# A Task-aware Dual Similarity Network for Fine-grained Few-shot Learning

Yan Qi[1], Han Sun[1](✉), Ningzhong Liu[1], and Huiyu Zhou[2]

[1] Nanjing University of Aeronautics and Astronautics, Jiangsu Nanjing, China
[2] School of Computing and Mathematical Sciences, University of Leicester, U.K
`sunhan@nuaa.edu.cn`

**Abstract.** The goal of fine-grained few-shot learning is to recognize sub-categories under the same super-category by learning few labeled samples. Most of the recent approaches adopt a single similarity measure, that is, global or local measure alone. However, for fine-grained images with high intra-class variance and low inter-class variance, exploring global invariant features and discriminative local details is quite essential. In this paper, we propose a Task-aware Dual Similarity Network(TDSNet), which applies global features and local patches to achieve better performance. Specifically, a local feature enhancement module is adopted to activate the features with strong discriminability. Besides, task-aware attention exploits the important patches among the entire task. Finally, both the class prototypes obtained by global features and discriminative local patches are employed for prediction. Extensive experiments on three fine-grained datasets demonstrate that the proposed TDSNet achieves competitive performance by comparing with other state-of-the-art algorithms.

**Keywords:** Fine-grained image classification · Few-shot learning · Feature enhancement

## 1 Introduction

As one of the most important problems in the field of artificial intelligence, fine-grained image classification [6, 8] aims to identify objects of sub-categories under the same super-category. Different from the traditional image classification task [21, 22], the images of sub-categories are similar to each other, which makes fine-grained recognition still a popular and challenging topic in computer vision.

Benefiting from the development of Convolution Neural Networks (CNNs), fine-grained image classification has made significant progress. Most approaches typically rely on supervision from a large number of labeled samples. In contrast, humans can identify new classes with only few labeled examples. Recently, some studies [25, 31] focus on a more challenging setting, which aims to recognize fine-grained images from few samples, and is called fine-grained few-shot learning(FG-FSL). Learning from fine-grained images with few samples brings two challenges. On the one hand, images in the same category are quite different due to poses, illumination conditions, backgrounds, etc. So how to capture

invariant features in limited samples is a particularly critical problem. On the other hand, it is complicated to distinguish subtle visual appearance clues on account of the small differences between categories. Therefore, we consider that the invariant global structure and the discriminative local details of objects are both crucial for fine-grained few-shot classification.

To effectively learn latent patterns from few labeled images, many approaches [7, 33] have been proposed in recent years. These methods can be roughly divided into two branches: the meta-learning methods and the metric learning ones. Metric learning has attracted more and more attention due to its simplicity and effectiveness, and our work will focus on such methods. Traditional approaches such as matching network [29] and relation network [27] usually utilize global features for recognition. However, the distribution of these image-level global features cannot be accurately estimated because of the sparseness of the samples. In addition, discriminative clues may not be detected only by relying on global features. CovaMNet [15] and DN4 [16] introduce the deep local descriptors which are exploited to describe the distribution with each class feature. Furthermore, although these methods learn abundant features, they deal with each support class independently and cannot employ the contextual information of the whole task to generate task-specific features. In conclusion, the importance of different parts changes with different tasks.

In this paper, we propose a Task-aware Dual Similarity Network(TDSNet) for fine-grained few-shot learning, which makes full use of both global invariant features and discriminative local details of images. More specifically, first, a local feature enhancement module is employed to activate discriminative semantic parts by matching the predicted distribution between objects and parts. Second, in the dual similarity module, the proposed TDSNet calculates the class prototypes as global invariant features. Especially, in the local similarity branch, task-aware attention is adopted to select important image patches for the current task. By considering the context of the entire support set as a whole, the key patches in the task are selected and weighted without paying too much attention to the unimportant parts. Finally, both global and local similarities are employed for the final classification. We conduct comprehensive experiments on three popular fine-grained datasets to demonstrate the effectiveness of our proposed method. Especially, our method can also have good performance when there is only one training image.

## 2   Related Work

**Few shot learning.** Few-shot learning aims at recognizing unseen classes with only few samples. The recently popular literature on few-shot learning can be roughly divided into the following two categories: meta-learning based methods and metric-learning based methods.

Meta-learning based methods attempt to learn a good optimizer to update model parameters. MAML [9] is dedicated to learning a good parameter initialization so that the model can adapt to the new task after training on few

samples. Ravi et al. [24] propose a meta-learner optimizer based on LSTM to optimize a classifier while also studying an initialization for the learner that contains task-aware knowledge.

Metric-learning based methods aim to measure the similarity by learning an appropriate metric that quantifies the relationship between the query images and support sets. Koch et al. [13] adopt a siamese convolutional neural network to learn generic image representations, which is performed as a binary classification network. Lifchitz et al. [18] directly predict classification for each local representation and calculates the loss. DN4 [16] employs k-nearest neighbors to construct an image-to-class search space that utilizes deep local representations. Unlike DN4, which is most relevant to our work, we argue that considering each support class independently may capture features shared among classes that are unimportant for classification. In this paper, task-aware local representations will be detected to explore richer information.

**Fine-grained image classification.** Because some early approaches [1,3] require a lot of bounding boxes or part annotations as supervision that needs a high cost of expert knowledge, more and more researchers are turning their attention to weakly supervised methods [20,23] that rely only on image-level annotations. Inspired by different convolutional feature channels corresponding to different types of visual modes, MC-Loss [5] proposes a mutual-channel loss that consists of a discriminality component and a diversity component to get the channels with locally discriminative regions for a specific class. TDSA-Loss [4] obtains multi-regional and multi-granularity features by constraining mid-level features with the attention generated by high-level features. Different from these methods, we consider that the discriminability of local features obtained only by the attention maps may not be guaranteed. In order to overcome this limitation, the proposed TDSNet activates the local representations with strong discriminability by matching the distribution between the global features and their sub-features, so that the discriminability of global features at fine-grained scales is improved.

## 3   Method

### 3.1   Problem Definition

In this paper, the proposed TDSNet also follows the common setup of other few-shot learning methods. Specifically, few-shot classification is usually formalized as N-way K-shot classification problems. Let $S$ denote a support set that contains $N$ distinct image classes, and each class contains $K$ labeled samples. Given a query set $Q$, the purpose of few-shot learning is classifying each unlabeled sample in $Q$ according to the support set $S$. However, limited samples in $S$ make it difficult to efficiently train a network. Therefore, auxiliary set $A$ is always introduced to learn transferable knowledge to improve classification performance. Note that $S$ and $A$ have their own distinct label spaces without intersections.

In order to learn transferable knowledge better, the episode training mechanism [29] is adopted in the training phase. Specifically, at each iteration, support

set $AS$ and query set $AQ$ are randomly selected from auxiliary set $A$ to simulate a new few-shot classification task. In the training process, multiple episodes are constructed to train the model.

### 3.2   Overview

The overall framework of our method is shown in Fig. 1. First of all, images are fed into the feature extractor which is usually implemented by CNN or ResNet [11] to get image embeddings. In this stage, a LFE module is designed to explore local details with strong discriminability by the supervision of global features. Next, the features are used as inputs to the metric module. Our metric module adopts dual similarity that is composed of global and local metric branches, which can not only exploit the intra-class invariance of global features but also explore rich clues hidden in local details. Specially, in the local similarity branch, the discriminative patches are reweighted to eliminate noise, i.e., patches shared in the task, and enhance the significant regions. The proposed TDSNet focuses on the relationships among local patches rather than isolated individuals. Finally, the mean value of global and local classification scores is reported as the final result. LFE module is only used during training.
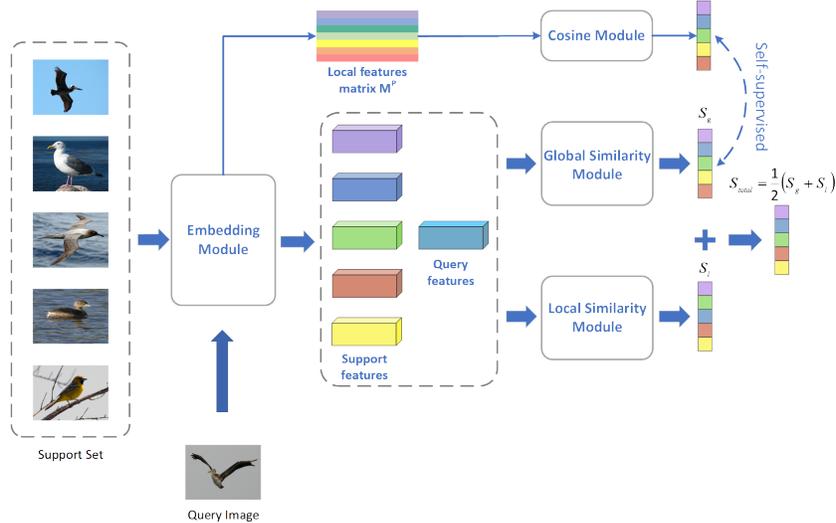


**Fig. 1.** Framework of the proposed TDSNet.

### 3.3   Local Feature Enhancement(LFE) Module

**Weakly supervised attention generation.** Parts of the objects are predicted first. In this paper, we explore discriminative regions in a weakly supervised

manner. Besides, instead of using a pre-trained convolutional neural network, we employ an attention generation strategy.

For the input image $X$, the feature $F \in R^{H \times W \times C}$ is explored through the feature embedding module $f_\varphi$, where $H$, $W$, and $C$ denote the height, width, and the number of channels of the feature respectively. Then, the attention maps $A^a \in R^{H \times W \times m}$ for each image can be determined by

$$A^a = f(F) = \bigcup_{k=1}^{m} A_k^a \tag{1}$$

$f(\cdot)$ represents a convolution operation and $A_k^a$ is the k-th local attention map. Similar to the bilinear pooling [10, 28], element-wise multiplication between $A^a$ and $F$ is performed to produce the part feature map $f$, which can be expressed as

$$f_k = g(A_k^a \odot F), k = 1, 2, \cdots, m \tag{2}$$

$f_k$ is the k-th local feature, $\odot$ denotes element-wise multiplication, and $g(\cdot)$ is the global average pooling operation. Finally, we stack these part maps to obtain the final part feature matrix. It can be represented as

$$M^P = \begin{pmatrix} f_1 \\ f_2 \\ \cdots \\ f_m \end{pmatrix} \tag{3}$$

**Local feature enhancement.** Weakly supervised attention generation is capable of activating some local parts of the objects, however, the discriminability of these local parts may not be guaranteed. Therefore, we propose a feature regularization method to constrain object representation, which extracts knowledge from global features to local features and guides the parts with strong discriminability to be encouraged. An effective way to achieve this effect is to match the prediction distributions between objects and their parts. Let $P_g$ and $P_a$ describe the predicted distributions of the global features and part features $M^P$ respectively. We optimize a KL divergence loss [19] that is applied for measuring the difference between two probability distributions as follows,

$$L_{KL(P_g \| P_a)} = -H(P_g) + H(P_g, P_a) \tag{4}$$

where $H(P_g) = -\sum P_g log P_g$, $H(P_g, P_a) = -\sum P_g log P_a$

This regularization loss forces the feature representation learning to focus on the discriminative details from a particular local region, by which we can further filter out unnecessary and misleading information to improve the discriminability of global features at fine-grained scales. Only the one with global features is used for final target prediction.

### 3.4   Dual Similarity Module

**Global similarity.** This branch adopts the global feature maps for classification, which employs the cosine distance as a metric function. The feature maps are

fed into two convolution blocks($h_{conv}$) to learn generic knowledge with global representations in the images, followed by a cosine similarity $cos(\cdot)$ to measure the similarity. For the query image $X_q$, the global similarity score corresponding to the support class $X_s$ is determined as follows,

$$S_g = cos(h_{conv}(\frac{1}{K}\sum_{s=1}^{K}f_\varphi(X_s)), h_{conv}(f_\varphi(X_q))) \quad (5)$$

We take the mean value of feature mappings from each support class as the class prototype, which is used to calculate the global structure so that the invariant features within the class can be learned.
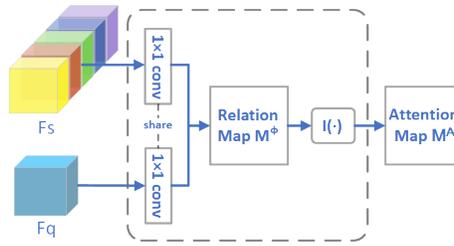


**Fig. 2.** The framework of task-aware attention.

**Task-aware local similarity.** Some recent works, such as DN4 and ConvM-Net, have shown that features based on local descriptors are richer than global. Specifically, local descriptors are able to capture local subtle information which is of greater benefit for fine-grained image recognition. For an image feature $F \in R^{H \times W \times C}$, it is regarded as a set of r(r=HW) C-dimensions local feature descriptors, which can be expressed as

$$f_\varphi(X) = [x_1, ..., x_r] \in R^{C \times r} \quad (6)$$

where $x_i$ denotes the i-th depth local descriptor. These local descriptors correspond to spatial local patches in the raw image. Basically, for each query image $X_q$, we get HW local descriptors to estimate its distribution, denoted as $L^q = f_\varphi(X_q) \in R^{C \times HW}$. Similarly, for support set, all local descriptors of class prototypes will be employed together as $L^s = f_\varphi(X_s) \in R^{C \times NHW}$. Next, we calculate the similarity matrix $M$ between query image and support set by

$$M_{i,j} = cos(L_i^q, L_j^s) \quad (7)$$

where $i \in 1, ..., HW$, $j \in 1, ..., NHW$, $cos(\cdot)$ represents cosine similarity. Each row in matrix $M$ represents the similarity of a specific query patch to all patches of the support set.

We then construct the task-aware attention map to reweight these parts. As shown in Fig. 2, we built another relation matrix $M^\phi$ for the next operation,

which is obtained by a convolution layer and a cosine similarity layer. We consider that local descriptors shared by multiple classes in the task do not contribute to the classification. For instance, if the query image patch has a high level of similarity to multiple patches in the support set, it has a minuscule contribution to classification. Therefore, we distract attention to make local descriptors that are shared in the task get relatively small attention values. The attention matrix $M^A$ is defined as

$$M_{i,j}^A = \frac{I(M_{i,j}^\phi)}{\sum_j I(M_{i,j}^\phi)} \tag{8}$$

$$I(x) = \begin{cases} x, \ if \ x > \beta \\ 0, \ otherwise. \end{cases} \tag{9}$$

$\beta$ is the threshold, which is set as the minimum of the top-k elements obtained by k-NN [2] from the relationship matrix $M^\phi$ to eliminate the noises. Since $I(\cdot)$ is indifferentiable, we approximate it by a variant function of sigmoid with a hyperparameter t as

$$I^*(x) = x/(1 + exp^{-t(x-\beta)}) \tag{10}$$

Theoretically, when t is large enough, it can be approximated as $I(\cdot)$. We perform element-wise multiplication between the weight matrix $M^A$ and the relation matrix $M$. Finally, the local similarity score for n-th class between the query image $X_q$ and the support class $X_s$ can be calculated by applying the attention map to the similarity matrix $M$ as follows:

$$S_l = \frac{1}{HW} \sum_{i=1}^{HW} \sum_{j=1}^{HW} (M^A \odot M)_{i,j} \tag{11}$$

The total classification score is formulated as follows which is used to make a final prediction:

$$S_{total} = \frac{1}{2}(S_g + S_l) \tag{12}$$

In particular, our local similarity branch only introduces a small number of parameters and the overfitting problem in the few-shot learning can also be alleviated to some extent.

### 3.5   Loss Function

In the training phase, the purpose is to learn a task agnostic network for classification. We can obtain $y_q^g$ and $y_q^l$ as two predicted results by global and local branches respectively. Then, predicted values are compared with the ground-truth label $y_q$ to calculate two classification losses.

$$L_q^g = \sum_{j}^{N} (y_{q,j}^g - y_{q,j})^2, q = 1, ..., |Q| \tag{13}$$

$$L_q^l = \sum_j^N (y_{q,j}^l - y_{q,j})^2, q = 1, ..., |Q| \tag{14}$$

In the end, the whole loss function can be written as:

$$L_{total} = L_q^g + L_q^l + \lambda L_{KL} \tag{15}$$

Where $\lambda$ is a trade-off parameter used to control the relative importance of the loss $L_{KL}$. Empirically, we set $\lambda = 0.4$.

## 4    Experiment

### 4.1    Datasets and Experimental Setting

We evaluate our method on three widely used fine-grained datasets, namely CUB-200-2011 [30], Stanford Dogs [12], and Stanford Cars [14]. We conduct experiments under the 5-way 1-shot and 5-way 5-shot settings. All images are resized to $84 \times 84$ before being fed into the feature extraction module. In the training process, episode training mechanism is used to train our model. For the three datasets, models are trained for 600 and 400 epochs corresponding to the 5-way 1-shot and 5-way 5-shot tasks, respectively. We use Adam optimizer to train the network with the initial learning rate of 0.001, decaying by half every 100,000 episodes. In the testing phase, the top-1 accuracy with 95% confidence interval will be reported by random sampling of 600 episodes from the test set.

**Table 1.** Comparison with typical FSL and FG-FSL methods on three fine-grained datasets. The best and the second best results are highlighted in red and green respectively.

| Dataset | CUB Birds | | Stanford Dogs | | Stanford Cars | |
|---|---|---|---|---|---|---|
| Setting | 1-shot | 5-shot | 1-shot | 5-shot | 1-shot | 5-shot |
| Matching Net [29] | 45.30±1.03 | 59.50±1.01 | 35.80±0.99 | 47.50±1.03 | 34.80±0.98 | 47.50±1.03 |
| Prototype Net [26] | 37.36±1.00 | 45.28±1.03 | 37.59±1.00 | 48.19±1.03 | 40.90±1.01 | 52.93±1.03 |
| Relation Net [27] | 59.58±0.94 | 77.62±0.67 | 43.05±0.86 | 63.42±0.76 | 45.48±0.88 | 60.26±0.85 |
| MAML [9] | 54.92±0.95 | 73.18±0.67 | 44.64±0.89 | 60.20±0.80 | 46.71±0.89 | 60.73±0.85 |
| PCM [32] | 42.10±1.96 | 62.48±1.21 | 28.78±2.33 | 46.92±2.00 | 29.63±2.38 | 52.28±1.46 |
| CovaMNet [15] | 52.42±0.76 | 63.76±0.64 | 49.10±0.76 | 63.04±0.65 | 56.65±0.86 | 71.33±0.63 |
| DN4 [16] | 46.84±0.81 | 74.92±0.64 | 45.41±0.76 | 63.51±0.62 | 59.84±0.80 | 88.65±0.44 |
| BSNet [17] | 65.89±1.00 | 78.48±0.65 | 51.68±0.95 | 67.93±0.75 | 54.39±0.92 | 73.37±0.77 |
| FOT [31] | 67.46±0.68 | 83.19±0.43 | 49.32±0.74 | 68.18±0.69 | 54.55±0.73 | 73.69±0.65 |
| ours | 69.34±0.89 | 80.34±0.59 | 54.48±0.87 | 69.45±0.69 | 62.14±0.91 | 75.64±0.72 |

### 4.2    Comparison with State-of-the-art Methods

To evaluate the validity of the proposed TDSNet, we conduct extensive experiments on three classic fine-grained datasets with 5-way 1-shot and 5-way 5-shot task settings and compare them with some SOTA methods.

As is demonstrated in Table. 1, the results compared with four classical few-shot methods and five fine-grained few-shot methods illustrate that our method achieves good performance on all three datasets. Specifically, the proposed TD-SNet performs better on the more challenging 1-shot task and shows high stability. The reason for this progress is that our approach focuses on the discriminative parts and gives them higher weights. Additionally, the changes in visual appearance may not affect our TDSNet because we also pay attention to invariant global structure.

### 4.3   Ablation Study

**Table 2.** Ablation study on the proposed components on CUB. LFE: local feature enhancement module, LS: local similarity module, att: task-aware attention.

| Method | 1-shot | 5-shot |
| --- | --- | --- |
| (a)Baseline | 63.33±1.01 | 77.64±0.67 |
| (b)Baseline+LFE | 65.20±0.99 | 78.77±0.67 |
| (c)Baseline+LFE+LS w/o att | 67.02±0.96 | 79.59±0.64 |
| (d)Baseline+LFE+LS w/ att | 69.34±0.89 | 80.34±0.59 |

**Effectiveness of local feature enhancement module.** We use the feature embedding module and the global similarity metric module of this paper as the baseline. Table. 2 shows that the addition of the local enhancement module makes the accuracy significantly improved by 1.87% and 1.13% respectively, which is mainly due to the activation of the features with strong discriminability. In this way, we can further filter out misleading information so as to improve the discrimination of features.

**Effectiveness of dual-similarity.** Compared with the results that only employ the global similarity measurement module, dual similarity demonstrates its superiority. It proves that only global features are hard to detect some detailed information that is suitable for fine-grained features, while only local features are sensitive to some intra-class changes.

**Effectiveness of task-aware attention.** We verify the effectiveness of task-aware attention on CUB, with 2.32% and 0.75% improvements respectively. Task-aware attention makes the network pay more attention to the features that are most relevant to the current task and reweight the key parts so that the features shared between classes will obtain less attention.

### 4.4   Visualization

As can be seen from Fig. 3, our TDSNet has less activation in the background and is more concentrated on the discriminative regions of objects, which demonstrates that the features can be semantically enhanced by our approach. It highlights the local details on raw images that represent different semantic patches with strong discriminability, such as the head and wings of a bird.
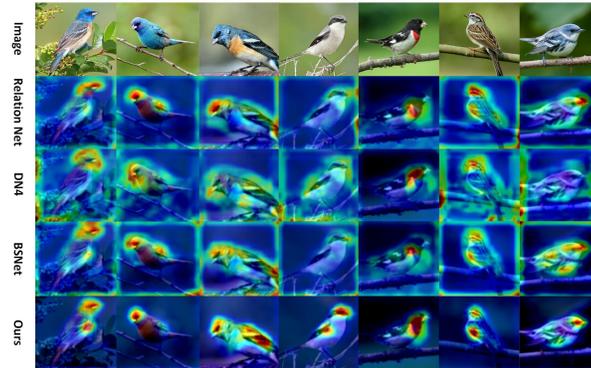
**Fig. 3.** Visualization of the features under Relation Network, DN4, BSNet and the proposed TDSNet on CUB. The redder the region, the more discriminative it is.

### 4.5 Number of Trainable Parameters

As Table. 3 shows, compared to the BSNet, which is also a bi-similarity method on the global feature, we have only half the number of parameters of this approach. This illustrates that although we adopt additional architectures to improve the performance, the proposed TDSNet only introduces a small number of the trainable parameters.

**Table 3.** Comparsion of the number of trainable parameters along on CUB.

| Model | Params | 5-way 5-shot |
|---|---|---|
| Prototype network | 0.113M | 45.28 |
| Relation network | 0.229M | 77.62 |
| DN4 | 0.113M | 74.92 |
| BSNet(R&C) | 0.226M | 78.84 |
| TDSNet(ours) | 0.114M | 80.34 |

## 5   Conclusion

In this paper, we propose a Task-aware Dual Similarity Network (TDSNet) for FG-FSL, which consists of two designed components, a local feature enhancement module and a measurement module that combines global and task-aware local similarity. Specifically, the former is designed to fully explore the discriminative details suitable for fine-grained classification, while the latter explores similarity by taking multiple perspectives, both global and local. Extensive experiments demonstrate that our proposed TDSNet achieves competitive results. In the future, we intend to reinforce the features of the foreground object and eliminate the negative effect of complicated backgrounds.

# References

1. Berg, T., Belhumeur, P.N.: Poof: Part-based one-vs.-one features for fine-grained categorization, face verification, and attribute estimation. In: Computer Vision & Pattern Recognition (2013)
2. Boiman, O., Shechtman, E., Irani, M.: In defense of nearest-neighbor based image classification. In: IEEE Conference on Computer Vision & Pattern Recognition (2008)
3. Branson, S., Horn, G.V., Belongie, S., Perona, P.: Bird species categorization using pose normalized deep convolutional nets. In: British Machine Vision Conference 2014 (2014)
4. Chang, D., Zheng, Y., Ma, Z., Du, R., Liang, K.: Fine-grained visual classification via simultaneously learning of multi-regional multi-grained features. arXiv **abs/2102.00367** (2021)
5. Chang, D., Ding, Y., Xie, J., Bhunia, A.K., Li, X., Ma, Z., et.al: The devil is in the channels: Mutual-channel loss for fine-grained image classification. IEEE Transactions on Image Processing **29**, 4683–4695 (2020)
6. Chang, D., Pang, K., Zheng, Y., Ma, Z., Song, Y.Z., Guo, J.: Your "flamingo" is my "bird": Fine-grained, or not. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 11471–11480 (2021)
7. Chen, H., Li, H., Li, Y., Chen, C.: Multi-scale adaptive task attention network for few-shot learning. CoRR **abs/2011.14479** (2020)
8. Ding, Y., Ma, Z., Wen, S., Xie, J., Chang, D., Si, Z., et.al: Ap-cnn: Weakly supervised attention pyramid convolutional neural network for fine-grained visual classification. IEEE Transactions on Image Processing **30**, 2826–2836 (2021)
9. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: International Conference on Machine Learning (2017)
10. Fu, J., Zheng, H., Tao, M.: Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In: IEEE Conference on Computer Vision & Pattern Recognition (2017)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
12. Khosla, A., Jayadevaprakash, N., Yao, B., Li, F.: L.: Novel dataset for fine-grained image categorization. CVPR workshop on fine-grained visual categorization (2013)
13. Koch, G., Zemel, R., Salakhutdinov, R.: Siamese neural networks for one-shot image recognition. In International Conference on Machine Learning **37** (2015)
14. Krause, J., Stark, M., Deng, J., Li, F.F.: 3d object representations for fine-grained categorization. In: IEEE International Conference on Computer Vision Workshops (2014)
15. Li, W., Xu, J., Huo, J., Wang, L., Luo, J.: Distribution consistency based covariance metric networks for few-shot learning. In: Thirty-Third AAAI Conference on Artificial Intelligence (2019)
16. Li, W., Wang, L., Xu, J., Huo, J., Gao, Y., Luo, J.: Revisiting local descriptor based image-to-class measure for few-shot learning. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 7260–7268 (2019)
17. Li, X., Wu, J., Sun, Z., Ma, Z., Cao, J., Xue, J.H.: Bsnet: Bi-similarity network for few-shot fine-grained image classification. IEEE Transactions on Image Processing **30**, 1318–1331 (2021)

18. Lifchitz, Y., Avrithis, Y., Picard, S., Bursuc, A.: Dense classification and implanting for few-shot learning. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
19. Luo, W., Yang, X., Mo, X., Lu, Y., Davis, L., Li, J., Yang, J., Lim, S.N.: Cross-x learning for fine-grained visual categorization. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 8241–8250 (2019)
20. Lz, A., Sh, B., Wei, L.A.: Learning sequentially diversified representations for fine-grained categorization. Pattern Recognition (2021)
21. Ma, Z., Lai, Y., Kleijn, W.B., Song, Y.Z., Wang, L., Guo, J.: Variational bayesian learning for dirichlet process mixture of inverted dirichlet distributions in non-gaussian image feature modeling. IEEE Transactions on Neural Networks and Learning Systems **30**(2), 449–463 (2019)
22. Ma, Z., Xie, J., Lai, Y., Taghia, J., Xue, J.H., Guo, J.: Insights into multiple/single lower bound approximation for extended variational inference in non-gaussian structured data modeling. IEEE Transactions on Neural Networks and Learning Systems **31**(7), 2240–2254 (2020)
23. Rao, Y., Chen, G., Lu, J., Zhou, J.: Counterfactual attention learning for fine-grained visual categorization and re-identification. In: IEEE/CVF International Conference on Computer Vision. pp. 1025–1034 (2021)
24. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: International Conference on Learning Representations (2017)
25. Shermin, T., Teng, S.W., Sohel, F., Murshed, M., Lu, G.: Integrated generalized zero-shot learning for fine-grained classification. Pattern Recognition **122**, 108246 (2022)
26. Snell, J., Swersky, K., Zemel, R.: Prototypical networks for few-shot learning. In: Advances in Neural Information Processing Systems. vol. 30 (2017)
27. Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P., Hospedales, T.: Learning to compare: Relation network for few-shot learning. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1199–1208 (2018)
28. Tao, H., Qi, H.: See better before looking closer: Weakly supervised data augmentation network for fine-grained visual classification. arXiv **1901.09891** (2019)
29. Vinyals, O., Blundell, C., Lillicrap, T., kavukcuoglu, k., Wierstra, D.: Matching networks for one shot learning. In: Advances in Neural Information Processing Systems. vol. 29 (2016)
30. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset. california institute of technology (2011)
31. Wang, C., Song, S., Yang, Q., Li, X., Huang, G.: Fine-grained few shot learning with foreground object transformation. Neurocomputing **466**, 16–26 (2021)
32. Wei, X., Wang, P., Liu, L., Shen, C., Wu, J.: Piecewise classifier mappings: Learning fine-grained learners for novel categories with few examples. IEEE Transactions on Image Processing **28**(12), 6116–6125 (2019)
33. Wu, Y., Zhang, B., Yu, G., Zhang, W., Wang, B., Chen, T., Fan, J.: Object-Aware Long-Short-Range Spatial Alignment for Few-Shot Fine-Grained Image Classification, pp. 107–115 (2021)