# Fine-grained Classification via Categorical Memory Networks

Weijian Deng, Joshua Marsh, Stephen Gould, and Liang Zheng

arXiv:2012.06793v1 [cs.CV] 12 Dec 2020

*Abstract*—**Motivated by the desire to exploit patterns shared across classes, we present a simple yet effective class-specific memory module for fine-grained feature learning. The memory module stores the prototypical feature representation for each category as a moving average. We hypothesize that the combination of similarities with respect to each category is itself a useful discriminative cue. To detect these similarities, we use attention as a querying mechanism. The attention scores with respect to each class prototype are used as weights to combine prototypes via weighted sum, producing a uniquely tailored response feature representation for a given input. The original and response features are combined to produce an augmented feature for classification. We integrate our class-specific memory module into a standard convolutional neural network, yielding a Categorical Memory Network. Our memory module significantly improves accuracy over baseline CNNs, achieving competitive accuracy with state-of-the-art methods on four benchmarks, including CUB-200-2011, Stanford Cars, FGVC Aircraft, and NABirds.**

*Index Terms*—**Fine-grained Classification, Categorical Memory Module, Inter-class Similarity**
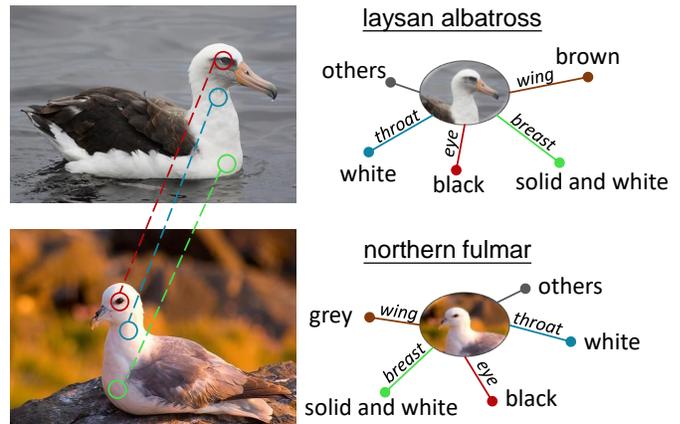


Fig. 1. Visual example of the correlation between fine-grained classes. Though the two birds belong to different subspecies, they are semantically and visually similar. Both share some common attributes (the first column), *e.g.*, eye color, breast pattern, and throat color. This similarity indicates their representations will share some underlying semantic cues.

## I. INTRODUCTION

**T**HIS article studies fine-grained classification, one of the most inherently challenging visual recognition tasks. Standard classification techniques are able to rely upon large-scale global features to differentiate classes. Meanwhile, fine-grained datasets simultaneously exhibit both subtle inter-class variation, such as global structure and visual appearance, and large intra-class variation, such as pose. This has led most methods to rely on *local differences* to provide discriminative cues [1], [2], [3], [4]. We propose an orthogonal research direction that instead seeks to exploit *inter-class similarities*, which is an essential characteristic of the fine-grained classification.

Convolutional neural networks (CNNs) tend to learn features that discard some semantic information in favour of more discriminative visual patterns. In the case of fine-grained classification where discriminative cues are particularly subtle, the loss of semantic information may lead to a lack of generalization. Our insight is that encouraging the network to represent similarities in addition to unique discriminative cues is beneficial for robust fine-grained feature learning. Consider the case study of bird classification. Different species share both semantic and visual attributes such as eye color, throat color, and breast pattern, as illustrated in Figure 1. Instead of discarding these shared semantics, we pose they could be used to produce more informative and robust features. A

W. Deng, J. Marsh, S. Gould, and L. Zheng (corresponding author) are with Australian National University, CBR, Australia.

Email: firstname.lastname@anu.edu.au

natural question arises: how can we use these shared inter-class semantics to improve the fine-grained classification?

Motivated by this, we propose a class-specific memory module that gives a network direct access to the statistics of class-specific feature representations. Our module is interpretable, with each slot in memory corresponding to the average feature of a given category, which we refer to as a class *prototype*. The class-specific memory module can be directly integrated into popular CNN architectures without any additional modifications. It is trained in an end-to-end supervised manner, with the memory module updated while training along with the rest of the network parameters. Given a labelled feature, we update the corresponding memory slot via moving average. Thus, each memory slot is encouraged to represent the prototypical information of its corresponding class.

Intuitively, each class should exhibit a specific distribution of similarities with respect to the other classes. We hypothesize that this distribution is itself a useful identifier of class. The rich semantic correlation between classes provides beneficial guidance for fine-grained feature learning. Specifically, we want to give neural networks the ability to perceive the shared semantics across classes while preserving the unique and subtle cues of the given input sample. To this end, we propose an attention-based method of addressing memory, where each prototype is retrieved based on its relevance to the current input feature. We define relevance as the attention score based on cosine similarity in the learned feature space. We use the

attention score to retrieve relevant class prototypes with respect to a given input, producing a response feature. Finally, the input feature and the response feature are combined, resulting in an augmented feature. This augmented feature is characterized by 1) expressing inter-class similarities of the input while 2) maintaining the individuality of input representations.

The proposed method is simple, effective, adds no additional parameters, and can be applied across non-rigid (e.g, birds) and rigid (e.g, cars) visual datasets. Furthermore, the entire memory module, in practice a matrix, is itself negligible in terms of computational overhead (including reading/writing) and adds no additional trainable parameters.

The key contributions of this paper include:

- We introduce a novel Categorical Memory Network (CMN) for fine-grain visual classification. Our method incorporates both feature learning and inter-class knowledge to yield more expressive features.
- To leverage the knowledge of all categories, we introduce an attention-based memory addressing method to adaptively fuse prototypes.
- CMN yields significant improvement over baselines, achieving competitive accuracy with the state-of-the-art accuracy on four fine-grained visual classification benchmarks.

## II. RELATED WORK

**Fine-grained Feature Learning.** Fine-grained visual categorization is a special case of image classification wherein there exists significantly more inter-class similarities than standard classification tasks (e.g ImageNet). Fine-grained datasets share the internal opposing dynamics of subtle inter-class and large intra-class variation, which together represent the central challenge for this task.

Deep Convolutional Neural Networks (CNNs) [5], [6], [7] are unable to sufficiently represent discriminative information required to classify fine-grained images. To obtain more expressive representations, several feature encoding methods have been proposed [8], [9], [10], [11], [12], [13]. Luo et al. [12] introduce a cross-layer regularizer to use multi-scale features for classification. Lin et al. [8] present a bi-linear pooling method to extract features with two Convolutional Networks. Gao et al. [9] introduce a bi-linear pooling in a kernelized framework, improving the computational efficiency of bi-linear pooling. Kong et al. [10] introduce a low-rank bi-linear classifier to decrease the computation time. Furthermore, Yu et al. [11] propose a hierarchical bi-linear pooling approach to fuse multi-layer features for fine-grained classification.

Another area of research focuses on finding discriminative regions in images. Early studies [14], [15], [16], [17], [18], [19] make use of bounding boxes and part annotations to localize discriminative regions. While such supervised annotations have proven beneficial to fine-grained feature learning, they are costly to obtain. More recently, weakly supervised methods have emerged as a promising research area [20], [21], [1], [2], [3], [22], [23], [4], [24]. These methods use image category labels to localize object parts and extract part features for classification. Yang et al. [1] use a navigator network to detect

informative regions. Zheng et al. [2] group ConvNet channels to detect various local patterns. Similarly, Fu et al. [3] recursively localizes parts at multiple scales while Sermanet et al. [22] uses an attention mechanism to select a series of regions. Sun et al. [23] introduce a one-squeeze multi-excitation module to learn region features. Moreover, Recasens et al. [4] recently proposed the use of saliency maps to localize informative regions.

In addition to the above approaches, some works focus on the small inter-class variations [25], [26], [27]. Zhou et al. [27] exploit the rich relationships among categories through bipartite-graph labels. Dubey et al. [25] introduce a pairwise confusion loss to prevent CNNs from over-fitting. The confusion loss attempts to bring class conditional probability distributions closer to each other. This enables the network to preserve some shared semantic among classes. Our CMN explicitly encourage the network to leverage shared prototypical information across classes for enhancing the input features. The class-specific memory module gives a network direct access to the statistics of class-specific feature representations. Thus, the network can adaptively retrieve relevant feature prototypes of different classes to augment input features. Moreover, our method focuses on feature combinations after obtaining image features, and as such is compatible with the existing approaches.

**Memory Networks.** Memory networks [28], [29], [30] are a unique variant of neural networks that allow a network to explicitly read and write information to an external memory module. Inputs are sequentially saved to memory and retrieved based upon the relevance (usually calculated via attention) to the current network input. Weston et al. [28] introduced memory-augmented networks for the task of question answering. Santoro et al. [31] use an external memory to tackle meta-learning tasks. Li et al. [32] and Kim et al. [33] use memory networks for data generation. The closely related concept of Neural Turing Machines [34] provides more complex functionality such as continuous memory representations and location and content-based access. We study memory networks in the context of fine-grained visual classification and introduce a categorical memory module to enhance the fine-grained feature learning.

In the class-specific memory module, each row of the module corresponds to the prototypical information of a single class in the form of a moving average. This practice is relevant to prototypical networks [35], [36], [37] in few-shot setting. However, our work is significantly different from them. These works use prototypes as class centers for the nearest-neighbor classifier. In comparison, we use relevant class prototypes to enhance the input features.

## III. CATEGORICAL MEMORY NETWORK

The proposed class-specific memory module can be directly integrated into an existing convolutional neural network without additional modifications. We refer to this network as a Categorical Memory Network (CMN). In practice, we insert the memory module after the last convolution layer of a convolutional neural network, e.g., ResNet-50 [6].

An overview of CMN is shown in Figure 2. We first pass the image through the convolutional layers and the final global average pooling, producing an original feature. We then use
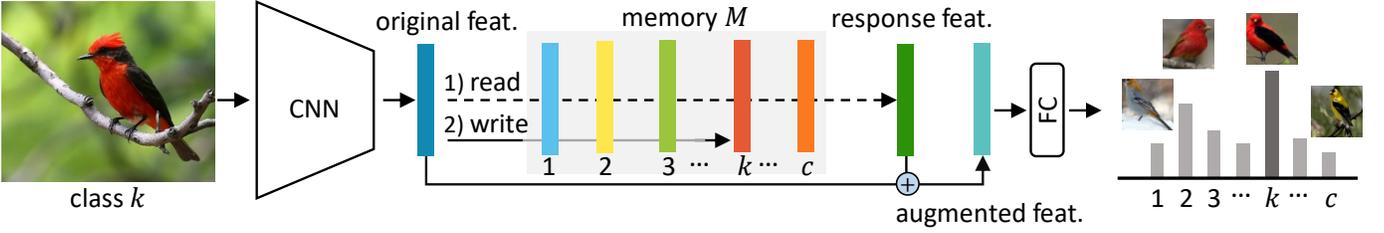
Fig. 2. Architecture of Categorical Memory Network (CMN) for fine-grained classification. CMN integrates a class-specific memory module into an existing convolution neural network architectures (*e.g.*, ResNet-50). The module contains memory slots to record category prototypes. Given an original feature produced by the CNN as a query, we retrieve and combine the relevant prototypes, and then output a response feature. We then update the corresponding memory slot based on the label of the original feature. Finally, we augment original feature with response feature. The resulting augmented feature is used for classification.

this original feature as the input to the class-specific memory module to retrieve the relevant prototypes, leading to a response feature. This response feature is then combined with the original feature via a weighted sum with learned coefficients, producing an inter-class similarity aware augmented feature. Finally, this augmented feature is fed through a single fully connected layer and a softmax to obtain class predictions.

The network is trained in an end-to-end manner via class-based supervised learning, *i.e.*, we use cross-entropy loss. In the following section, we introduce the categorical memory module and detail its reading and writing procedures.

### A. Class-specific Memory Module

Motivated by the intuition that giving a network access to leverage class-specific information is beneficial to fine-grain feature learning, we introduce a memory module to store the mean feature vector for each class. We refer to this as the feature prototype of a class. The memory module primarily consists of two operations, 1) writing the feature prototypes to memory module, and 2) computing response feature based on the current input feature.

During *training*, the memory slots, i.e rows, are read from based on the relevance to the given input. We then update memory slots based on the categories of given inputs in a mini-batch. During *testing*, the memory is not updated and used only for reading.

*1) Memory Writing:* We implement memory as a matrix $\mathbf{M} \in \mathbb{R}^{\mathbf{C} \times \mathbf{D}}$, where $\mathbf{C}$ is the number of classes and $\mathbf{D}$ is the dimension of the prototype. Each row $\mathbf{m}_i$ of $\mathbf{M}$ represents the feature prototype of the $i$-th class. Instead of updating the prototypes with respect to the entire training set, we update by moving average per mini-batch[1]. Given a feature $\mathbf{f}_i$ from $i$-th category, the corresponding prototype $\mathbf{f}_i$ is computed as,

$$\mathbf{m}_i \leftarrow \mathbf{m}_i + \beta(\mathbf{f}_i - \mathbf{m}_i), \tag{1}$$

where $\beta \in (0, 1)$ is a coefficient that controls the update rate of prototype $\mathbf{m}_i$.

*2) Memory Reading:* Our goal is to borrow some prototypical knowledge from relevant classes to augment the input feature, yielding more informative and comprehensive representations.

[1]We perform the moving average update with the current mini-batch samples after memory reading.
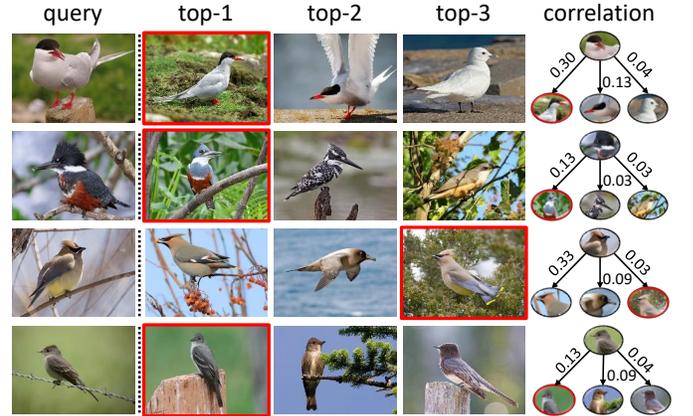


Fig. 3. Visualization of memory addressing. Our module is able to retrieve the relevant class prototypes for a given query image (the first column). In this figure, we show exemplar images from the three most relevant prototypes (the second to fourth columns) with respect to a given query image. We highlight red the prototype of the same class. In the last column, we show the attention scores between the input image and each retrieved prototype.

We use attention as our addressing mechanism. Given an input feature as a query, our module produces a response feature that is a uniquely tailored linear combination of all class prototypes in memory. Specifically, we calculate the attention score of the input feature $\mathbf{f}$ with respect to each prototype $\mathbf{m}_i$ (for $i \in \{1, \ldots, C\}$). These attention scores are used as weights to combine the prototypes, resulting in a response feature. We define the memory addressing process as,

$$\widetilde{\mathbf{m}} = \sum_{i=1}^{C} w_i \mathbf{m}_i, \tag{2}$$

where $w_i$ is the attention score for the prototype $\mathbf{m}_i$, and $\widetilde{\mathbf{m}}$ is the response feature.

**Attention Score.** Intuitively, if a prototype is close to input feature in the representation space, the prototype is more likely to contain relevant information and is weighted accordingly. In practice, we use cosine similarity distance as our similarity metric. We normalize the similarity scores by a softmax function, giving the final attention scores. This process is defined as,

$$w_i = \frac{\exp(\mathbf{f} \cdot \mathbf{m}_i / \tau)}{\sum_{j=1}^{C} \exp(\mathbf{f} \cdot \mathbf{m}_j / \tau)}, \tag{3}$$
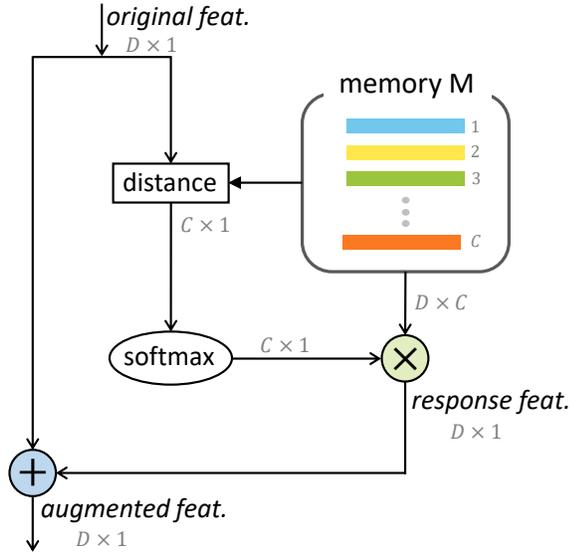
Fig. 4. Pipeline of memory reading process. Given an original feature, we calculate its Euclidean distance to each class prototype, which we then normalize to produce attention scores. After that, we use attention scores as weights to combine prototypes, resulting a response feature (response feat.). Finally, we combine the original feature and the response feature by using a weighted summation, yielding an augmented feature (augmented feat.) that is used for classification. In this figure, $\otimes$ is matrix multiplication, and $\oplus$ is broadcast element-wise addition.

where $w_i$ is the normalized attention score between input feature $\mathbf{f}$ and prototype $\mathbf{m}_i$ of class $i$, the temperature $\tau$ is empirically set as 0.1. In our case, each prototype represents one fine-grained class. Classes with strong similarities to an input image should receive higher attention scores and hence be weighted higher in the memory response.

An illustration of this is shown in Figure 3. We show exemplary images from the three most relevant prototypes for a given input image query. The query image and the images from the classes represented by the retrieved prototypes share visual and semantic similarities, *e.g.*, color, wing pattern, bill shape, and breast pattern. Moreover, we also observe that the input image is more likely to select the prototype corresponding to its class. Based on this, we think that our attention-based memory addressing method is effective to detect relevant class prototypes for the given image.

**Feature Augmentation.** The original feature, which is instance-specific, and response feature, which is class-specific, are combined to produce an augmented feature representation that is used for classification. Specifically, we use the following summation:

$$\mathbf{f}_{aug} = \mathbf{f} + \widetilde{\mathbf{m}}. \tag{4}$$

We also attempted more sophisticated methods to combine two features such as element-wise scaling and fully-connected layers, but these did not improve performance. An example of memory reading pipeline is illustrated in Figure 4. The reading process can be simply done by matrix operations.

### B. Discussion

**Working mechanism.** The learning mechanism of our method is that it makes the network adaptively exploit the
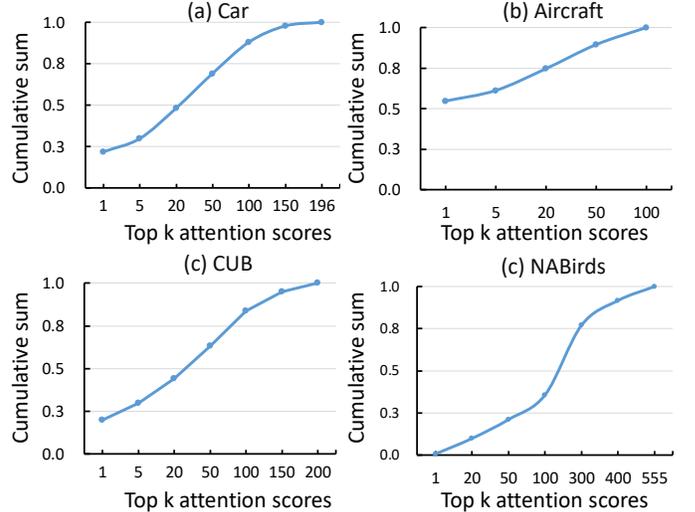


Fig. 5. Cumulative sum of attention scores for top-k highest matching prototypes in memory averaged over all samples in the (a) Car, (b) CUB, (c) Aircraft, and (d) NABirds datasets. Notice that on the right end of each graph, the attention scores of all classes sum to one.

knowledge of each class to generate more expressive representations. The advantage of the mechanism is two-fold, 1) it encourages neural networks to preserve shared semantic information, enabling a comprehensive understanding of fine-grained classes; 2) it allows networks to utilize external knowledge from relevant classes to better represent a given image, resulting in more robust and informative representation. The result is a final feature that captures both instance specific characteristics and the inter-class semantic relationships for the given input. Specifically, for the most indistinguishable classes, CMN will learn to focus more on instance specific characteristics. Thus, the subtle differences essential for classifying these classes can be better represented.

It is worth noting that our method does not require any side information such as structural labels [38], [39] or shared attributes [27]. Instead, our network can learn to adaptively select shared cues across all categories for a given input. Moreover, our method is compatible with several part-localization based methods such as [40], [4], and hence has the potential to further improve recognition.

**Interpretability.** The physical meaning of the proposed class-specific memory module is very clear. The module is used to store and share representative characteristics between classes. Specifically, each row of the module corresponds to the prototypical information of a single class in the form of a moving average. In addition, the attention-based memory reading is based on the relevance between given inputs and prototypes and hence is also interpretable. Note that, in few-shot setting, class prototypes are used as class centers for constructing the nearest-neighbor classifier [35], [36], [37]. Our work is significantly different from them: we use relevant class prototypes to enhance the input features.

**Attention distribution.** We visualize the cumulative sum of attention scores for top-k most similar prototypes in memory averaged over all samples in a dataset. The cumulative sum

TABLE I
STATISTICS OF THREE FINE-GRAINED CLASSIFICATION DATASETS.

| Dataset | #Train Set | #Test Set | #Category |
|---|---|---|---|
| FGVC Aircraft [41] | 6,667 | 3,333 | 100 |
| CUB-200-2011 [42] | 5,994 | 5,794 | 200 |
| Stanford Cars [43] | 8,144 | 8,041 | 196 |
| NABirds [44] | 23,929 | 24,633 | 555 |

TABLE II
COMPARISON OF VARIOUS METHODS IN TERMS OF TOP-1 ACCURACY (%)
ON *rigid* FGVC-AIRCRAFT. * DENOTES OUR RE-IMPLEMENTATION.

| Methods | Backbone | Accuracy (%) |
|---|---|---|
| BilinearCNN [47] | VGG-16 | 84.1 |
| Low-rank B-CNN [10] | VGG-16 | 87.3 |
| Kernel-Pooling [48] | VGG-16 | 86.9 |
| DFL-CNN [24] | VGG-16 | 91.1 |
| MA-CNN [2] | VGG-19 | 89.9 |
| Kernel-Pooling [48] | ResNet-50 | 85.7 |
| FT-ResNet [6]* | ResNet-50 | 90.0 |
| DFL-CNN [24] | ResNet-50 | 91.7 |
| NTS-Net [1] | ResNet-50 | 91.4 |
| DCL [49] | ResNet-50 | 93.0 |
| Cross-X [12] | ResNet-50 | 92.6 |
| LIO [50] | ResNet-50 | 92.7 |
| ACNet [51] | ResNet-50 | 92.4 |
| CMN (ours) | ResNet-50 | 93.8 |

curve directly represents the distribution of attention scores on a specific dataset. We show the cumulative sum curves of four datasets in Figure 5. We observe that approximately 80% of the response feature is made up of only half the prototypes. This is consistent with our intuition that more relevant prototypes contribute more information to augment the input feature.

**Efficiency.** Our method introduces a memory module in the form of a $C \times D$ matrix. It involves an additional matrix multiplication, which is computationally negligible. Thus, it has very similar inference time with baseline even when the number of classes $C$ is large. CMN requires only a single forward pass, while the most localization-based methods [1], [40] require multiple forward passes. Comparing with them, the computational cost of CMN is much less.

## IV. EXPERIMENT

### A. Datasets

We comprehensively evaluate our method on four datasets, including CUB-200-2011 [42], Stanford Cars [43], FGVC Aircraft [41], and NABirds [44], respectively.

**CUB-200-2011** contains 200 birds categories with roughly 30 training images per category. It has 5994 training images and 5794 testing images. **FGVC-Aircraft** is comprised of 10000 images over 100 categories, with a ratio of the training set to testing set of roughly $2:1$. **Stanford Car** is made up of 196 categories of cars, with 8144 examples in training set and 8041 examples in the testing set. **NABirds** is a large scale dataset, consisting of 555 classes. It is made up of $23,929$ training images and $24,633$ testing images.

In the experiment, we follow the same train/test splits in Table I. Our method does not utilize any auxiliary supervision methods such as extra annotations and relies solely on image labels for the learning supervision. We report the top-1 accuracy in all experiments.

### B. Implementation Details

**Implementation Details.** We implement our method in Pytorch [45] and train all models on a single Tesla P-100 GPU. We adopt ResNet-50 [6] pre-trained on ImageNet [46] as our backbone network. The proposed method is trained using the SGD optimizer with momentum 0.9, weight decay 1e-4 and a mini-batch size of 16. We use a learning rate 4e-3 for the backbone and $5\times$ multiplier for the newly added layer. We train the network for 60 epochs and decay the learning rate by 0.1 after 40 epochs. We preprocess images to $448 \times 448$ and use random horizontal flipping with a probability of 0.5 for data augmentation. Moreover, we set coefficient $\beta$ in Eq. 1 as

0.9. The entire network is trained in an end-to-end manner via the supervision of image labels. The network *does not* require any special initialization, multiple training stages, and part or bounding box annotations. We report the top-1 classification accuracy from the last epoch.

### C. Performance Comparison

**Compared methods.** We compare the proposed method with ResNet-50 based methods and also include the best results of VGG based methods. For a fair comparison, we do not include approaches using additional data or annotations. The baseline is fine-tuning from ResNet-50 (FT-ResNet) [6], which is used in NTS-Net [1] and TASN [40].

**Comparison on FGVC-Aircraft.** In Table II, we compare our method (CMN) with the existing methods on Aircraft. CMN achieves 93.8% accuracy, which is $+3.8\%$ higher over baseline (FT-ResNet). The achieved accuracy is highest on the Aircraft. This indicates that the proposed class-specific memory module is beneficial for improving feature representations. Moreover, CMN outperforms LIO [50] and DCL [49] by 1.1% and 0.8%, respectively. In addition, our method surpasses ACNet [51] by 1.4%. Compared with TASN [40] and NTS-Net [1], CMN also gains competitive accuracy, showing its effectiveness for Aircraft classification.

**Comparison on Stanford Cars.** As shown in Table III, CMN gains 2.4% improvement over baseline (FT-ResNet). Compared with the state-of-the-art methods, CMN can achieve competitive accuracy. For example, it is 0.3% and 0.4% higher than ACNet [51] and DCL [24], respectively. Moreover, our methods is also higher than other VGG based approaches, such as MA-CNN [2] and Kernel-Pooling [48]. The above comparisons demonstrate our method is effective for the two rigid datasets with a significant structural variation. It is worthy of noting that CMN *does not* require multiple forward passes like the most localization-based methods [1], [40]. Therefore, we think the proposed method is efficient and effective for fine-grained feature learning.

TABLE III
COMPARISON OF VARIOUS METHODS IN TERMS OF TOP-1 ACCURACY (%)
ON *rigid* **STANFORD CARS**. * DENOTES OUR RE-IMPLEMENTATION.

| Methods | Backbone | Accuracy (%) |
|---|---|---|
| BilinearCNN [47] | VGG-16 | 91.3 |
| Low-rank B-CNN [10] | VGG-16 | 90.9 |
| Kernel-Pooling [48] | VGG-16 | 92.4 |
| DFL-CNN [24] | VGG-16 | 93.3 |
| RA-CNN [3] | VGG-19 | 92.5 |
| MA-CNN [2] | VGG-19 | 92.8 |
| Kernel-Pooling [48] | ResNet-50 | 91.1 |
| MAMC [23] | ResNet-50 | 92.8 |
| FT-ResNet [6]* | ResNet-50 | 92.5 |
| DT-RAM [52] | ResNet-50 | 93.1 |
| DFL-CNN [24] | ResNet-50 | 93.1 |
| NTS-Net [1] | ResNet-50 | 93.9 |
| TASN [40] | ResNet-50 | 93.8 |
| DCL [49] | ResNet-50 | 94.5 |
| Cross-X [12] | ResNet-50 | 94.6 |
| LIO [50] | ResNet-50 | 94.5 |
| ACNet [51] | ResNet-50 | 94.6 |
| CMN (ours) | ResNet-50 | 94.9 |

TABLE IV
COMPARISON OF VARIOUS METHODS IN TERMS OF TOP-1 ACCURACY (%)
ON *non-rigid* **CUB-200-2011**. * DENOTES OUR RE-IMPLEMENTATION.

| Methods | Backbone | Accuracy (%) |
|---|---|---|
| BilinearCNN [47] | VGG-16 | 84.1 |
| Low-rank B-CNN [10] | VGG-16 | 84.2 |
| Kernel-Pooling [48] | VGG-16 | 86.2 |
| DFL-CNN [24] | VGG-16 | 85.8 |
| RA-CNN [3] | VGG-19 | 85.3 |
| MA-CNN [2] | VGG-19 | 86.5 |
| FT-ResNet [6]* | ResNet-50 | 85.8 |
| DFL-CNN [24] | ResNet-50 | 87.4 |
| NTS-Net [1] | ResNet-50 | 87.5 |
| TASN [40] | ResNet-50 | 87.9 |
| DCL [49] | ResNet-50 | 87.8 |
| Cross-X [12] | ResNet-50 | 87.7 |
| LIO [50] | ResNet-50 | 88.0 |
| ACNet [51] | ResNet-50 | 88.1 |
| CMN (ours) | ResNet-50 | 88.2 |

TABLE V
COMPARISON OF VARIOUS METHODS ON *non-rigid* **NABIRDS** DATASET.
HERE, * DENOTES OUR RE-IMPLEMENTATION. OUR METHOD ACHIEVES
COMPETITIVE ACCURACY WITH THE-STATE-OF-THE-ART METHODS.

| Method | Backbone | Accuracy (%) |
|---|---|---|
| BilinearCNN [47] | VGG-16 | 79.4 |
| FT-ResNet [6]* | ResNet-50 | 84.0 |
| MaxEnt [26] | ResNet-50 | 69.2 |
| MaxEnt [26] | DenseNet-161 | 83.0 |
| Cross-X [12] | ResNet-50 | 86.2 |
| CMN (ours) | ResNet-50 | 87.8 |

**Comparison on CUB-200-2011.** We additionally report the results on non-rigid CUB-200-2011 shown in Table III. We show that our CMN achieves 2.4% improvement over baseline (FT-ResNet) and arrives at 88.2%. This result is also comparable with the current state-of-the-art method ACNet [51] (88.2% vs. 88.1%). It is worth noting that CMN outperforms ACNet by +1.4% and +0.3% on Aircraft and Car, respectively. Compared with TASN [40], our method is also competitive (88.2% vs. 87.9%). In addition, TASN and NTS-Net [1] require multiple forward passes to extract region features, making them significantly more computationally prohibitive. CMN requires only a single forward pass through the network, and thus is relatively efficient. Moreover, we think CMN would be compatible with these part-localization based methods, and further achieve improvements. Our method also gains comparable accuracy with Cross-X [12]. Cross-X introduces a cross-layer regularizer to leverage multi-scale features from different layers, which we believe if combined with our method could also result in an improvement in classification accuracy.

**Comparison on NABirds.** The architecture of the purposed method is simple and effective. It can easily be applied to large-scale datasets. We also conduct experiments on another non-rigid dataset NABirds and report results in Table V. In this dataset, we observe that our method also achieves +3.8% improvement over baseline and arrives at 87.8 %. This improvement also validate the effectiveness of our method. Moreover, CMN is competitive with the recent state-of-the-art method Cross-X [12]. More specific, our method is 1.6% higher than Cross-X in the same settings. We think our method considers shared patterns across classes (inter-class similarities) that is a characteristic of the fine-grained classification, and thus achieves improvements over baselines.

### D. Component Analysis

**Effect of attention-based memory reading.** Given an input feature, we adopt an attention-based method of retrieving relevant prototypes (as detailed in Section III-A2). To validate its effectiveness, we introduce two variants for comparison, 1) equal selection, where all attention scores are equal to $1/C$ ($C$ is the number of categories); 2) predicted selection [32], which uses a fully connected layer followed by a softmax to predict attention scores for the given input feature. We compare these two variants and our method in Figure 8. They both improve the accuracy compared with our baseline. This demonstrates that leveraging class prototypes is beneficial for fine-grained feature learning. Meanwhile, equal selection delivers worse performance than predicted selection (87.1% vs. 87.5%), indicating that the retrieved prototypes should be dependent on the input feature. Our attention-based method outperforms equal section and predicted selection by 1.1% and 0.7%, respectively. This is because the attention-based method is based on Euclidean distance, which is a direct and effective way to represent the feature correlation.

**Importance of learning in the class-specific module.** Class-specific memory module records the moving average of each class and is updated during training. To validate the importance of learning in the module, we test a random module for comparative purposes, *i.e.*, we randomly initialize the memory matrix and fix it during training. In Table VI, we
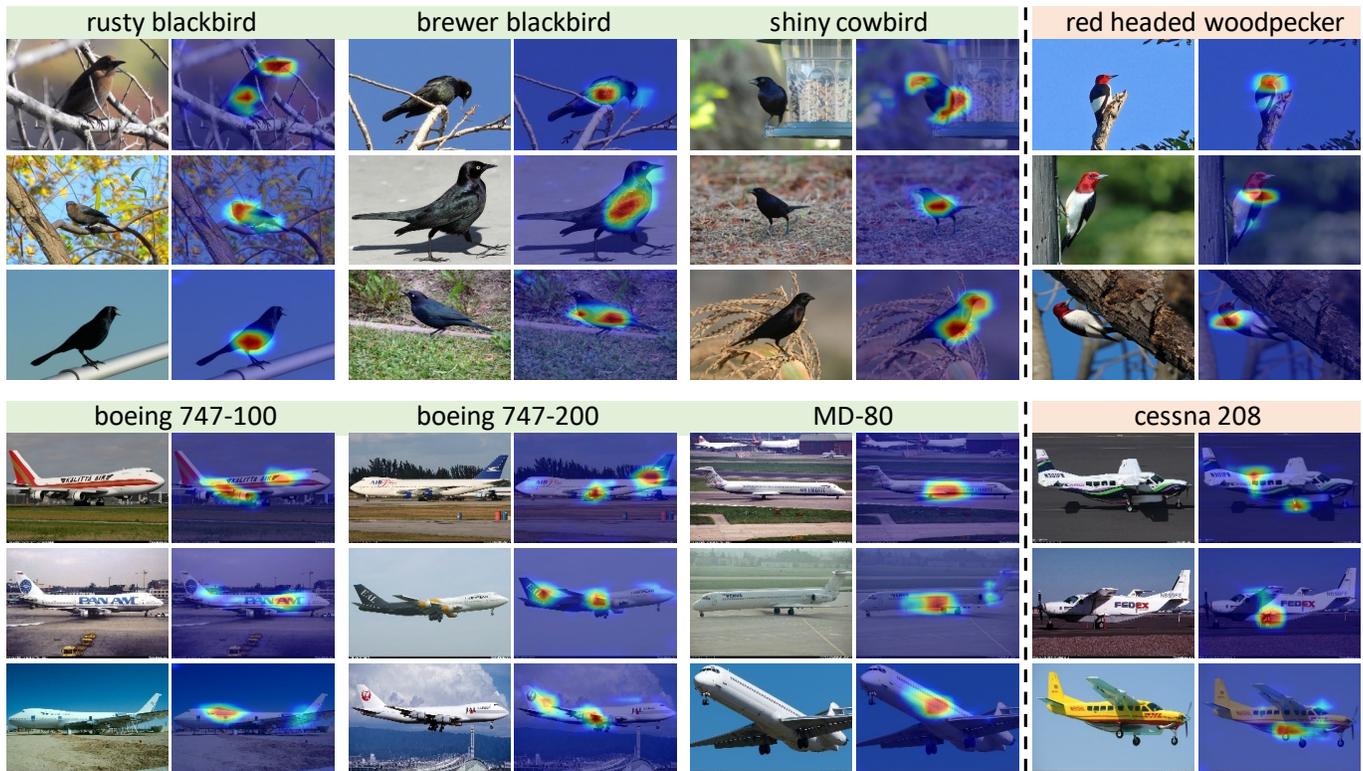
Fig. 6. Visualization of the learned visual cues on the *rigid* dataset Aircraft and the *non-rigid* dataset CUB. Based on the attention scores computed by our method, we select images from three visually similar classes (the first three columns) and one dissimilar class (the last column) for each dataset. The similar classes share some visual patterns, *e.g.*, "breast" and "bill" for birds, and "under-wing fuselage", "tail", and "upperdeck" for aircraft. Each class contains unique and subtle visual cues. Best viewed in color.
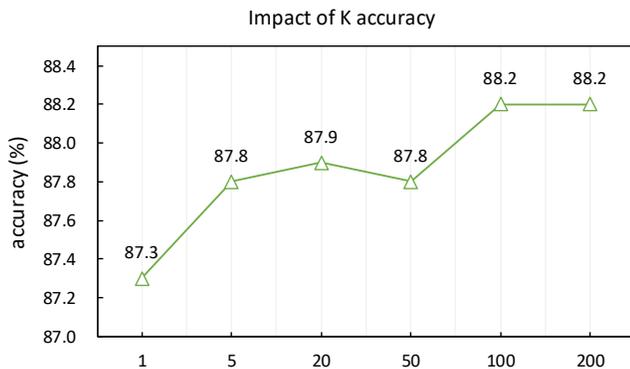


Fig. 7. Importance of prototype diversity. We show the accuracy (%) of using only the prototypes with top-K attention scores to form a memory response feature. When K=200, the response is given access to all prototypes. The results are on the CUB-200-2011 dataset.
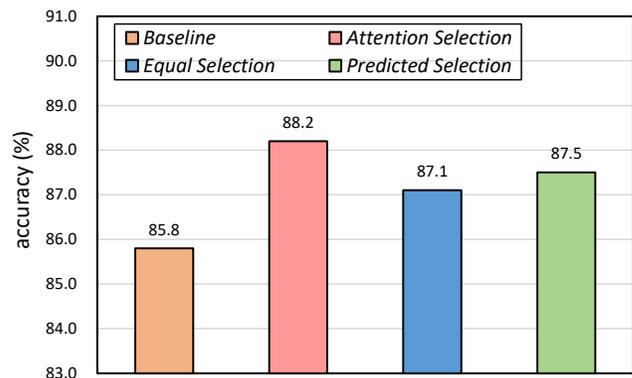


Fig. 8. Comparison of three attention score mechanisms for memory reading. We test attention-based selection *i.e.* our technique, equal weighting for all prototypes, and predicting scores based on the given image via a fully-connected layer. The results are on the CUB-200-2011 dataset. The proposed attention-based selection is the best way to combine class prototypes.

compare the two modules on Aircraft and CUB. We observe that the random module is on par with our baseline. This indicates that 1) naively increasing network capacity is not beneficial and 2) that the module should record useful and meaningful information and be updated during network learning.

**Importance of prototype diversity.** We leverage all prototypes to compute a response feature (Eq. 2). In Figure 8, we show the effect on accuracy of using only the prototypes with the top-K attention scores on CUB. We observe that allowing the response to rely on more than just the most

similar prototypes improves accuracy. Moreover, we observe that the accuracy does improve using more than 100 prototypes. More specifically, the accuracy does not improve whrn using more than 100 prototypes. This is because the response feature is most formed by the half relevant prototypes, which is consistent with our observation in attention distribution (Figure 5): the most relevant prototypes contribute more information to enhance the input feature.

TABLE VI
IMPORTANCE OF LEARNING IN THE CLASS-SPECIFIC MODULE. THE
CLASS-SPECIFIC MODULE DENOTES RECORDING THE MOVING AVERAGE OF
EACH CLASS. THE RANDOM MODULE MEANS THAT MEMORY MATRIX IS
RANDOMLY INITIALIZED AND FIXED DURING TRAINING.

| Methods | FGVC Aircraft | CUB-200-2011 |
|---|---|---|
| Baseline | 90.0 | 85.8 |
| Random module | 90.0 | 85.8 |
| Class-specific module | 93.8 | 88.2 |

TABLE VII
EFFECT OF MEMORY MODULE ON CIFAR-100. THE MEMORY MODULE
DOES NOT BRING ANY IMPROVEMENT AND IN FACT, SEEMS TO HINDER
PERFORMANCE. THIS SUGGESTS THAT OUR METHOD IS SPECIFICALLY
TAILORED FOR FINE-GRAINED RECOGNITION, WHERE IMAGES ARE
SEMANTICALLY AND VISUALLY SIMILAR.

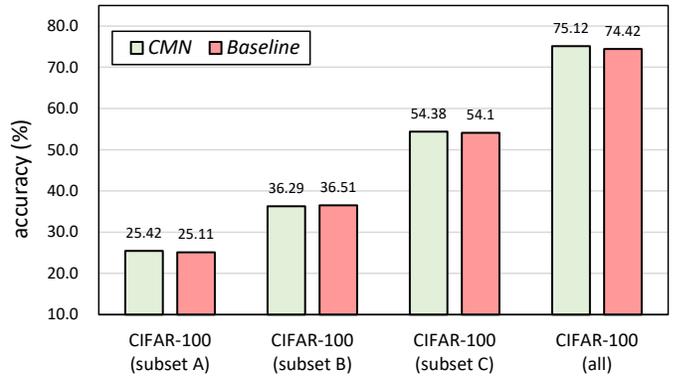| Networks | Memory module | Top-1 Accuracy (%) |
|---|---|---|
| ResNet-18 |  | 75.42 |
| ResNet-18 | ✓ | 75.12 |
| ResNet-34 |  | 75.54 |
| ResNet-34 | ✓ | 75.51 |
| ResNet-50 |  | 77.10 |
| ResNet-50 | ✓ | 76.70 |



Fig. 9. Comparison between ResNet-18 and CMN under different training sets. CIFAR-100 contains all 50,000 training images; CIFAR-100 (subset A) has 5,000 training images, CIFAR-100 (subset B) has 10,000 training images, CIFAR-100 (subset C) has 25,000 training images. We find that both methods are on par with each other under four training sets. This indicates CMN is not suitable for general datasets where the inter-class similarity is relatively small.

**Feature visualization.** We average the attention scores for the images of each class. For a particular class, a higher average attention score indicates that the corresponding class is semantically similar. In Figure 6, we visualize the feature maps of last convolutional layer. We select the three most relevant classes (the first three columns) as defined by our method and one irrelevant class (the last column). We observe that images of correlated classes tend to focus of similar patterns, *e.g.*, *Shiny Cowbird*, *Brewer Blackbird*, and *Rusty Blackbird* have the same visual evidence of breast pattern. This shared pattern can effectively distinguish them from irrelevant *Red-headed Woodpecker*. Moreover, each class has it own distinctive cues, *e.g.*, *Boeing 747-200* has a discriminative evidence on "tail" part, and *Cessna 208* has an unique pattern on "wheel" part.

*E. Further Evaluation.*

We also conduct experiments on more general classification dataset CIFAR-100 [53], where the inter-class similarity is relatively small. In Table VII, we validate the effect of the memory module on several networks. We observe that our CMN does not bring any improvement on this dataset. This is because the classes in CIFAR-100 have low class-class similarity. Furthermore, our method could add some unnecessary class-class cues to the input feature, the learning difficulty may be increased, thus having slightly lower accuracy.

In fine-grained classification datasets, the number of training images per class is relatively small. In comparison, CIFAR-100 contains adequate training data for each class. In addition, the class-specific prototypes have shown its effectiveness for few-shot learning [35], [36], [37]. This motivates us to further study the working mechanism of CMN under the few-example setting. Specifically, we study whether CMN mainly helps when training data per class is scarce on CIFAR-100.

To validate this, we reduce the number of training images in CIFAR-100 and report the experimental results in Figure 9. In practice, we create three subsets, 1) CIFAR-100 (subset A) has 5,000 training images; 2) CIFAR-100 (subset B) has 10,000 training images; 3) CIFAR-100 (subset C) has 25,000 training images. As shown in Figure 9, we observe that both CMN and original ResNet-18 achieve almost the same top-1 accuracy under the different CIFAR-100 subsets. This also indicates that CMN might not suitable for general classification dataset that contains relatively small inter-class similarity. Thus, we believe the advantage of CMN lies in that it has ability to leverage the shared semantic information across classes to well represent the input image. In summary, CMN is specifically tailored for fine-grained classification, where images are both semantically and visually similar.

## V. CONCLUSIONS

We introduce an efficient, effective, and interpretable class-specific memory module for fine-grained visual classification. By storing the prototypical information of all classes, the memory module allows the network to share semantic information between classes. We use an attention-based method to retrieve relevant class prototypes for a given input and generate a specially tailored response feature to augment the original feature. We integrate the memory module with a convolutional neural network, yielding a Categorical Memory Network (CMN). Our method leverages relevant prototypes to produce representations that express inter-class semantics while preserving instance-specific discriminative cues. Moreover, encouraging the network to represent inter-class similarities leads to more robust fine-grained representations.

Our method significantly increases accuracy at effectively no additional computational cost, achieving accuracy competitive with state-of-the-art methods on four fine-grained classification benchmarks. In the future, we plan to explore unsupervised categorical memory networks where the aggregation of prototypes may allow for the discovery of unsupervised sub-categories.

## REFERENCES

[1] Z. Yang, T. Luo, D. Wang, Z. Hu, J. Gao, and L. Wang, "Learning to navigate for fine-grained classification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 420–435.

[2] H. Zheng, J. Fu, T. Mei, and J. Luo, "Learning multi-attention convolutional neural network for fine-grained image recognition," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 5209–5217.

[3] J. Fu, H. Zheng, and T. Mei, "Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4438–4446.

[4] A. Recasens, P. Kellnhofer, S. Stent, W. Matusik, and A. Torralba, "Learning to zoom: a saliency-based sampling layer for neural networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 51–66.

[5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.

[6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[7] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.

[8] T.-Y. Lin, A. RoyChowdhury, and S. Maji, "Bilinear cnn models for fine-grained visual recognition," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1449–1457.

[9] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell, "Compact bilinear pooling," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 317–326.

[10] S. Kong and C. Fowlkes, "Low-rank bilinear pooling for fine-grained classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 365–374.

[11] C. Yu, X. Zhao, Q. Zheng, P. Zhang, and X. You, "Hierarchical bilinear pooling for fine-grained visual recognition," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 574–589.

[12] W. Luo, X. Yang, X. Mo, Y. Lu, L. S. Davis, J. Li, J. Yang, and S.-N. Lim, "Cross-x learning for fine-grained visual categorization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8242–8251.

[13] D. Chang, Y. Ding, J. Xie, A. K. Bhunia, X. Li, Z. Ma, M. Wu, J. Guo, and Y.-Z. Song, "The devil is in the channels: Mutual-channel loss for fine-grained image classification," *IEEE Transactions on Image Processing*, vol. 29, pp. 4683–4695, 2020.

[14] N. Zhang, J. Donahue, R. Girshick, and T. Darrell, "Part-based r-cnns for fine-grained category detection," in *European conference on computer vision*. Springer, 2014, pp. 834–849.

[15] D. Lin, X. Shen, C. Lu, and J. Jia, "Deep lac: Deep localization, alignment and classification for fine-grained recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1666–1674.

[16] S. Huang, Z. Xu, D. Tao, and Y. Zhang, "Part-stacked cnn for fine-grained visual categorization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1173–1182.

[17] H. Zhang, T. Xu, M. Elhoseiny, X. Huang, S. Zhang, A. Elgammal, and D. Metaxas, "Spda-cnn: Unifying semantic part detection and abstraction for fine-grained recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1143–1152.

[18] C. Goring, E. Rodner, A. Freytag, and J. Denzler, "Nonparametric part transfer for fine-grained recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 2489–2496.

[19] X. He, Y. Peng, and J. Zhao, "Which and how many regions to gaze: Focus discriminative regions for fine-grained visual categorization," *International Journal of Computer Vision*, pp. 1–21, 2019.

[20] Y. Peng, X. He, and J. Zhao, "Object-part attention model for fine-grained image classification," *IEEE Transactions on Image Processing*, vol. 27, no. 3, pp. 1487–1500, 2017.

[21] 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 Transactions on Image Processing*, vol. 25, no. 4, pp. 1713–1725, 2016.

[22] P. Sermanet, A. Frome, and E. Real, "Attention for fine-grained categorization," *arXiv preprint arXiv:1412.7054*, 2014.

[23] M. Sun, Y. Yuan, F. Zhou, and E. Ding, "Multi-attention multi-class constraint for fine-grained image recognition," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 805–821.

[24] Y. Wang, V. I. Morariu, and L. S. Davis, "Learning a discriminative filter bank within a cnn for fine-grained recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4148–4157.

[25] A. Dubey, O. Gupta, P. Guo, R. Raskar, R. Farrell, and N. Naik, "Pairwise confusion for fine-grained visual classification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 70–86.

[26] A. Dubey, O. Gupta, R. Raskar, and N. Naik, "Maximum-entropy fine grained classification," in *Advances in Neural Information Processing Systems*, 2018, pp. 637–647.

[27] F. Zhou and Y. Lin, "Fine-grained image classification by exploring bipartite-graph labels," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1124–1133.

[28] J. Weston, S. Chopra, and A. Bordes, "Memory networks," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.

[29] S. Sukhbaatar, J. Weston, R. Fergus *et al.*, "End-to-end memory networks," in *Advances in neural information processing systems*, 2015, pp. 2440–2448.

[30] A. Graves, G. Wayne, M. Reynolds, T. Harley, I. Danihelka, A. Grabska-Barwińska, S. G. Colmenarejo, E. Grefenstette, T. Ramalho, J. Agapiou *et al.*, "Hybrid computing using a neural network with dynamic external memory," *Nature*, vol. 538, no. 7626, p. 471, 2016.

[31] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, "Meta-learning with memory-augmented neural networks," in *International conference on machine learning*, 2016, pp. 1842–1850.

[32] C. Li, J. Zhu, and B. Zhang, "Learning to generate with memory," in *International Conference on Machine Learning*, 2016, pp. 1177–1186.

[33] Y. Kim, M. Kim, and G. Kim, "Memorization precedes generation: Learning unsupervised gans with memory networks," *arXiv preprint arXiv:1803.01500*, 2018.

[34] A. Graves, G. Wayne, and I. Danihelka, "Neural turing machines," *arXiv preprint arXiv:1410.5401*, 2014.

[35] K. Hsu, S. Levine, and C. Finn, "Unsupervised learning via meta-learning," *arXiv preprint arXiv:1810.02334*, 2018.

[36] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," in *Advances in neural information processing systems*, 2017, pp. 4077–4087.

[37] D. Wertheimer and B. Hariharan, "Few-shot learning with localization in realistic settings," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6558–6567.

[38] S. Xie, T. Yang, X. Wang, and Y. Lin, "Hyper-class augmented and regularized deep learning for fine-grained image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2645–2654.

[39] D. Wang, Z. Shen, J. Shao, W. Zhang, X. Xue, and Z. Zhang, "Multiple granularity descriptors for fine-grained categorization," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2399–2406.

[40] H. Zheng, J. Fu, Z.-J. Zha, and J. Luo, "Looking for the devil in the details: Learning trilinear attention sampling network for fine-grained image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5012–5021.

[41] S. Maji, E. Rahtu, J. Kannala, M. Blaschko, and A. Vedaldi, "Fine-grained visual classification of aircraft," *arXiv preprint arXiv:1306.5151*, 2013.

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

[43] J. Krause, M. Stark, J. Deng, and L. Fei-Fei, "3d object representations for fine-grained categorization," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2013, pp. 554–561.

[44] G. Van Horn, S. Branson, R. Farrell, S. Haber, J. Barry, P. Ipeirotis, P. Perona, and S. Belongie, "Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 595–604.

[45] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.

[46] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[47] T. Lin, A. Roy Chowdhury, and S. Maji, "Bilinear convolutional neural networks for fine-grained visual recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 6, pp. 1309–1322, 2018.

[48] Y. Cui, F. Zhou, J. Wang, X. Liu, Y. Lin, and S. Belongie, "Kernel pooling for convolutional neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2921–2930.

[49] Y. Chen, Y. Bai, W. Zhang, and T. Mei, "Destruction and construction learning for fine-grained image recognition," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[50] M. Zhou, Y. Bai, W. Zhang, T. Zhao, and T. Mei, "Look-into-object: Self-supervised structure modeling for object recognition," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[51] R. Ji, L. Wen, L. Zhang, D. Du, Y. Wu, C. Zhao, X. Liu, and F. Huang, "Attention convolutional binary neural tree for fine-grained visual categorization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 468–10 477.

[52] Z. Li, Y. Yang, X. Liu, F. Zhou, S. Wen, and W. Xu, "Dynamic computational time for visual attention," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1199–1209.

[53] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," Citeseer, Tech. Rep., 2009.