

Piecewise classifier mappings: Learning fine-grained learners for novel categories with few examples

Xiu-Shen Wei, *Member, IEEE*, Peng Wang, Lingqiao Liu, Chunhua Shen, *Member, IEEE*
Jianxin Wu, *Member, IEEE*

Abstract—Humans are capable of learning a new fine-grained concept with very little supervision, *e.g.*, few exemplary images for a species of bird, yet our best deep learning systems need hundreds or thousands of labeled examples. In this paper, we try to reduce this gap by studying the fine-grained image recognition problem in a challenging few-shot learning setting, termed few-shot fine-grained recognition (FSFG). The task of FSFG requires the learning systems to build classifiers for novel fine-grained categories from few examples (only one or less than five). To solve this problem, we propose an end-to-end trainable deep network which is inspired by the state-of-the-art fine-grained recognition model and is tailored for the FSFG task.

Specifically, our network consists of a bilinear feature learning module and a classifier mapping module: while the former encodes the discriminative information of an exemplar image into a feature vector, the latter maps the intermediate feature into the decision boundary of the novel category. The key novelty of our model is a “piecewise mappings” function in the classifier mapping module, which generates the decision boundary via learning a set of more attainable sub-classifiers in a more parameter-economic way. We learn the exemplar-to-classifier mapping based on an auxiliary dataset in a meta-learning fashion, which is expected to be able to generalize to novel categories. By conducting comprehensive experiments on three fine-grained datasets, we demonstrate that the proposed method achieves superior performance over the competing baselines.

I. INTRODUCTION

Fine-grained recognition tasks such as identifying the species of birds [1], dogs [2] and cars [3], have been popular in applications of computer vision. Since the categories are all similar to each other, different categories can only be distinguished by slight and subtle differences, which makes fine-grained recognition a challenging problem. Over the past decade, fine-grained recognition has attracted tremendous attention and observed rapid performance boost thanks to the integration of the sophisticated deep network structures with large annotated training datasets [4], [5], [6], [7], [8], [9], [10].

However, the large-scale fine-grained data volume required to train such classification algorithms limits the ranges where they can be successfully applied to, *e.g.*, very sparse training

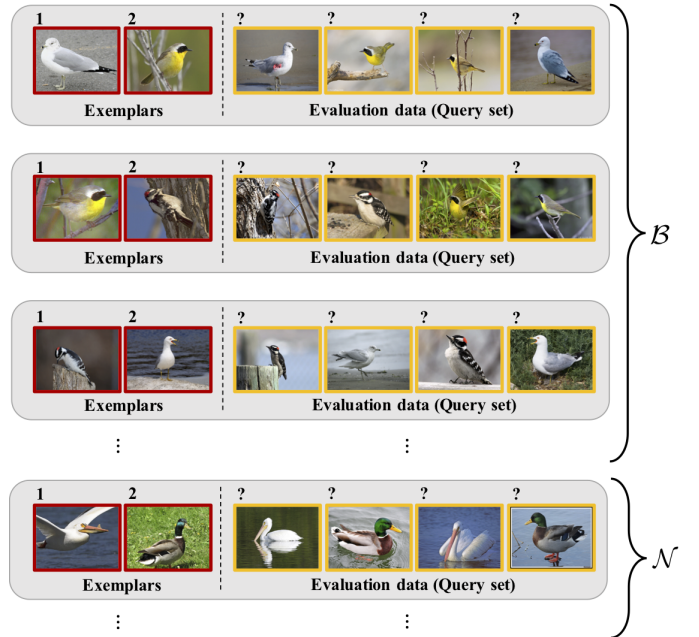


Figure 1. Illustration of the few-shot fine-grained image recognition (FSFG) task. The aim is to learn the classifier for a fine-grained category, bird species in this example, from few exemplars. We train the exemplar-to-classifier mapping based on an auxiliary dataset \mathcal{B} and test the FSFG performance on another dataset \mathcal{N} . There are no category overlaps between these two sets.

samples can be collected for some rare bird species. Humans, in contrast, are capable of learning a new fine-grained concept with very little supervision. To mimic this human ability, in this work, we study the fine-grained image recognition in a more practical and challenging few-shot setting, that is, we aim to learn the classifiers of novel fine-grained categories from very few labeled training examples (*a.k.a.* exemplars, usually 1 or 5).

Learning a classifier for a fine-grained category identified by few exemplars is a challenging problem, as satisfactory classification performance can be expected only when the learned classifiers can capture the subtle differences between categories and is able to generalize beyond the very limited supervisions. To realize such exemplar-to-classifier mapping, we propose an end-to-end trainable network which is inspired by state-of-the-art fine-grained recognition model and is tailored for the FSFG task. Specifically, the network consists of a bilinear feature learning module and a classifier mapping module. While the former encodes the discriminative information of exemplar image into a feature vector, the latter, as the key part of the network, maps the intermediate image features

• The first two authors contributed equally to this work. X.-S. Wei is with Megvii Research Nanjing, Megvii Technology, China. P. Wang is with University of Wollongong, Australia. L. Liu and C. Shen are with the University of Adelaide, Australia. J. Wu is with the National Key Laboratory for Novel Software Technology, Nanjing University, China.

• This work is supported by National Natural Science Foundation of China (61772256), DE170101259, and the program A for Outstanding Ph.D. candidate of Nanjing University (201702A010). C. Shen’s participation was in part supported by the CRC GeoVision Project.

• Email: {weixs.gm, wujx2001}@gmail.com, {lingqiao.liu, chunhua.shen}@adelaide.edu.au, p.wang6@hotmail.com.

into the category-level decision boundaries. Two problems remain to succeed with such mappings. On one hand, the distribution of the image-level representation can be complex which poses a great challenge for the mapping. On the other hand, the feature generated from bilinear pooling is very high dimensional, which further impedes the mapping due to the risk of parameter explosion.

The key novelty of our model to mitigate these problems is a “piecewise mappings” function in the classifier mapping module, which generates the decision boundary via learning a set of more attainable sub-classifiers in a much more parameter-economic way. Due to the outer product computation in bilinear pooling, the feature obtained, by nature, can be viewed as a set of sub-vectors, each of which implicitly attends to part of the image. We perform the sub-vector to sub-classifier mapping resorting to highly non-linear mappings. Then, these sub-classifiers are recombined into a global classifier so that it can tell samples from different categories. Intuitively, we learn the feature-to-classifier mapping based on the implicit “part” which may encode simpler and purer information and consequently makes the mapping easier. As a by-product, the piecewise mappings significantly reduce the number of model parameters and enable a more efficient computation. We learn the exemplar-to-classifier mapping using an auxiliary dataset in a meta-learning fashion as shown in Fig. 1. The aim in the meta-training phase is to learn a “mapping paradigm” which is expected to be able to generalize to novel categories.

In experiments, we perform the proposed FSFG method on three fine-grained benchmark datasets, *i.e.*, *CUB Birds* [1], *Stanford Dogs* [2], *Stanford Cars* [3]. Empirical results show that our FSFG model significantly outperforms competing baseline methods, including exemplar SVM [11], k -nearest neighbor and state-of-the-art generic few-shot learning methods [12], [13], [14]. Furthermore, we also conduct extensive ablation studies about our proposed method. These results could validate the effectiveness and efficiency of our FSFG model.

In summary, our major contributions are three-fold:

- We study the fine-grained image recognition in a challenging few-shot setting and propose a novel meta-learning strategy to address the FSFG problem.
- We devise a novel exemplar-to-classifier mapping strategy, named piecewise mappings, which resorts to the special structure of the bilinear CNN features to learn a discriminative classifier in a parameter-economic way.
- We conduct comprehensive experiments on three fine-grained benchmark datasets, and our proposed model achieves superior performance over competing solutions on all these datasets.

II. RELATED WORK

In this section, we briefly review the related work of both fine-grained image recognition and generic few-shot learning.

A. Fine-grained image recognition

Fine-grained recognition is a challenging problem and has recently emerged as an active topic [2], [3], [1]. Over the past decade, fine-grained recognition has achieved high performance

levels thanks to the integration of powerful deep learning techniques with large annotated training datasets. A number of effective fine-grained recognition methods have been developed in the literature [15], [4], [5], [6], [7], [8], [9]. Among them, some work, *e.g.*, [6], [8], attempted to learn a more discriminative feature representation by developing powerful deep models. Some methods aligned the objects in fine-grained images to eliminate pose variations and the influence of camera position, *e.g.*, [15], [7]. Moreover, some of them relied on localizing discriminative parts with/without strong supervisions, *e.g.*, [4], [5], [7].

However, current fine-grained recognition systems assume a set of categories known *a priori*, despite the obviously dynamic and open nature of the visual world [16], [17], [18]. Compared with previous work, we are studying fine-grained image recognition in a challenging few-shot learning setting where the model is required to recognize novel fine-grained categories by only a few labeled images.

B. Generic few-shot image recognition

Nowadays, few-shot image recognition (*a.k.a.* few-shot learning or low-shot learning) [16], [19], [20], [21] has attracted more and more attentions in computer vision and pattern recognition. This line of research explores the possibility of endowing learning systems the ability of rapid learning for novel categories from a few examples. More specifically, these systems are able to learn new concepts on the fly, from few or even a single example as in one-shot learning. Few-shot image recognition is usually tackled by using generative models [22], [23] or, in a discriminative setting, using ad-hoc solutions such as exemplar support vector machines [11]. While recently, many methods solved it in a learning-to-learn formulation [13], [24], [25], [18], [17], [26], [27], [28].

Specifically, in recent years, Vinyals et al. [19] proposed Matching Networks, which uses an attention mechanism over a learned embedding of the labeled set of examples (the support set) to predict classes for the unlabeled points (the query set). It can be interpreted as a weighted nearest-neighbor classifier applied within an embedding space. Later, Snell et al. [13] developed Prototypical Networks for generic few-shot learning. [13] further improved [19] by considering there exists an embedding in which points cluster around a single prototype representation for each class. It achieved better classification accuracy than [19] in the few-shot learning setting. In [17], it learned a regression network that maps from small-sample model parameters (*i.e.*, small-sample decision boundary) to large-sample model parameters (*i.e.*, large-sample decision boundary). Meanwhile, the method of [17] was also performed in a meta-learning fashion. Additionally, in [27], the authors reformulated the parameter update into an LSTM and achieved this via a meta-learner. To solve new learning tasks with few samples, the method in [24] designed a so called model-agnostic meta-learning scheme, the essential idea of which is to require the parameters to be able to perform well on new task via one or few gradient steps on this task. The method in [28] took a step further by updating the model parameters as well as the learning rate in a uniform meta-learning framework.

More similar to our works are [20], [21], which attempted to train parameter predictors for novel categories also from activations. However, the most differences between ours and [20], [21] are two-fold. 1) Our novel categories classifiers are learnt in a meta-learning fashion, while, [20], [21] employ traditional learning strategy. 2) More importantly, our proposed method can leverage the bilinear structure of powerful image representations for fine-grained objects. Besides, we also develop a novel classifier learning paradigm, *i.e.*, piecewise classifier mappings (a.k.a. sub-classifier mapping), which can not only prevent overfitting caused by high-dimensionality of bilinear, but also have a good motivation for the few-shot fine-grained recognition task. Experimental results validate our proposal and prove our learning strategy design.

Additionally, many previous few-shot image recognition studies all focused on generic images (*e.g.*, images of the ImageNet [29] and CIFAR [30] datasets) or generic patterns (*e.g.*, characters of the Omniglot [31] dataset). In fact, some generic few-shot learning methods, *e.g.*, [17] and [25], did consider fine-grained recognition scenarios and evaluate on fine-grained datasets. However, compared with those tasks, we specifically consider a novel few-shot image recognition topic, *i.e.*, few-shot fine-grained image recognition. The most different point of our topic from the generic few-shot image recognition is that, fine-grained recognition relies on more subtle image cues which makes it considerably more challenging. We demonstrate that the proposed model, especially our piecewise mappings component, can cater to the desire of capturing the subtle differences in a fine-grained scenario from limited training data, even one-shot.

III. LEARNING FEW-SHOT FINE-GRAINED LEARNERS

In this section, we firstly present our learning strategy for FSFG and introduce the relevant notations. Then, a detailed elaboration of various aspects of our method will be followed in the subsequent sections.

A. Learning strategy and notations

Our work is built upon the framework of meta-learning which treats the classifier generation process as a mapping function from the few labeled training samples of a category, called “exemplars” hereafter, to their corresponding category classifier. Fig. 2 shows the key idea of this learning scheme. This *exemplar-to-classifier* mapping is learned on an auxiliary training set \mathcal{B} . It contains N labeled training images $\mathcal{B} = \{(\mathcal{I}_1, y_1), (\mathcal{I}_2, y_2), \dots, (\mathcal{I}_N, y_N)\}$, where \mathcal{I}_i is an example image and $y_i \in \{1, 2, \dots, C_{\mathcal{B}}\}$ is its corresponding label. Once the mapping function is learned, it will be applied on another testing set \mathcal{N} to evaluate its performance, where \mathcal{N} contains images of novel categories that do not appear in \mathcal{B} .

To train the mapping function, we randomly sample a set of “meta-training sets” from \mathcal{B} . Each meta-training set (corresponding to a training episode) contains $C_{\mathcal{E}} < C_{\mathcal{B}}$ randomly chosen categories and a few images associated with them. A meta-training set is composed of an “exemplar set” \mathcal{E} and a “query set” \mathcal{Q} to mimic the scenario at the testing stage. Specifically, \mathcal{E} contains N_e (*e.g.*, 1 or 5) exemplar images per

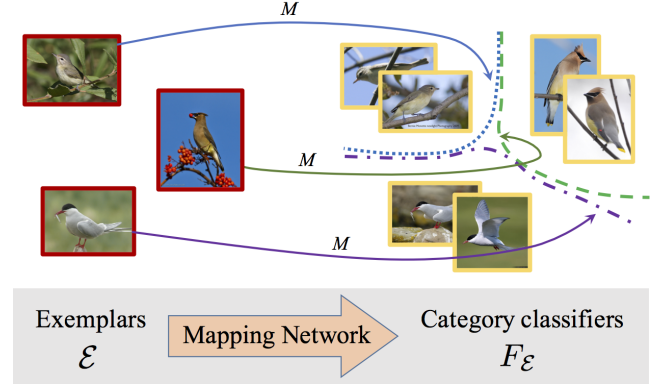


Figure 2. Key idea of the proposed FSFG model. In each episode, we sample an exemplar set \mathcal{E} from \mathcal{B} , which is composed of a subset of categories (three categories in this example) and each category contains few exemplars (the images with red border). We wish to learn a mapping M that can map these exemplars into their corresponding category classifiers (the dashed lines). The mapping parameters are learned so that these classifiers can correctly distinguish the query images (the images with yellow border).

category. The query set \mathcal{Q} is coupled with \mathcal{E} (has the same categories), but has no overlapped images. Each category of \mathcal{Q} contains N_q query images. During training, \mathcal{E} will be fed into the to-be-learned mapping function M to generate the category classifiers $F_{\mathcal{E}}$:

$$\mathcal{E} \xrightarrow{M} F_{\mathcal{E}}. \quad (1)$$

Then, $F_{\mathcal{E}}$ are subsequently applied to \mathcal{Q} for evaluating the classification loss. The training objective then amounts to learning the mapping function by minimizing the classification loss. This process is formally written as follows:

$$\min_{\lambda} \mathbb{E}_{\{\mathcal{E}, \mathcal{Q}\} \sim \mathcal{B}} \{\mathcal{L}(F_{\mathcal{E}} \circ \mathcal{Q})\}, \quad (2)$$

where λ denotes the model parameters of the mapping function M (from \mathcal{E} to $F_{\mathcal{E}}$), and \mathcal{L} is the loss function. $F_{\mathcal{E}} \circ \mathcal{Q}$ denotes applying the category classifiers $F_{\mathcal{E}}$ generated by the exemplar set \mathcal{E} on the query set \mathcal{Q} .

B. Model

We implement the above exemplar-to-classifier mapping by adopting a trainable neural network. Fig. 3 shows the overall architecture of the network. As we can see, the network is composed of two modules: a representation learning module and a classifier mapping module. While the former adopts a bilinear CNN structure to encode the discriminative information of an exemplar image into a high-dimensional feature vector, the latter, as the key part of the network, maps the intermediate image representation into a category classifier. In the next two sub-sections, we elaborate these two modules in more details.

1) *Representation learning*: We employ a bilinear CNN (BCNN) structure [8] to learn the image representation considering its state-of-the-art performance in fine-grained image recognition. BCNN consists of two feature extractors whose outputs are multiplied using outer product at each location of the image and pooled to obtain an image representation. Concretely, given two convolutional networks (A and B) as two streams of BCNN, we assume their outputs are re-organized

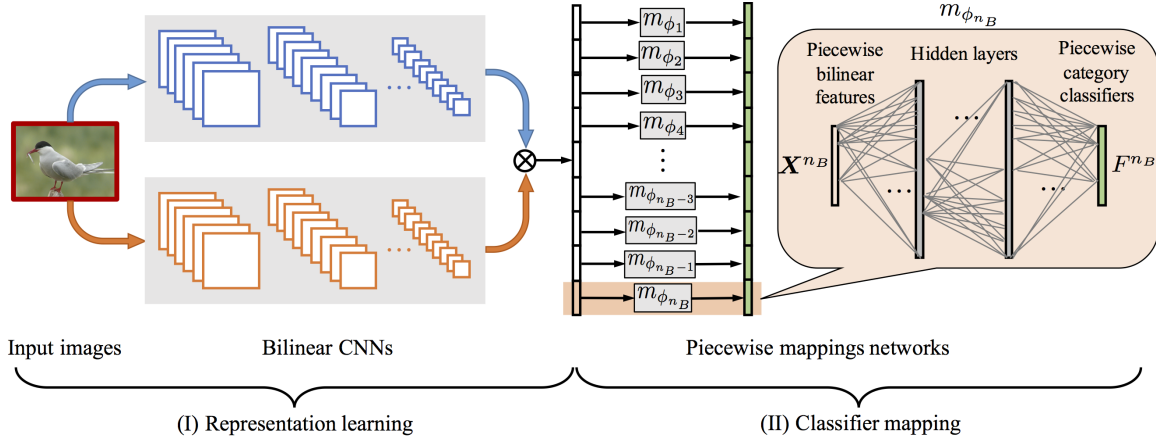


Figure 3. Overview structure of our proposed FSFG model. On the left, it is the first component (the bilinear pooling module) for representation learning. On the right, the second component (the classifier mapping module) maps the intermediate image features into the category classifiers.

into $f_A(\mathcal{I}) \in \mathbb{R}^{n_A \times L}$ and $f_B(\mathcal{I}) \in \mathbb{R}^{n_B \times L}$, where n_A, n_B denotes the dimensionality of the outputs and L denotes the spatial locations. Then, at location l , the bilinear representation will be $\mathbf{b}_l \in \mathbb{R}^{n_A \times n_B}$,

$$\mathbf{b}_l = f_A(l, \mathcal{I}) f_B(l, \mathcal{I})^\top. \quad (3)$$

The vectorized versions of $\{\mathbf{b}_l\}$ will be pooled over the entire image to derive the image representation $\mathbf{x} \in \mathbb{R}^{D \times 1}$ (for interpretation simplicity we let $D = n_A \times n_B$), that is,

$$\mathbf{x}(\mathcal{I}) = \sum_{l=1}^L \text{vec}(\mathbf{b}_l). \quad (4)$$

With the outer product computation, bilinear structure modulates one feature stream with another. Thus, the BCNN feature \mathbf{x} can be viewed as a set of n_B sub-vectors \mathbf{x}^t :

$$\mathbf{x} = [\mathbf{x}^1; \mathbf{x}^2; \dots; \mathbf{x}^t; \dots; \mathbf{x}^{n_B}], \forall t: \mathbf{x}^t \in \mathbb{R}^{n_A \times 1}, \quad (5)$$

where \mathbf{x}^t is the modulated feature of f_A by the t -th feature of f_B . This is similar to the multiplicative feature interactions in attention mechanisms [8]. From the observation that each modulated feature map tends to focus on an implicit “part” of an object (cf. Fig. 4), and thus, \mathbf{x}^t can be viewed as the feature description for that “part”. In our implementation, we train the bilinear CNN by performing the same procedure in [8] and use it as the image representation extractor.

To represent a set of N_e exemplar images belonging to category k , we simply compute the mean image representation as the category-level representation by:

$$\mathbf{X}_k = \frac{1}{N_e} \sum_{i=1}^{N_e} \mathbf{x}_i, \quad (6)$$

where $\{\mathbf{x}_i\}$ are samples with $y_i = k$.

2) *Classifier mapping*: Now that the information of each category identified by few exemplars has been encoded into a bilinear feature vector, the task of the classifier mapping module is to map these intermediate category-level representations into their corresponding category classifiers. Mathematically, this module computes a D -dimensional classifier $F_k \in \mathbb{R}^D$ for each category through a mapping $M: \mathbb{R}^D \rightarrow \mathbb{R}^D$.



(a) CUB Birds



(b) Stanford Dogs



(c) Stanford Cars

Figure 4. “Parts” of an object specifically correspond to the meaningful regions of the fine-grained objects, e.g., the beak of birds, the foot of dogs and the wheel of cars, etc. In the figures, by following [8], we show the patches with the highest activation for several random filters of the BCNN models used in our experiments on three datasets, respectively.

A straightforward solution to realize this mapping is via a global mapping, either linear or nonlinear. For example, a linear mapping can be:

$$F_k = \mathbf{W}_g \mathbf{X}_k + \mathbf{b}_g, \quad (7)$$

where $\mathbf{W}_g \in \mathbb{R}^{D \times D}$ and $\mathbf{b}_g \in \mathbb{R}^D$ denote the parameters of the global mapping. However, this mapping strategy suffers from two drawbacks. First, as the feature \mathbf{X}_k is supposed to encode the category-level information, the distribution of which can be highly complex. This poses a great challenge for the

global mapping to find a decision boundary in such a complex feature space. Second, since the bilinear feature tends to be high dimensional, this mapping may result in parameter explosion, which will make the network training hard or infeasible.

To mitigate these problems, we propose a novel “piecewise mappings” strategy, which exploits the structure of the bilinear features. As analyzed in Sec. III-B1, the bilinear feature \mathbf{X}_k can be viewed as a set of sub-vectors \mathbf{X}_k^t with each sub-vector describes an implicit “part” of the object. Intuitively, we can test if an object falls into the category described in the exemplars by checking whether each “part” of it is compatible with the exemplars. This motivates us to apply a piecewise mapping to first map each sub-vector \mathbf{X}_k^t into its corresponding sub-classifier F_k^t , and then combine these sub-classifiers together to generate the global category classifier. Fig. 3 shows this mapping with more details.

Concretely, a sub-vector \mathbf{X}_k^t is firstly mapped into a sub-classifier F_k^t via a nonlinear multilayer perceptron (MLP) $m_{\phi_t}(\cdot)$ as

$$F_k^t = m_{\phi_t}(\mathbf{X}_k^t). \quad (8)$$

We learn n_B such MLPs $\{m_{\phi_t}(\cdot)\}$ to derive n_B sub-classifiers $\{F_k^t\}$, and then these sub-classifiers are concatenated together to generate the global category classifier F_k :

$$F_k = [F_k^1; F_k^2; \dots; F_k^{n_B}]. \quad (9)$$

Essentially, our model simplifies the global mapping approach by assuming that the classifier for the t -th sub-vector is solely determined by the information from the t -th sub-vector in the exemplar set. Despite resulting more restrictive mapping function, this assumption makes the network much easier to train. Note that, this mapping scheme will significantly reduce the model parameters involved in classifier generation. Taking one-layer mapping for example, let’s assume $n_A = n_B = 512$. For the global mapping, it requires more than 512^4 parameters. For the proposed piecewise mappings, however, the number is reduced to about 512^3 . In addition, although there are parameter-economy variants of BCNN [32], our piecewise classifier mappings still show better performance. This suggests that the proposed classifier mapping function brings benefits more than merely reducing the model size (cf. Table II).

3) *Network training*: Given a query sample \mathbf{x} with label $y = c$, we compute its prediction distribution via softmax as:

$$p_M(y = c | \mathbf{x}) = \frac{\exp(F_c \cdot \mathbf{x})}{\sum_{c'} \exp(F_{c'} \cdot \mathbf{x})}. \quad (10)$$

The model parameters are trained via minimizing the negative log-likelihood $\mathcal{J}(\mathbf{x}, y) = -\log(p_M(c | \mathbf{x}))$. With this, we can now summarize the training in an episode as follows. First, we select an exemplar set \mathcal{E} from \mathcal{B} and learn/generate the classifiers $F_{\mathcal{E}}$. Then, we establish a query set \mathcal{Q} . The model parameters are optimized by minimizing $\mathcal{J}(\mathcal{Q})$. Algorithm 1 illustrates the training process in more details.

IV. EXPERIMENTS

In this section, we first describe the experimental settings, and then present the main results. Later, Ablation studies are given to further evaluate the effectiveness of our proposed classifier mapping strategy.

Algorithm 1 Training episode loss computation for the proposed piecewise mappings.

Require: \mathcal{B} is an auxiliary training set with N images belonging to $C_{\mathcal{B}}$ categories; \mathcal{B}_c denotes a subset of \mathcal{B} containing all images belonging to the c -th category; $C_{\mathcal{E}}$ denotes the number of categories in an exemplar set \mathcal{E} as well as a query set \mathcal{Q} for an episode; \mathcal{E}_k denotes the elements $(\mathbf{x}_i, y_i = k)$ in \mathcal{E} with element size N_e ; \mathcal{Q}_k denotes the elements $(\mathbf{x}_j, y_j = k)$ in \mathcal{Q} with element size N_q ; n denotes the number of piecewise mappings; $\text{RandomSample}(\mathcal{T}, N)$ denotes a set of N elements chosen uniformly at random from set \mathcal{T} , without replacement; \mathcal{S} denotes a category set and S_i denotes its i -th element.

- 1: Select a category subset \mathcal{S} for an episode
 $\mathcal{S} \leftarrow \text{RandomSample}(\{1, 2, \dots, C_{\mathcal{B}}\}, C_{\mathcal{E}})$;
- 2: **for** k in $\{1, 2, \dots, C_{\mathcal{E}}\}$ **do**
- 3: Select $\mathcal{E}_k \leftarrow \text{RandomSample}(\mathcal{B}_{S_k}, N_e)$;
- 4: Compute the category-level representation \mathbf{X}_k following Eq. 6;
- 5: Generate the category classifier F_k by Eq. 8 and Eq. 9;
- 6: Select $\mathcal{Q}_k \leftarrow \text{RandomSample}(\mathcal{B}_{S_k} \setminus \mathcal{E}_k, N_q)$;
- 7: **end for**
- 8: Initialize loss $\mathcal{J} \leftarrow 0$;
- 9: **for** k in $\{1, 2, \dots, C_{\mathcal{E}}\}$ **do**
- 10: **for** (\mathbf{x}, y) in \mathcal{Q}_k **do**
- 11: $\mathcal{J} \leftarrow \mathcal{J} + \mathcal{J}(\mathbf{x}, y)$;
- 12: **end for**
- 13: **end for**
- 14: $\mathcal{J} = \frac{\mathcal{J}}{C_{\mathcal{E}} \times N_q}$
- 15: Update model parameters by minimizing \mathcal{J} ;
- 16: **return** n piecewise mappings $[m_{\phi_1}; \dots; m_{\phi_n}]$.

Table I
CATEGORY SPLIT FOR THREE DATASETS. C_{total} DENOTES THE TOTAL NUMBER OF CATEGORIES IN A DATASET, $C_{\mathcal{B}}$ DENOTES THE NUMBER OF CATEGORIES IN \mathcal{B} AND $C_{\mathcal{N}}$ DENOTES THE NUMBER OF CATEGORIES IN \mathcal{N} .

# category	CUB Birds	Stanford Dogs	Stanford Cars
C_{total}	200	120	196
$C_{\mathcal{B}}$	150	90	147
$C_{\mathcal{N}}$	50	30	49

A. Datasets, setups and implementation details

Our experiments are conducted on three fine-grained benchmark datasets, i.e., *CUB Birds* (200 categories of birds, 11,788 images) [1], *Stanford Dogs* (120 categories of dogs, 20,580 images) [2], *Stanford Cars* (196 categories of cars, 16,185 images) [3]. For each dataset, we randomly split its original image categories into two disjoint subsets: one as the auxiliary training set \mathcal{B} , and the other as the FSFG testing set \mathcal{N} . Table I presents the details of the category split. For each category in \mathcal{B} , we follow the raw splits provided by these datasets to split the data into training and validation. While the former is used to train the parameters, the latter is used to monitor the learning process.

To mimic the testing condition, in each training episode, we set the category size of the exemplar set \mathcal{E} to be same as the number of categories in the testing set \mathcal{N} , i.e., $C_{\mathcal{E}} = C_{\mathcal{N}}$. Further we set $N_e = 1$ ($N_e = 5$) for one-shot learning (five-shot learning) and N_q is set to be 20 in all settings (by following the protocol in [13]). Similarly, during the testing phase, for each category in \mathcal{N} , we randomly choose one exemplar (five exemplars) for one-shot learning (five-shot

learning), and another 20 samples are randomly selected to evaluate the recognition performance. We repeat this evaluation process twenty times, and the mean classification accuracy is used as the evaluation criterion.

In theory, we can choose any network structures as the base network for our bilinear feature learning module. Since our key contribution is in the classifier mapping scheme, we choose AlexNet [33] as the two streams in BCNN, considering the trade off between its representation capacity and computational efficiency. Specifically, we adopt the AlexNet model pre-trained on the Places 205 database [34] to initialize the representation learning parameters. The reason why we use the Place dataset [34] instead of ImageNet [29] is to avoid the FGFS testing categories to be present in the pre-training dataset. We fine-tune the bilinear feature learning module on the auxiliary training set first and freeze it during the classifier learning process. For the classifier mapping module, without otherwise stated, we choose the mapping function m_{ϕ_t} to be a three-layer MLP, where 1024 hidden units are adopted in each layer and Exponential Linear Units (ELU) [35] is used in each layer as the non-linear activation function. SGD is used to optimize the parameters with learning rate of 0.1. We implement our model using the open-source library PyTorch.

B. Main results

We present the main results of FSFG by firstly introducing some baseline methods and then reporting the empirical results on these three datasets.

1) *Comparison methods*: In our experiments, we compare our proposed model to the following competitive baselines. Apart from the original bilinear CNN, we also implement a compact bilinear CNN [32] as the image feature extractor to facilitate the comparison, which enables much lower feature dimensionality but keeps almost the same classification discriminative ability [32]. For compact bilinear pooling, we follow the optimal settings suggested in [32]. The dimensionality of compact bilinear pooling representations is 8,192-d (much less than 65,536-d of fully bilinear pooling). In our empirical results, the results of compact bilinear pooling are denoted as “CB” in Table II, and the results of fully bilinear pooling are denoted as “FB”. Note that, most existing methods for generic few-shot learning are not applicable to our problem due to the formidable computation cost on high-dimensional bilinear features.

- **k -NN** (k -nearest neighbors): Following the testing setting introduced in Sec. IV-A, we choose one sample (five samples) for each category in \mathcal{N} as exemplar(s) and 20 samples in the same category for evaluation. We use the BCNN (either original or compact version) fine-tuned on \mathcal{B} as the image representation extractor, and nearest neighbor is adopted as the classifier to categorize the evaluation images. Specifically, the image representations are first ℓ_2 -normalized and cosine distance is used as the distance metric. Note that, for five-shot learning, the representations of five exemplars are averaged before normalization to serve as the category-level representation. This process will be repeated twenty times as for our

method. (This applies to all other baselines, so we omit this when introducing the following baselines.)

- **SVM** (support vector machine): After obtaining the bilinear representations for exemplars of the testing categories in \mathcal{N} , we train a classifier for each category based on these representations. In particular, for one-shot learning, this baseline becomes exemplar-SVMs [11].
- **Siamese-Net** [12]: As a standard metric-learning strategy, Siamese-Net is a competitive solution for few-shot learning. It learns a feature space in which images of the same category are close but images belonging to different categories are separated apart. We train a Siamese-Net based on \mathcal{B} by sampling pair-wise examples and the corresponding binary labels (“1” presents examples are from the same category and “0” is not.) Similar to [12], the regularized cross-entropy loss on the binary classifier is used. During evaluation, Siamese-Net could rank similarities between exemplars and testing data.
- **Prototypical Network** [13] is one of state-of-the-art generic few-shot learning methods. It learns a metric space via the meta-learning fashion. In the learned metric space, classification can be performed by computing distances to prototype representations of each class. Here, we compare it as a strong baseline in our few-shot fine-grained setting.
- **Relation Network** [14] is recently proposed for dealing with the few-shot generic image recognition problem. It develops a novel meta learning paradigm for few-shot learning. Specifically, a Relation Network is able to classify images of few classes by computing relation scores between query images and the few examples of each new class. Different from the other previous few-shot learning methods whose learning process occurs in the feature embedding, Relation Network can be seen as both learning a deep embedding and learning a deep non-linear metric (*i.e.*, a similarity function).
- **Global mapping**: As aforementioned in Sec. III-B2, an alternative solution to our proposed piecewise classifier mappings is global mapping. It follows the idea of the global feature to global classifier mapping by applying the mapping function directly on the category-level representation.

2) *Comparison results*: Table II presents the average accuracy rates of FSFG on the novel categories of three fine-grained datasets. For each dataset, we report both one-shot and five-shot recognition results. As shown in that table, our proposed model consistently and significantly outperforms the other baseline methods on these datasets.

Generally, we see the simple baseline k -NN performs well and it even outperforms other more sophisticated baselines on some settings, *e.g.*, on *Stanford Dogs*. This is due to the discriminative capacity of the bilinear CNN features. SVM observes more obvious advantage comparing to k -NN when exploiting five training exemplars. Siamese-Net, as another discriminative method, achieves comparable performance to SVM but is outperformed by our method. In addition, our proposed method also outperforms Prototypical Networks and Relation Networks by a large margin. This reflects our meta-

Table II

COMPARISON RESULTS (MEAN \pm STD.) ON THREE FINE-GRAINED DATASETS. THE HIGHEST AVERAGE ACCURACY OF EACH COLUMN IS MARKED IN BOLD. “•/•” DENOTES THAT OUR PROPOSED MODEL PERFORMS SIGNIFICANTLY BETTER/WORSE THAN THE CORRESPONDING METHOD BY THE PAIRWISE t -TEST WITH CONFIDENCE LEVEL 0.05. “FB” STANDS FOR USING THE FULLY BILINEAR POOLING REPRESENTATIONS, AND “CB” IS FOR USING COMPACT BILINEAR POOLING.

Method	<i>CUB Birds</i>		<i>Stanford Dogs</i>		<i>Stanford Cars</i>	
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
k -NN (FB)	38.85 \pm 3.43 •	55.58 \pm 0.84 •	24.53 \pm 2.36 •	40.30 \pm 2.34 •	26.99 \pm 2.91 •	43.40 \pm 1.68 •
k -NN (CB)	24.52 \pm 1.80 •	41.85 \pm 1.51 •	18.31 \pm 1.81 •	32.37 \pm 1.15 •	21.25 \pm 1.78 •	39.42 \pm 1.57 •
SVM (FB)	34.47 \pm 1.93 •	59.19 \pm 1.28 •	23.37 \pm 3.18 •	39.50 \pm 1.07 •	25.66 \pm 1.53 •	51.07 \pm 1.51 •
SVM (CB)	24.94 \pm 1.97 •	41.93 \pm 1.69 •	18.25 \pm 2.83 •	30.50 \pm 1.76 •	21.34 \pm 1.94 •	39.43 \pm 1.46 •
Siamese-Net (FB) [12]	37.38 \pm 1.53 •	57.73 \pm 1.38 •	23.99 \pm 1.66 •	39.69 \pm 1.17 •	25.81 \pm 1.67 •	48.95 \pm 1.31 •
Siamese-Net (CB) [12]	26.58 \pm 2.47 •	43.51 \pm 1.53 •	19.28 \pm 2.60 •	31.49 \pm 1.22 •	22.41 \pm 1.55 •	40.07 \pm 1.88 •
Prototypical Network (FB) [13]	38.96 \pm 1.43 •	58.62 \pm 1.65 •	25.05 \pm 1.34 •	40.42 \pm 1.54 •	25.33 \pm 1.87 •	49.03 \pm 1.60 •
Prototypical Network (CB) [13]	28.88 \pm 1.41 •	44.28 \pm 1.57 •	21.40 \pm 1.24 •	32.99 \pm 2.11 •	24.48 \pm 1.67 •	42.91 \pm 1.18 •
Relation Network (FB) [14]	39.68 \pm 1.19 •	59.39 \pm 1.50 •	26.11 \pm 1.14 •	41.55 \pm 1.88 •	25.98 \pm 1.30 •	49.66 \pm 1.19 •
Relation Network (CB) [14]	30.01 \pm 1.11 •	45.19 \pm 1.25 •	22.96 \pm 1.58 •	33.81 \pm 1.69 •	25.74 \pm 1.77 •	44.09 \pm 1.53 •
Global mapping (FB-)	24.12 \pm 1.39 •	34.59 \pm 1.77 •	20.55 \pm 1.48 •	30.93 \pm 1.91 •	20.50 \pm 1.60 •	30.58 \pm 1.82 •
Global mapping (CB)	25.42 \pm 2.22 •	36.37 \pm 1.04 •	20.77 \pm 2.75 •	32.33 \pm 2.11 •	20.24 \pm 1.94 •	32.66 \pm 1.86 •
Ours	42.10\pm1.96	62.48\pm1.21	28.78\pm2.33	46.92\pm2.00	29.63\pm2.38	52.28\pm1.46

Table III

COMPARISON RESULTS OF GLOBAL MAPPING AND PIECEWISE MAPPINGS (OUR PROPOSAL) ON THREE DATASETS. THE HIGHEST AVERAGE ACCURACY OF EACH COLUMN IS MARKED IN BOLD. “•” DENOTES THAT THE PIECEWISE MAPPINGS OUTPERFORM THE GLOBAL MAPPING WITH CONFIDENCE LEVEL 0.05 BY THE PAIRWISE t -TEST.

Method	<i>CUB Birds</i>		<i>Stanford Dogs</i>		<i>Stanford Cars</i>	
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
Global mapping	27.36 \pm 1.64 •	38.05 \pm 1.55 •	19.55 \pm 2.27 •	32.53 \pm 2.35 •	16.06 \pm 2.06 •	26.17 \pm 1.02 •
Piecewise mappings (Ours)	31.00\pm2.85	48.80\pm2.33	23.07\pm3.24	41.02\pm2.50	18.98\pm2.18	31.51\pm1.38

learning strategy can better generalize to unseen/novel fine-grained categories. For the global mapping, because BCNN generates image representation of ultra-high dimensionality (*i.e.*, 65,536 in our case), it is infeasible to learn a global mapping on such high-dimensional feature vectors. In order to realize the global mapping, we apply an additional linear mapping to first reduce 65,536-d features into 8,192-d feature vectors, and based on the low-dimensional features, we conduct the global mapping. It is denoted as “Global mapping (FB-)” in Table II. Specifically, the global mapping is also implemented as a three-layer networks. As seen, our proposed piecewise mappings significantly outperforms the global mapping. In ablation studies, we will further compare these two types of mapping schemes.

Another interesting observation here is that the few-shot recognition performance gap between FB and CB is large. Note that, both FB and CB are trained on the same training set and achieve comparable classification performance on the validation set. This phenomenon may be explained as that the CB feature is not suitable for similarity matching (*i.e.*, the experimental case of the testing set). It is an open problem worth future explorations.

In addition, we further investigate whether the proposed piecewise mapping idea works for existing generic few-shot learning approaches. Concretely, we apply our piecewise mapping module on the popular Prototypical Networks by learning a set of prototype features from the sub-vectors of the bilinear features. And the final classification of a sample is achieved by fusing the prediction scores from all the sub-features. Using similar hyper-parameters as the original Prototypical Networks, the modified Prototypical Networks achieve 40.16% \pm 1.37%

(60.18% \pm 1.43%), 26.98% \pm 1.32% (42.55% \pm 1.65%), 27.41% \pm 1.35% (51.49% \pm 1.37%) recognition accuracy in the one-shot (five-shot) setting on *CUB Birds*, *Stanford Dogs* and *Stanford Cars*, respectively. By comparing with Table II, the performance of the Prototypical Networks with piecewise mappings consistently outperforms original Prototypical Networks on all splits, which shows the effectiveness of our piecewise manner.

C. Ablation studies

To further inspect our piecewise mappings strategy for FSFG, we conduct ablation experiments on two aspects. First, we compare the global mapping and piecewise mappings on a fairer setting. Second, we investigate the influence of the mapping function m_{ϕ_t} variations on the FSFG performance. Finally, we also change the number of piecewise mappings (*i.e.*, n_B) to show its stability.

1) *Piecewise mappings vs. global mapping*: As aforementioned, due to high-dimensionality of bilinear feature, it is infeasible to learn a non-linear (even a simple linear) global mapping on the original bilinear features (*e.g.*, 65,536 dimensionality) in practice. To perform the global mapping, we modify the original AlexNet structure by reducing the number of units of the last convolution layer from 256 to 64. By doing this, the bilinear feature becomes $64 \times 64 = 4096$ -dimensionality, which is feasible to learn a non-linear global mapping. In experiments, a three-layer MLP acts as the global mapping. The hidden units number is selected via cross-validation based on a set of {4096, 8192, 16384, 20480}. Finally, 16,384 hidden units are selected because of its optimal performance.

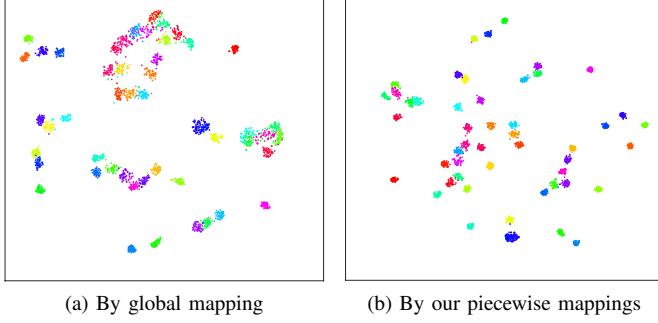
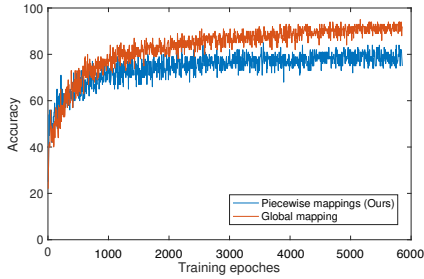
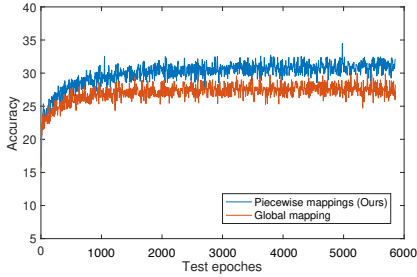


Figure 5. Visualization of the category classifiers generated by global mapping and piecewise mappings in 2D space by t-SNE [36]. Each dot denotes a generated classifier and different colors represent different categories. For each category, fifty classifiers are shown, each of which is obtained via randomly sampled five exemplars. This visualization is based on *CUB Birds*. (The figures are best viewed in color.)



(a) Learning curves of the training phrase.



(b) Learning curves of the test phrase.

Figure 6. Comparisons of learning curves of our proposed piecewise mappings with the global mapping. The blue curves indicate the learning behaviors of our proposed piecewise mappings, and the red curves are the global mapping.

For our proposed piecewise mappings, based on the modified BCNN, the piecewise mappings function is applied to 64-d sub-vectors. Totally, there are 64 piecewise mappings. Each of them is implemented as a three-layer network whose hidden layers contain 256 hidden units. ELU [35] is used as the activation function for both global mapping and piecewise mappings.

Table III demonstrates the comparison results of piecewise mappings vs. global mapping. Still the piecewise mappings significantly outperform the global mapping on all the three datasets. These observations can serve as a stronger evidence for the superiority of our proposed method.

Apart from the above quantitative evaluation, we present some qualitative results by visualizing the 4,096-d category classifiers generated by global mapping and piecewise mappings in the 2D space in Fig. 5. The dots with the same color

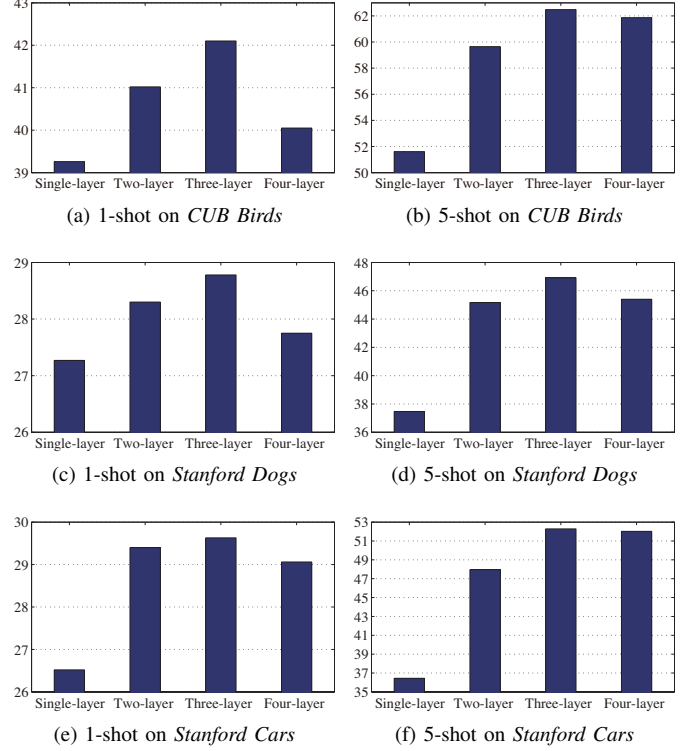


Figure 7. Ablation study on m_{ϕ_t} with different number of layers. In each sub-figure, the horizontal axis is the number of layers and the vertical axis represents the accuracy rate.

denote the classifiers generated from different exemplar images of the same category in \mathcal{N} . Different colors represent classifiers of different categories. We randomly select 250 exemplars per category to conduct five-shot recognition. Thus, one category contains 50 versions of classifiers (50 dots in the same one color). As shown in the figure, the classifiers generated by piecewise mappings exhibit better category-separability and more centralized intra-category aggregation. This, in some sense, reflects that the classifiers generated by our method tend to capture the essence of the corresponding categories and maintain better distinguishing capacity.

On the other hand, in theoretical aspect, the global mapping by fully connected layers is capable to learn the mapping function learned by our piecewise mappings method. In other words, the global mapping should have a larger representation capacity than the proposed piecewise mappings. But, why the global mapping performs worse like above? We hereby show the learning curves of the global mapping and our piecewise mappings in Fig. 6. It is clear to see that the global mapping (*i.e.*, the red curves) achieves higher training accuracy, while it gets worse test accuracy. The observation shows the global mapping has a lower generalization ability, which proves the global mapping is overfitting due to its larger representation capacity. Besides, this looks related to the regularization which constraints the feature mapping happening only in a subset of feature instead of the whole representation. While, thanks to the parameter economy brought by our piecewise mappings, it alleviates overfitting of high dimensional BCNN features.

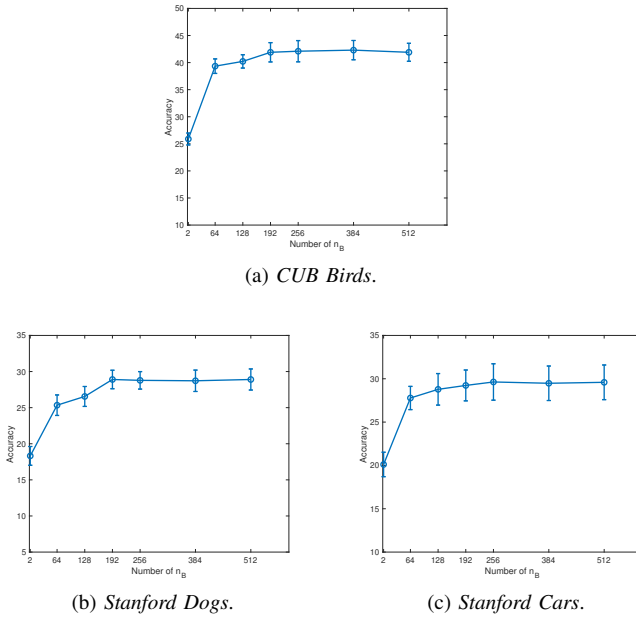


Figure 8. Comparisons of one-shot accuracy on *CUB Birds*, *Stanford Dogs* and *Stanford Cars* with different numbers of n_B .

2) m_{ϕ_t} with different numbers of layers: We implement the mapping functions m_{ϕ_t} in our classifier mapping module as MLPs. Since the depth plays an important role in determining the modeling capacity of MLPs, in this part, we investigate how the FSFG performance changes w.r.t. different number of layers in m_{ϕ_t} . Specifically, we change the number of layers from 1 to 4. The ablation study results are shown in Fig. 7.

Generally, we can see that a single-layer mapping leads to worst performance. This is due to its so limited modeling capacity that cannot realize the complex feature-to-classifier mapping. FSFG performance rises when adding another layer and peaks when three-layer mappings are used. Beyond that point, continuing to increase the depth of the mapping functions will do harm to the recognition performance, especially in the one-shot scenario. This study necessitates the need to apply a highly non-linear mapping to learn a satisfactory classifier.

3) Different numbers of n_B : In this section, we change the numbers of our piecewise mapping functions as the elements from a set of $\{2, 64, 128, 192, 256, 384, 512\}$. Meanwhile, we fix the number of n_A as 256. The comparisons are conducted on *CUB Birds*, *Stanford Dogs* and *Stanford Cars* in the one-shot setting. As shown in Fig. 8, it is obvious to observe that, when the number of n_B is large than 128, there is no significant accuracy gap with the number of 192, 256, 384 and 512. Except for those, when the number of n_B equals 2, there is a performance drop due to the extremely small representation ability. The observations of Fig. 8 also reveal the stability of our proposed method.

V. CONCLUSION

In this paper, we have presented the study on fine-grained image recognition in a practical and challenging few-shot learning setting, which requires to learn the classifier for a fine-grained category identified by few exemplars. To address this

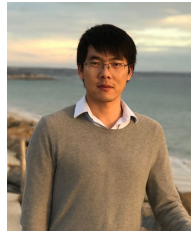
problem, we proposed an end-to-end trainable network which was inspired by the bilinear CNN model and was tailored for fine-grained few-shot learning. The key novelty of our network was the piecewise classifiers mapping module. By considering the special structure of bilinear CNN features, it decomposed the exemplar-to-classifier mapping into a set of more attainable “part”-to-“part classifier” mappings. As a by-product, it significantly reduced the model parameters. Through comprehensive experiments on three popular fine-grained image datasets, our method showed promising results.

In the future, it appears promising to use transfer learning techniques by leveraging the already gained experience (e.g., the classifiers of the known categories) based on the base set for generalizing the learning ability upon the novel set.

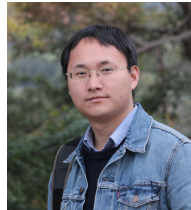
REFERENCES

- [1] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, “The Caltech-UCSD birds-200-2011 dataset,” *Technique Report CNS-TR-2011-001*, 2011.
- [2] A. Khosla, N. Jayadevaprakash, B. Yao, and L. Fei-Fei, “Novel dataset for fine-grained image categorization,” in *CVPR Workshop on Fine-Grained Visual Categorization*, Colorado Springs, CO, Jun. 2011, pp. 806–813.
- [3] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, “3D object representations for fine-grained categorization,” in *ICCV Workshop on 3D Representation and Recognition*, Sydney, Australia, Dec. 2013, pp. 554–561.
- [4] J. Fu, H. Zheng, and T. Mei, “Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition,” in *CVPR*, Honolulu, HI, Jul. 2017, pp. 4438–4446.
- [5] S. Huang, Z. Xu, D. Tao, and Y. Zhang, “Part-stacked CNN for fine-grained visual categorization,” in *CVPR*, Las Vegas, NV, Jun. 2016, pp. 1173–1182.
- [6] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu, “Spatial transformer networks,” in *NIPS*, Montréal, Canada, Dec. 2015, pp. 2008–2016.
- [7] D. Lin, X. Shen, C. Lu, and J. Jia, “Deep LAC: Deep localization, alignment and classification for fine-grained recognition,” in *CVPR*, Boston, MA, Jun. 2015, pp. 1666–1674.
- [8] T.-Y. Lin, A. RoyChowdhury, and S. Maji, “Bilinear Convolutional Neural Networks for Fine-Grained Visual Recognition,” *IEEE TPAMI*, vol. 40, no. 6, pp. 1309–1322, 2017.
- [9] Y. Zhang, X.-S. Wei, J. Wu, J. Cai, J. Lu, V.-A. Nguyen, and M. N. Do, “Weakly supervised fine-grained categorization with part-based image representation,” *IEEE TIP*, vol. 25, no. 4, pp. 1713–1725, 2016.
- [10] X.-S. Wei, J.-H. Luo, J. Wu, and Z.-H. Zhou, “Selective convolutional descriptor aggregation for fine-grained image retrieval,” *IEEE TIP*, vol. 26, no. 6, pp. 2868–2881, 2017.
- [11] T. Malisiewicz, A. Gupta, and A. A. Efros, “Ensemble of exemplar-SVMs for object detection and beyond,” in *ICCV*, Barcelona, Spain, Nov. 2011, pp. 89–96.
- [12] G. Koch, R. Zemel, and R. Salakhutdinov, “Siamese neural networks for one-shot image recognition,” in *ICML*, New York, NY, Jun. 2016, pp. 1–8.
- [13] J. Shell, K. Swersky, and R. Zemel, “Prototypical networks for few-shot learning,” in *NIPS*, Long Beach, CA, Dec. 2017, pp. 4077–4087.
- [14] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, “Learning to compare: Relation network for few-shot learning,” in *CVPR*, Salt Lake City, UT, Jun. 2018, pp. 1199–1208.
- [15] S. Branson, G. V. Horn, S. Belongie, and P. Perona, “Bird species categorization using pose normalized deep convolutional nets,” in *BMVC*, Nottingham, England, Sept. 2014, pp. 1–14.
- [16] L. Bertinetto, J. F. Henriques, J. Valmadre, P. H. S. Torr, and A. Vedaldi, “Learning feed-forward one-shot learners,” in *NIPS*, Barcelona, Spain, Dec. 2016, pp. 523–531.
- [17] Y.-X. Wang and M. Hebert, “Learning to learn: Model regression networks for easy small sample learning,” in *ECCV*, Amsterdam, Netherlands, Oct. 2016, pp. 616–634.
- [18] Y. Wang and M. Hebert, “Model recommendation: Generating object detectors from few samples,” in *CVPR*, Boston, MA, Jun. 2015, pp. 1619–1628.
- [19] O. Vinyals, C. Blundell, T. Lillicrap, K. Kavukcuoglu, and D. Wierstra, “Matching networks for one shot learning,” in *NIPS*, Barcelona, Spain, Dec. 2016, pp. 3630–3638.

- [20] S. Qiao, C. Liu, W. Shen, and A. Yuille, "Few-shot image recognition by predicting parameters from activations," in *CVPR*, Salt Lake City, UT, Jun. 2018, pp. 7229–7238.
- [21] H. Qi, M. Brown, and D. G. Lowe, "Low-shot learning with imprinted weights," in *CVPR*, Salt Lake City, UT, Jun. 2018, pp. 5822–5830.
- [22] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, "Human-level concept learning through probabilistic program induction," *Science*, vol. 350, no. 6266, pp. 1332–1338, 2015.
- [23] D. J. Rezende, S. Mohamed, I. Danihelka, K. Gregor, and D. Wierstra, "One-shot generalization in deep generative models," in *ICML*, New York, NY, Jun. 2016, pp. 1521–1529.
- [24] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *ICML*, Sydney, Australia, Aug. 2017, pp. 1–10.
- [25] Y.-X. Wang and M. Hebert, "Learning from small sample sets by combining unsupervised meta-training with CNNs," in *NIPS*, Barcelona, Spain, Dec. 2016, pp. 244–252.
- [26] S. Yeung, V. Ramanathan, O. Russakovsky, L. Shen, G. Mori, and L. Fei-Fei, "Learning to learn from noisy web videos," in *CVPR*, Honolulu, HI, Jul. 2017, pp. 1–9.
- [27] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," in *ICLR*, Toulon, France, Apr. 2017, pp. 1–11.
- [28] Z. Li, F. Zhou, F. Chen, and H. Li, "Meta-SGD: Learning to learn quickly for few shot learning," *arXiv preprint: 1707.09835*, 2017.
- [29] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet large scale visual recognition challenge," *IJCV*, vol. 115, no. 3, pp. 211–252, 2015.
- [30] A. Krizhevsky and G. E. Hinton, "Convolutional deep belief networks on CIFAR-10," *Technique Report*, 2010.
- [31] B. M. Lake, R. Salakhutdinov, J. Gross, and J. B. Tenenbaum, "One shot learning of simple visual concepts," in *Proceedings of Annual Meeting of the Cognitive Science Society*, Boston, MA, 2011, pp. 1–6.
- [32] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell, "Compact bilinear pooling," in *CVPR*, Las Vegas, NV, Jun. 2016, pp. 317–326.
- [33] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *NIPS*, Lake Tahoe, NV, Dec. 2012, pp. 1097–1105.
- [34] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, "Learning deep features for scene recognition using places database," in *NIPS*, Montréal, Canada, Dec. 2014, pp. 487–495.
- [35] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (ELUs)," in *ICLR*, San Juan, Puerto Rico, May. 2016, pp. 1–14.
- [36] L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *JMLR*, vol. 9, pp. 2579–2605, 2008.



Peng Wang received the B.S. and M.S. degrees from School of Electronic and Information Engineering, Beijing Jiaotong University, in 2009 and 2012, respectively, and the Ph.D. degree from the School of Information Technology and Electrical Engineering, The University of Queensland. He is currently a lecturer with University of Wollongong. Prior to this, he was a Post-Doctoral Researcher with The University of Adelaide. His research interests include image classification, video analytics, and deep learning.



Lingqiao Liu received his Ph.D. degree from the Australian National University in 2014. He then joined the University of Adelaide as a research fellow. He is now a DECRA research fellow and Lecturer with the school of computer science at the University of Adelaide. He is the recipient of the Discovery Early Career Researcher Award from the Australia Research Council, and the University of Adelaide Research Fellowship. His current research includes deep learning and its application in computer vision and natural language processing. He served as an associate editor for IEEE Robotics and Automation Letters and a reviewer/PC member for multiple international journals (e.g., TPAMI, TNN, TIP, TCSVT) and conferences (e.g., CVPR, ICCV, IJCAI, AAAI).



Chunhua Shen is a Professor at School of Computer Science, University of Adelaide. He studied at Nanjing University, at Australian National University, and received his PhD degree from the University of Adelaide. From 2012 to 2016, he held an Australian Research Council Future Fellowship.



Xiu-Shen Wei (M'18) received his BS degree in computer science, and his Ph.D. degree in computer science and technology from Nanjing University. He is now the Research Lead of Megvii Research Nanjing, Megvii Technology, China. He has published about twenty academic papers on the top-tier international journals and conferences, such as IEEE TPAMI, IEEE TIP, IEEE TNNLS, Machine Learning Journal, CVPR, ICCV, IJCAI, etc. He achieved the first place in the Apparent Personality Analysis competition (in association with ECCV 2016), the

first place in the iNaturalist competition (in association with CVPR 2019) and the first runner-up in the Cultural Event Recognition competition (in association with ICCV 2015) as the team director. His research interests are computer vision and machine learning. He has served as a PC member of ICCV, CVPR, ECCV, NIPS, IJCAI, AAAI, etc.



Jianxin Wu (M'09) received his BS and MS degrees in computer science from Nanjing University, and his PhD degree in computer science from the Georgia Institute of Technology. He is currently a professor in the Department of Computer Science and Technology at Nanjing University, China, and is associated with the National Key Laboratory for Novel Software Technology, China. He has served as an area chair for CVPR, ICCV and AAAI. His research interests are computer vision and machine learning.