally shifted neurons. In: ICML (2018)
34. Nguyen, K., Todorovic, S.: Feature weighting and boosting for few-shot segmentation. In: ICCV (2019)
35. Nichol, A., Achiam, J., Schulman, J.: On first-order meta-learning algorithms. arXiv:1803.02999 (2018)
36. Qi, C., Yingwei, P., Ting, Y., Chenggang, Y., Tao, M.: Memory matching networks for one-shot image recognition. In: CVPR (2018)
37. Qiao, L., Zhao, Y., Li, Z., Qiu, X., Wu, J., Zhang, C.: Defrcn: Decoupled faster r-cnn for few-shot object detection. In: ICCV (2021)
38. Rakelly, K., Shelhamer, E., Darrell, T., Efros, A., Levine, S.: Conditional networks for few-shot semantic segmentation. In: ICLR Workshop (2018)
39. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: ICLR (2017)
40. Ren, M., Triantafillou, E., Ravi, S., Snell, J., Swersky, K., Tenenbaum, J.B., Larochelle, H., Zemel, R.S.: Meta-learning for semi-supervised few-shot classification. In: ICLR (2018)
41. Rizve, M.N., Khan, S., Khan, F.S., Shah, M.: Exploring complementary strengths of invariant and equivariant representations for few-shot learning. In: CVPR (2021)
42. Rusu, A.A., Rao, D., Sygnowski, J., Vinyals, O., Pascanu, R., Osindero, S., Hadsell, R.: Meta-learning with latent embedding optimization. In: ICLR (2019)
43. Sarlin, P.E., DeTone, D., Malisiewicz, T., Rabinovich, A.: Superglue: Learning feature matching with graph neural networks. In: CVPR (2020)
44. Shaban, A., Bansal, S., Liu, Z., Essa, I., Boots, B.: One-shot learning for semantic segmentation. In: BMVC (2018)
45. Siam, M., Oreshkin, B.N., Jagersand, M.: AMP: Adaptive masked proxies for few-shot segmentation. In: ICCV (2019)
46. Snell, J., Swersky, K., Zemel, R.: Prototypical networks for few-shot learning. In: NeurIPS (2017)
47. Su, J.C., Maji, S., Hariharan, B.: When does self-supervision improve few-shot learning? In: ECCV (2020)
48. Sun, B., Li, B., Cai, S., Yuan, Y., Zhang, C.: Fscf: Few-shot object detection via contrastive proposal encoding. In: CVPR (2021)

49. Sun, J., Shen, Z., Wang, Y., Bao, H., Zhou, X.: Loftr: Detector-free local feature matching with transformers. In: CVPR (2021)
50. Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P.H., Hospedales, T.M.: Learning to compare: Relation network for few-shot learning. In: CVPR (2018)
51. Tian, Y., Wang, Y., Krishnan, D., Tenenbaum, J.B., Isola, P.: Rethinking few-shot image classification: a good embedding is all you need? In: ECCV (2020)
52. Tian, Z., Zhao, H., Shu, M., Yang, Z., Li, R., Jia, J.: Prior guided feature enrichment network for few-shot segmentation. TPAMI (2020)
53. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: ICML (2021)
54. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: NeurIPS (2017)
55. Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al.: Matching networks for one shot learning. In: NeurIPS (2016)
56. Wang, H., Zhang, X., Hu, Y., Yang, Y., Cao, X., Zhen, X.: Few-shot semantic segmentation with democratic attention networks. In: ECCV (2020)
57. Wang, K., Liew, J.H., Zou, Y., Zhou, D., Feng, J.: Panet: Few-shot image semantic segmentation with prototype alignment. In: ICCV (2019)
58. Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In: ICCV (2021)
59. Wang, X., Huang, T.E., Darrell, T., Gonzalez, J.E., Yu, F.: Frustratingly simple few-shot object detection. arXiv:2003.06957 (2020)
60. Wang, Y., Chao, W.L., Weinberger, K.Q., van der Maaten, L.: Simpleshot: Revisiting nearest-neighbor classification for few-shot learning. arXiv:1911.04623 (2019)
61. Wang, Y.X., Ramanan, D., Hebert, M.: Meta-learning to detect rare objects. In: ICCV (2019)
62. Wu, J., Liu, S., Huang, D., Wang, Y.: Multi-scale positive sample refinement for few-shot object detection. In: ECCV (2020)
63. Wu, Z., Li, Y., Guo, L., Jia, K.: Parn: Position-aware relation networks for few-shot learning. In: ICCV (2019)
64. Xiao, Y., Marlet, R.: Few-shot object detection and viewpoint estimation for objects in the wild. In: ECCV (2020)
65. Xiao, Y., Marlet, R.: Few-shot object detection and viewpoint estimation for objects in the wild. In: ECCV (2020)
66. Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., Luo, P.: Segformer: Simple and efficient design for semantic segmentation with transformers. In: NeurIPS (2021)
67. Xu, C., Fu, Y., Liu, C., Wang, C., Li, J., Huang, F., Zhang, L., Xue, X.: Learning dynamic alignment via meta-filter for few-shot learning. In: CVPR (2021)
68. Yan, X., Chen, Z., Xu, A., Wang, X., Liang, X., Lin, L.: Meta r-cnn: Towards general solver for instance-level low-shot learning. In: ICCV (2019)
69. Yang, B., Liu, C., Li, B., Jiao, J., Ye, Q.: Prototype mixture models for few-shot semantic segmentation. In: ECCV (2020)
70. Yang, Z., Wang, Y., Chen, X., Liu, J., Qiao, Y.: Context-transformer: tackling object confusion for few-shot detection. In: AAAI (2020)
71. Ye, H.J., Hu, H., Zhan, D.C., Sha, F.: Few-shot learning via embedding adaptation with set-to-set functions. In: CVPR (2020)
72. Ye, H.J., Hu, H., Zhan, D.C., Sha, F.: Few-shot learning via embedding adaptation with set-to-set functions. In: CVPR (2020)
73. Zhang, C., Cai, Y., Lin, G., Shen, C.: Deepemd: Few-shot image classification with differentiable earth mover's distance and structured classifiers. In: CVPR (2020)

74. Zhang, C., Lin, G., Liu, F., Guo, J., Wu, Q., Yao, R.: Pyramid graph networks with connection attentions for region-based one-shot semantic segmentation. In: ICCV (2019)
75. Zhang, C., Lin, G., Liu, F., Yao, R., Shen, C.: CANet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In: CVPR (2019)
76. Zhang, D., Zhang, H., Tang, J., Wang, M., Hua, X., Sun, Q.: Feature pyramid transformer. In: ECCV (2020)
77. Zhang, X., Wei, Y., Yang, Y., Huang, T.S.: SG-one: Similarity guidance network for one-shot semantic segmentation. IEEE Transactions on Cybernetics (2020)
78. Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., Fu, Y., Feng, J., Xiang, T., Torr, P.H., et al.: Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In: CVPR (2021)
79. Zhengyu, C., Jixie, G., Heshen, Z., Siteng, H., Donglin, W.: Pareto self-supervised training for few-shot learning. In: CVPR (2021)
80. Zhiqiang, S., Zechun, L., Jie, Q., Marios, S., Kwang-Ting, C.: Partial is better than all:revisiting fine-tuning strategy for few-shot learning. In: AAAI (2021)
81. Zhong, Z., Zheng, L., Kang, G., Li, S., Yang, Y.: Random erasing data augmentation. In: AAAI (2020)
82. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable detr: Deformable transformers for end-to-end object detection. In: ICLR (2020)