

# Fully Convolutional Attention Networks for Fine-Grained Recognition

Xiao Liu, Tian Xia, Jiang Wang, Yi Yang, Feng Zhou and Yuanqing Lin  
Baidu Research

{liuxiao12,xiatian,wangjiang03,yangyi05, zhoufeng09, linyuanqing}@baidu.com

## Abstract

*Fine-grained recognition is challenging due to its subtle local inter-class differences versus large intra-class variations such as poses. A key to address this problem is to localize discriminative parts to extract pose-invariant features. However, ground-truth part annotations can be expensive to acquire. Moreover, it is hard to define parts for many fine-grained classes. This work introduces **Fully Convolutional Attention Networks (FCANs)**, a reinforcement learning framework to optimally glimpse local discriminative regions adaptive to different fine-grained domains. Compared to previous methods, our approach enjoys three advantages: 1) the weakly-supervised reinforcement learning procedure requires no expensive part annotations; 2) the fully-convolutional architecture speeds up both training and testing; 3) the greedy reward strategy accelerates the convergence of the learning. We demonstrate the effectiveness of our method with extensive experiments on four challenging fine-grained benchmark datasets, including CUB-200-2011, Stanford Dogs, Stanford Cars and Food-101.*

## 1. Introduction

Fine-grained recognition refers to the task of distinguishing sub-ordinate categories, such as bird species [1], dog breeds [2], car models [3], flower categories [4], food dishes [5], etc. With the great potential in rivaling human experts, it has shown tremendous applications in real world ranging from e-commerce [6, 7] to education [8, 9]. Although great success has been achieved for basic-level recognition in the last few years [10, 11, 12, 13], fine-grained recognition still faces two challenges. First, it is more difficult and time-consuming to gather a large amount of labeled fine-grained data because it calls for experts with specialized domain knowledge. In addition, the difference between fine-grained classes is very subtle. The most discriminative features are often not based on the global shape or appearance variation but contained in the mis-alignment of local parts or patterns. For instance, as shown in Fig. 1,

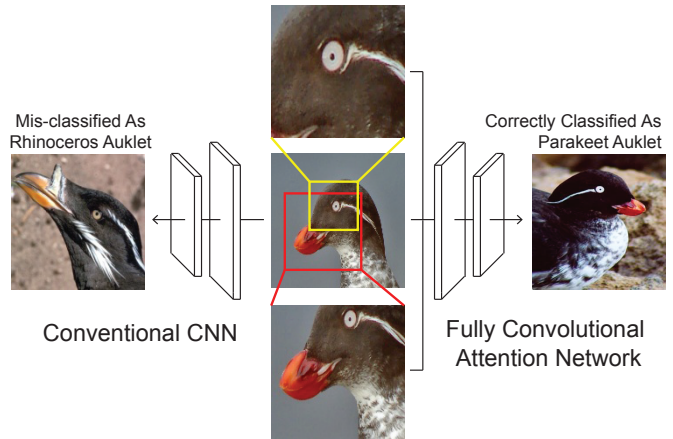


Figure 1. Conventional CNN approach (left) finds difficulty in differentiating similar fine-grained categories with subtle local variations (eg., Rhinoceros Auklet against Parakeet Auklet). In contrast, our proposed fully convolutional attention networks (right) is able to automatically and efficiently localize parts (eg., bird’s eye and beak) given only weakly supervised fine-grained class labels.

the eye texture and beak shape are crucial to differentiate between Parakeet Auklet and Rhinoceros Auklet.

To that end, the main body of previous research has focused on devising more discriminative features by detecting and aligning object parts. Nevertheless, most conventional methods [14, 15] utilize manually defined parts to localize the regions, such as “the head of a bird”, for fine-grained recognition. Relying on manually defined parts has several drawbacks: 1) The precise part annotations are usually expensive to acquire. 2) The strongly supervised part-based model might fail if some parts are occluded. 3) For some fine-grained categories, it is very difficult to manually define parts for them. For example, it is very difficult to define parts for food recognition, as suggested in Fig. 2. 4) Most importantly, there is no clue that manually defined parts are optimal for all fine-grained recognition tasks.

To overcome these problems, we propose a visual attention framework called *Fully Convolutional Attention Networks (FCANs)* for fine-grained recognition without part



Figure 2. Fine-grained recognition often involves part localization, i.e. (a) localizing head and breast to distinguish birds, and (b) localizing brand to classify car makes. It is relatively easy to define parts for structured objects like birds and cars. However, it is hard to define rigid parts for unstructured classes, such as (c) food. This may be solved by attention models.

annotation. Given only image label, our framework utilizes reinforcement learning to simultaneously localize object parts and classify the object within the scene. Intuitively, the framework simulates human visual system that proceeds object recognition via a series of *glimpse* on object parts. At each glimpse, it strives to find the most discriminative location that can differentiate object’s category given the previous observations. Similar to previous visual attention models [16, 17], we employ the REINFORCE algorithm during training [18], where the action is the location of each glimpse, the state is the image and the locations of the previous glimpses, and the reward measures the classification correctness. The whole framework can be trained only by an image classification loss, thus requiring no manual part annotations. The visual attention approach is demonstrated to perform well on fine-grained recognition without requiring manually labeled object parts [17].

Compared to the previous reinforcement learning-based visual attention frameworks [16, 17], the FCANs enjoy better computational efficiency as well as higher classification accuracy in fine-grained recognition. More concretely, our proposed framework improves the attention models in three ways:

- **Computational Efficiency:** The previous frameworks run a convolutional neural network individually on each image crop, which is computationally expensive during both training and testing. In contrast, our method re-uses the same feature maps (computed by a fully convolutional neural network [12, 11]) during each glimpse in a way similar to Fast-RCNN [19]. This makes training and prediction computationally more

efficient because of the fully convolutional neural network architecture and feature sharing technique.

- **Multiple Part Localization:** During testing, our model is able to simultaneously locate multiple parts of adaptive sizes, while the previous frameworks [16, 17] generally only locate one part at each iteration.
- **Faster Training Convergence:** Instead of assigning a delayed reward at the end of attention iterations as previous methods [16, 17], we apply a new greedy reward strategy at every step of attention, which is crucial to both the convergence speed of training and the accuracy of prediction.

As a result, our proposed approach improves the recognition accuracy over previous reinforcement learning based methods [16, 17] while being computationally more efficient. Our method also achieve competitive results against other state-of-the-art methods on multiple fine-grained datasets.

## 2. Related Work

Fine-grained recognition has been extensively studied in recent years [5, 9, 20, 21, 22, 3, 2, 14, 4]. We review the three most relevant directions in this section.

### 2.1. Representation Learning

Since the seminal work of AlexNet [10], we are witnessing a fast-pacing transition from hand-crafted feature to end-to-end convolutional neural networks in representation learning [11, 12, 13]. Most of the current state-of-the-art fine-grained recognition algorithms are also based on deep CNN representation to distinguish the subtle difference [23, 24, 25]. Branson *et al.* [15] claim that integrating lower-level layer and higher-level layer features learns more discriminative representation for fine-grained recognition. Lin *et al.* [25] propose a bilinear architecture to model local pairwise feature interactions for fine-grained recognition, where convolutional features from two models are combined in a translation invariant manner. Qian *et al.* [26] propose a multi-stage metric learning framework to learn a distance metric that pulls data points of the same class close and pushes data points from different classes far apart. Wang *et al.* [27] combine saliency-aware object detection approach and object-centric sampling scheme to extract more robust and discriminative features for large-scale fine-grained car classification. In parallel to these efforts, our method combines representation learning with part detection in a unified framework.

### 2.2. Part Models

Since 70’s, early cognitive research study [28] has shown that subordinate-level recognition is based on comparing

the appearance details of object parts. Drawing inspiration from this fact, various pose normalization methods [29, 30, 31, 14, 32] have been proposed to focus on the important regions. However, these methods are strongly supervised ones, heavily relying on manually pre-defined parts modeled by Poselet [33] or DPM [34]. Due to this limit, most of recent efforts were spent on how to automatically discover critical parts in a weaker setting. For instance, Berg *et al.* [35] use data mining techniques to learn a set of intermediate features that can differentiate two classes based on the appearance of a particular part. Yang *et al.* [36] propose a template model to discover the common geometric patterns of object parts and the co-occurrence statistics of the patterns. Similarly, Gavves *et al.* [37] and Chai *et al.* [38] segment images and align the image segments in an unsupervised fashion. The aligned image segments are utilized for feature extraction separately. Recently, Simon and Rodner [39] propose neural activation constellations, an approach that is able to learn part models in an unsupervised manner. Compared to our method, however, these methods require tedious and ad-hoc tuning of individual components.

### 2.3. Attention Models

One of the main drawbacks of part-based models is the need for a strong motivation in part definition (either by hand or by data-mining method), which may lack for many non-structured objects such as food dishes [40]. On the other hand, several works introduce attention-based models for task-driven object/part localization. For instance, Mnih *et al.* [16] present a recurrent neural network model for object detection by adaptively selecting a sequence of attention regions and extract appearance representations in these regions. Since this model is non-differentiable, it is trained with reinforcement learning technique to learn task-specific policies. Ba *et al.* [41] extend [16] and successfully achieve good results on a more challenging multi-digit recognition task. Despite the remarkable contributions in theory, the recurrent attention models still suffer from several drawbacks in practice. First, they only result in small performance improvement. For instance, Sermanet *et al.* [17] further extend [41] to fine-grained recognition but only achieve 76.8% mean accuracy percentage (with 3 glimpses) on Stanford Dogs dataset while the result of GoogLeNet baseline [12] is 75.5%. Second, the computational burden is high. Calculating features at each glimpse in [17] requires forwarding GoogLeNet three times, leading to very slow training and testing. An exceptional work is Spatial Transformer Networks [42], which build on a differentiable attention mechanism that does not need reinforcement learning for training. As an alternative approach, we show reinforcement learning can still be effective and efficient in improving fine-grained recognition.

## 3. Fully Convolutional Attention Networks

Fig. 3 illustrates the architecture of the Fully Convolutional Attention Networks (FCANs) with three main components: the feature network, the attention network, and the classification network.

**Feature Map Extraction:** The feature network contains a fully convolutional network that extracts features from the input image and its subsequent attention crops. These feature maps are shared for both part attention and fine-grained classification. During experiment, we adopt one of the popular CNN architectures (e.g., VGG-16 [11], GoogLeNet [12] or ResNet [13]) as the basis fully convolutional network, pre-trained on ImageNet dataset [27] and fine-tuned on the target fine-grained dataset. During testing, the image and all attention crops are resized to a canonical size before feature extraction, similar to [16]. Hence the amount of computation it performs can be controlled independently of the input image size.

During training, although cropping local image regions can achieve good performance, it requires us to perform multiple forward and backward passes of a deep convolutional network in one batch, where the time complexity for feature extraction depends on the number of parts and number of attention regions sampled for each part. In practice, this is too time-consuming. Thus, we extract feature maps from the original image at multiple scales and re-use them across all time steps. The features for each part is obtained by selecting the corresponding region in the convolutional feature maps, so that the receptive field of the selected region is the same as the size of the part. As a result, we only need to run the forward pass once in one training batch.

**Fully Convolutional Part Attention:** The attention network localizes multiple parts by generating multiple part score maps from the basis convolutional feature maps. Each score map is generated using two stacked convolutional layers and one spatial softmax layer. The first convolutional layer uses  $64 \ 3 \times 3$  kernels, and the second one uses one  $3 \times 3$  kernels to output a single-channel confidence map. The spatial softmax layer converts the confidence map into probability. During testing, the model selects the attention region with the highest probability as the part location. During training, the model samples attention regions multiple times according to the probability map. The same process is applied for a fixed number of time steps for multiple part locations. Each time step generates the location for a particular part. We will detail this step in the following sections.

**Fine-Grained Classification:** The classification network contains a convolutional network for each part as well as the whole image. The classification network for each part is a fully convolutional layer followed by a softmax layer. Different parts might have different sizes, and a local image region is cropped around each part location according to its size. The final prediction score is the average of all the

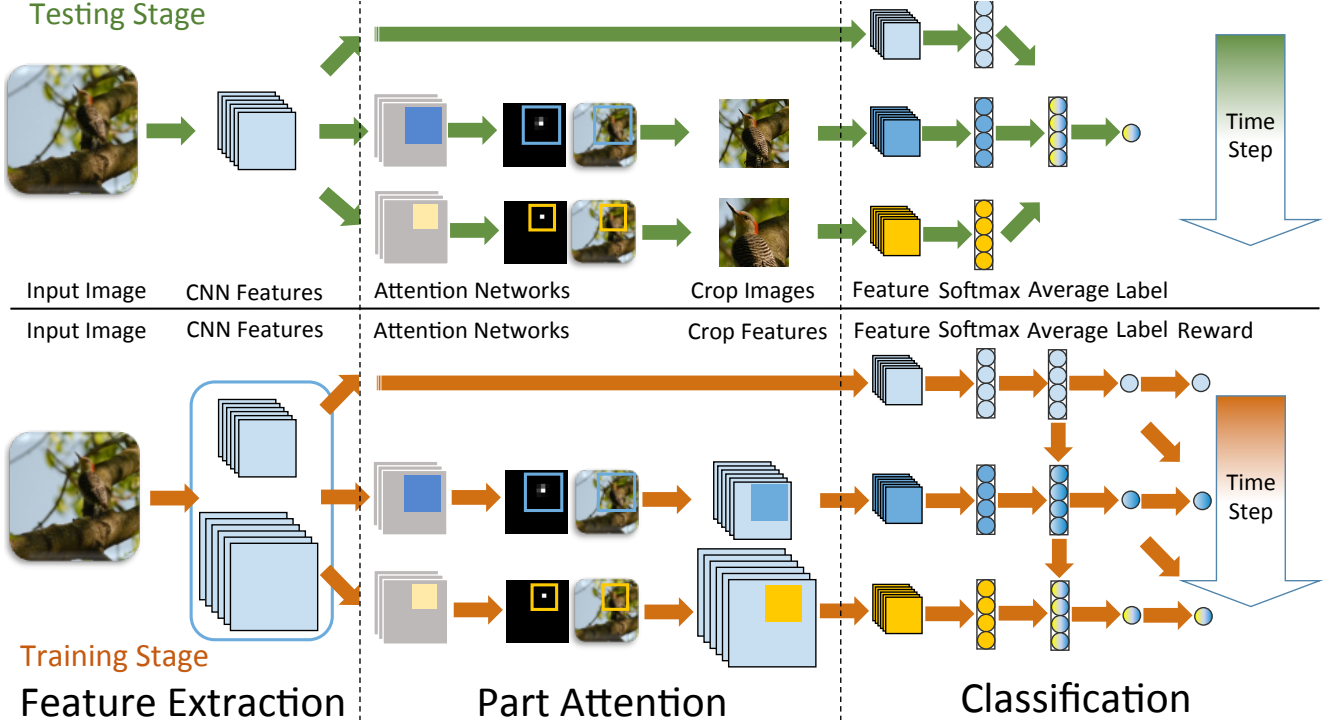


Figure 3. The architecture of our FCANs framework. In this example, the attention network finds two parts of different sizes (the blue region and the yellow region). The upper part shows the architecture for testing, and the lower part shows the architecture for training. During testing, we crop all corresponding part patches from the high resolution image for classification. During training, we re-use the convolutional features in the attention networks for classification. Note that during testing, we can compute all part attentions simultaneously, which makes the model computationally more efficient than traditional recurrent attention models.

prediction scores from the individual classifiers. In order to discriminate the subtle visual differences, each local image region is cropped at high resolution.

### 3.1. Model

The entire attention problem is formulated into a Markov Decision Process (MDP). During each time step of MDP, the FCANs work as an agent to perform an action based on the observation and receives a reward. In our work, the action corresponds to the location of the attention region, the observation is the input image and the crops of the attention regions and the reward measures the quality of the classification using the attention region. The target of our learning is to find the optimal decision policy to generate actions from observations, characterized by the parameters of the FCANs, to maximize the sum expected reward across all time steps.

We define the input image as  $x$  and the feature network (parameterized by  $\theta_f$ ) computes the feature maps as  $\phi(x, \theta_f)$ . The attention network outputs  $T$  attention locations  $\{l^1, \dots, l^T\}$  with each location  $l^t \sim \pi(l^t | \phi, \theta_l^t)$ , where  $\pi$  is the policy for attention selection, parameterized by  $\theta_l = \{\theta_l^t\}_{t=1 \dots T}$ . At time step  $t$ , the classification component crops an image region at location  $l^t$ , extracts a new

feature  $\phi(l^t)$  and predicts classification score  $s_t$  with the classification network (parameterized by  $\theta_c = \{\theta_c^t\}_{t=1 \dots T}$ ). It then computes the final classification score  $S_t$  as the average of all prediction scores until time  $t$

$$S_t = \frac{1}{t} \sum_{\tau=1}^t s_\tau(\phi(l^\tau), \theta_c^\tau) \quad (1)$$

Note that in FCANs, both  $\theta_l$  and  $\theta_c$  have different sets of parameters  $\{\theta_l^t, \theta_c^t\}$  at different time steps. Only the parameters in the feature network  $\theta_f$  are shared across all time steps. This is different from the original recurrent attention models [16, 41] where all parameters are shared across multiple time steps. The reward  $r^t$  for the  $t$ -th step measures how the output of  $S_t$  matches the ground truth label  $y$ .

### 3.2. Training

Since there are no ground-truth annotations to indicate where to select attention regions and each attention is a non-differentiable function, we adopt reinforcement learning to learn the network parameters.

Given a set of training images with ground truth labels  $(x_n, y_n)_{n=1 \dots N}$ , we jointly optimize the three components

to maximize the following objective function:

$$\max_{\theta} J(\theta) = \max_{\theta_f, \theta_l, \theta_c} R(\theta_f, \theta_l) - L(\theta_f, \theta_c) \quad (2)$$

where  $\theta = \{\theta_f, \theta_l, \theta_c\}$  are the parameters of the feature networks, the attention networks and the classification networks respectively.

$$L(\theta_f, \theta_c) = \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T L_n^t(x_n, y_n, \theta_f, \theta_c) \quad (3)$$

is the average cross-entropy classification loss over  $N$  training samples and  $T$  time steps.

$$R(\theta_f, \theta_l) = \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T \mathbb{E}_{\theta} [r_n^t] \quad (4)$$

is the average expected reward over  $N$  training samples and  $T$  time steps.

$$\mathbb{E}_{\theta} [r_n^t] = \sum_{l_n^t} \pi(l_n^t | x_n, \theta_f, \theta_l^t) r_n^t \quad (5)$$

is the expected reward of the  $t$ -th selected attention region from the  $n$ -th sample.  $\theta_l^t$  is the parameters of the  $t$ -th attention network,  $\pi(l_n^t | x_n, \theta_f, \theta_l^t) = \pi(l_n^t | \phi(x_n), \theta_l^t)$  is the probability of selecting  $l_n^t$  as the attention region. The reward function  $r_n^t$  is crucial for developing an efficient learning algorithm. We describe the design of the reward function in the following section.

### 3.3. Reward Strategy

A straightforward reward strategy is to measure the quality of the attention region selection policy as a whole using the final classification result, i.e.,  $r_n^t = 1$  if  $t = T$  and  $y_n = \arg \max_y P(y | S_n^T)$ , and 0 otherwise. Although MDP with such a reward strategy can learn in a recurrent way [16], it confuses the effect of the selected regions in different time steps, and it might lead to the problem of convergence difficulty.

We consider an alternative reward strategy, namely *greedy reward*:

$$r_n^t = \begin{cases} 1 & t = 1 \wedge y_n = \arg \max_y P(y | S_n^1) \\ 1 & t > 1 \wedge y_n = \arg \max_y P(y | S_n^t) \wedge L_n^t < L_n^{t-1} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

where  $L_n^t$  is the classification loss for the  $n$ -th sample at  $t$ -th step. If the image is classified correctly in the first step, the attention network immediately receives a reward. In other time steps, we reward the corresponding attention network only if the image is classified correctly and the classification loss decreases with regards to the last time step. Otherwise, the attention network receives zero reward. Since the attention network immediately receives a reward when an image is correctly classified with the current attention region, the convergence of training is much easier.

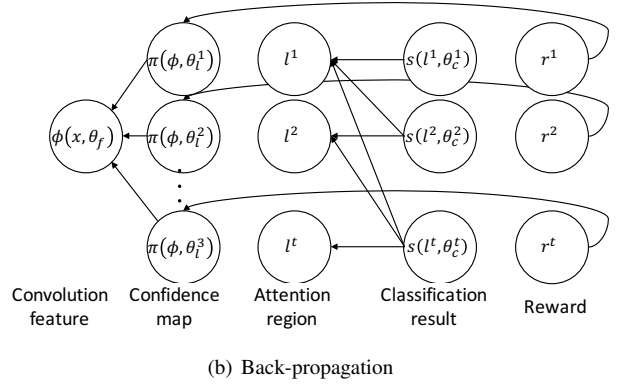
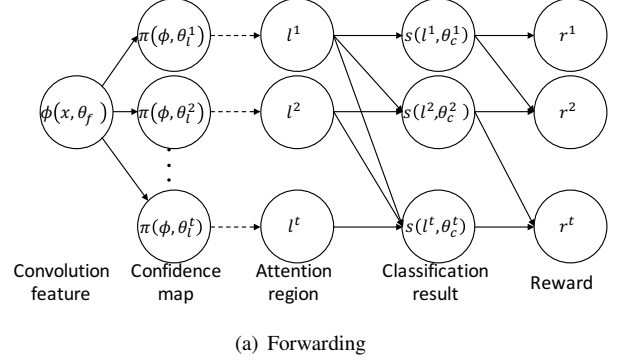


Figure 4. The forward (a) and back-propagation (b) processes for training attention networks as MDPs. The dashed lines indicate the sampling procedures.

### 3.4. Optimization

It is difficult to directly compute the gradient of  $\mathbb{E}_{\theta} [r_n^t]$  over  $\theta$  because it requires evaluating exponentially many possible part locations during training. Hence we employ REINFORCE algorithm and approximate the gradient in a Monte Carlo way [43].

$$\nabla_{\theta} \mathbb{E}_{\theta} [r_n^t] \approx \frac{1}{K} \sum_{k=1}^K \nabla_{\theta} (\log \pi(l_{nk}^t | \phi(x_n), \theta_l^t)) r_{nk}^t \quad (7)$$

where  $l_{nk}^t \sim \pi(l_n^t | \phi(x_n), \theta_l^t)$  is sampled according to a multinomial distribution parameterized by the output confidence map of the  $t$ -th attention network.

The forward process of training the attention networks as MDPs is shown in Fig. 4. Given the basis convolutional feature maps  $\phi(x)$  as input, the attention networks output the confidence map  $\pi(\phi, \theta_l^t)$  at different time step  $t$ . Each  $\pi(\phi, \theta_l^t)$  forms a multinomial distribution, and the location of attention region  $l^t$  is sampled under the distribution. The sampling procedure is repeated for  $K$  times. We then use them for classification network and further get the reward  $r^t$ .

During back-propagation, the gradient  $\nabla_{\theta} L(\theta_f, \theta_c)$  can be obtained by back-propagating the classification networks. The gradient  $\nabla_{\theta} R(\theta_f, \theta_l)$  is calculated using policy gradient as shown in Equation 7. Notice that when the reward is 0, we can just ignore the sample.

### 3.5. Implementation Details

**Step-wise training:** Although jointly training the entire model is possible, we develop a 3-step algorithm for the sake of training speed. In the first step, we initialize and fine-tune the CNN model to extract the basis convolutional feature maps for attention and classification. In the second step, we fix and cache the basis convolutional feature maps from the first step, and train the attention networks separately. In the third step, we fix and cache the selected attention regions from the second step, and fine-tune the final classification model. Through feature caching, repeated feature calculating is avoided. Notice that the convolutional neural networks for attention and the final classification is different, though they are initialized similarly during pre-training described below. We repeat these steps several times until convergence.

**Fast-RCNN approximation:** During training, although we can compute a multi-scale feature maps to obtain the features for high resolution region crops. It could still be time-consuming when the image resolution is large. Thus, we employ an approximated feature extraction method that is similar to Fast-RCNN [19], where we only compute a feature map from the input image at one scale. The convolutional features for each part is obtained by selecting the corresponding region in the convolutional feature map of the whole image, so that the receptive field of the selected region is the same as the size of the part. This further accelerates the training of attention networks.

Note that since we adopt a 3-step training, the Fast-RCNN approximation are only utilized during attention network training. The final classification networks are still trained given the features extracted from cropped high resolution images.

### 3.6. Discussion

Our attention component is inspired from the recurrent visual attention model [16]. However, instead of building a recurrent attention network that share parameters over different time steps, our model uses multiple convolutional networks with different parameters to model the temporal effect. During testing, these attention networks work like independent part detectors that share the same basis image feature. It is even possible to combine all the attention networks into a single convolutional network to compute part attentions simultaneously. This makes inference much faster.

Dataset	#Class	#Train	#Test	BBox	Part
Stanford Dogs [2]	120	12,000	8,580	✓	
Stanford Cars [3]	196	8,144	8,041	✓	
CUB-200-2011 [1]	200	5,994	5,794	✓	✓
Food-101 [5]	101	75,750	25,250		

Table 1. Statistics for the four fine-grained benchmark datasets.

## 4. Experiments

We conduct extensive experiments on four benchmark datasets, including CUB-200-2011 [1], Stanford Dogs [2], Stanford Cars [3] and Food-101 [5]. Table 1 shows the statistics of the four datasets.

### 4.1. Experimental Setup

We use the ResNet-50 [13] for feature extraction. During pre-training, we first resize all images to  $512 \times 512$  resolution, and fine-tune the ResNet-50 with randomly cropped  $448 \times 448$  patches. For each input image, ResNet-50 outputs a  $2048 \times 16 \times 16$  `res_5c` feature map. We then use the feature map to train the attention networks to find two parts. The first part selects a  $4 \times 4$  region in the feature map (corresponding to a  $128 \times 128$  patch in the resized image), and the second one selects a  $8 \times 8$  region (corresponding to a  $256 \times 256$  patch in the resized image). We then crop the two result attention patches and resize to  $512 \times 512$  to train ResNet-50 prediction models in the final classification stage.

All models are trained using RMSProp with batch size of 512 and 90 epochs. The initial learning rate is 0.01 and multiplied by 0.1 every 30 epochs. Our implementation is based on Caffe [44].

### 4.2. Computational Time

On Stanford Dogs dataset, our FCANs take 3 hours to train on a single Tesla K40 GPU, significantly faster than a conventional recurrent attention model [17] that takes about 30 hours to converge in our implementation. Fine-tuning the convolutional features requiring additional training time for both models. For an image with  $512 \times 512$  resolution, our testing time is  $\sim 150$ ms. The cost of attention selection is negligible compared with the feature calculation time. Compared with recurrent attention models [17] that takes  $\sim 250$ ms during testing, our method is faster.

### 4.3. Comparison with State-of-the-Arts

We compare our framework with all previous methods and summarize the results from Table 3 to Table 5.

On CUB-200-2011, our recognition accuracy (84.3%) is comparable with all state-of-the-art methods [25, 42, 24] without using ground-truth bounding boxes during testing.

CUB-200-2011	Accuracy(%)	Acc w. Box(%)
Zhang <i>et al.</i> [32]	73.9	76.4
Branson <i>et al.</i> [15]	75.7	85.4*
Simon <i>et al.</i> [39]	81.0	-
Krause <i>et al.</i> [22]	82.0	82.8
Lin <i>et al.</i> [25]	84.1	85.1
Jaderberg <i>et al.</i> [42]	84.1	-
Kong <i>et al.</i> [24]	84.2	-
Our Model	<b>84.3</b>	84.7

Table 2. Comparison to related work on CUB-200-2011 dataset. \* Testing with both ground truth box and parts.

Stanford Dogs	Accuracy(%)	Acc w. Box(%)
Gavves <i>et al.</i> [37]	-	50.1
Simon & Rodner [39]	68.1	-
Sermanet <i>et al.</i> [17]	76.8	-
Zhang <i>et al.</i> [45]	79.9	-
Krause <i>et al.</i> [40]	82.6	-
Our Model	<b>88.9</b>	-

Table 3. Comparison to related work on Stanford Dogs dataset.

Stanford Cars	Accuracy(%)	Acc w. Box(%)
Chai <i>et al.</i> [46]	78.0	-
Gosselin <i>et al.</i> [47]	82.7	87.9
Girshick <i>et al.</i> [48]	88.4	-
Lin <i>et al.</i> [25]	91.3	-
Wang <i>et al.</i> [49]	-	92.5
Krause <i>et al.</i> [22]	92.6	92.8
Our Model	91.5	<b>93.1</b>

Table 4. Comparison to related work on Stanford Cars dataset.

Method	Accuracy(%)	Acc w. Box(%)
L. Bossard <i>et al.</i> [5]	50.8	-
A. Myers <i>et al.</i> [50]	79.0	-
Our Model	<b>86.3</b>	-

Table 5. Experimental results on Food-101 dataset.

On Stanford Dogs, Stanford Cars and Food-101, our model is also very competitive. For example, we obtain 93.1% accuracy on Stanford Cars test set with bounding box during testing, which is so far the best result. Note that our baseline method Sermanet *et al.* [17] uses reinforcement learning based recurrent attention models, which is similar to our approach. Our method improves them by more than 12% on Stanford Dogs, suggesting the FCANs as an effective framework for fine-grained recognition.

Method	Dogs	Cars	Birds	Foods
Finetune baseline	87.3	89.7	82.0	82.1
+ Random regions	87.9	90.1	82.3	83.0
+ Center regions	87.5	90.6	82.4	82.7
+ Attention regions	<b>88.9</b>	<b>91.5</b>	<b>84.3</b>	<b>86.3</b>

Table 6. Experimental comparison on the effect of attentions.

Method	Dogs	Cars	Birds	Foods
Finetune baseline	87.3	87.5	82.0	82.1
One attention only	88.1	84.2	80.4	79.9
+ One attention	88.5	90.2	83.3	85.5
+ Two attentions	88.9	91.5	<b>84.3</b>	86.3
+ More attentions	<b>89.0</b>	<b>91.6</b>	<b>84.3</b>	<b>86.5</b>

Table 7. Experimental comparison on the number of attentions.

#### 4.4. Ablation Study

**Effect of Attention:** Since our approach is roughly three times (full image + two attention regions) more expensive than a single model during testing, we conduct two additional model-fusion baselines to demonstrate its superiority. One is the random region experiment, where we augment the baseline single image model with two random cropped regions. The second baseline is the center region experiment, where we crop two center regions in the image. The sizes of the two crops in both experiments are the same as the sizes of the parts in the attention model. Table 6 summarizes the results. When costing the same amount of testing time, the attention networks clearly outperform random region and center region models.

**Number of Attentions:** Table 7 summarizes the results of how the number of attentions affects the final classification accuracy. Take Stanford Dogs as an example, after fine-tuning the baseline ResNet-50 achieves 87.3% accuracy. Combining one  $8 \times 8$  attention region with the prediction results of original image improves significantly to 88.5%. Combining one  $8 \times 8$  region, one  $4 \times 4$  region and the original image together further improves the results to 88.9%. We find adding more than two attentions (i.e. 3 attentions) only improves the performance slightly at the expense of more computations. Hence throughout the experiments we fix the number of attentions as two.

**Reward Strategy:** Table 8 illustrates the effectiveness of our training reward strategy. Compared against the traditional reward setting which only assigns a reward after all attention iterations, our greedy reward strategy works significantly better. We hypothesize that the greedy reward helps the reinforcement learning to quickly converge to discriminative sub-regions.



Figure 5. Qualitative comparison between our method (left) and recurrent attention [17] (right) on different datasets. On the left, we plot the first two attention regions regenerated by FCAN, which corresponds to  $4 \times 4$  and  $8 \times 8$  attention regions respectively (lighter color indicates higher score). On the right, we also show the first two selected regions by [17] using our implementation.

Method	Dogs	Cars	Birds	Foods
Baseline reward	88.1	90.5	82.9	84.7
Greedy reward	<b>88.9</b>	<b>91.5</b>	<b>84.3</b>	<b>86.3</b>

Table 8. Experimental comparison on the reward strategy. The baseline reward strategy only assigns a reward after all attention iterations.

#### 4.5. Qualitative Results

We qualitatively compare the attention regions selected by our model and the recurrent attention model [17] in Fig. 5. Both models contain attention mechanisms and apply reinforcement learning to train to focus on local discriminative regions. We observe that in our model different attentions correspond to different image regions, while the attention regions generated in [17] focus on only one re-

gion. Our attention map is also more diverse than the attention map in [17]. This illustrates how our attention model outperforms the previous reinforcement learning based attention work.

## 5. Conclusion

In this paper, we present Fully Convolutional Attention Networks (FCANs) for fine-grained recognition. With the fully convolutional architecture, our model is much faster than previous reinforcement learning based visual attention models during both training and testing. We conduct extensive experiments on four different fine-grained benchmark datasets and show its competitive performance against state-of-the-art methods.



## References

- [1] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, “The caltech-ucsd birds-200-2011 dataset,” 2011. 0, 5
- [2] A. Khosla, N. Jayadevaprakash, B. Yao, and F-F. Li, “Novel dataset for fine-grained image categorization: Stanford dogs,” in *Proc. CVPR Workshop on Fine-Grained Visual Categorization (FGVC)*, vol. 2, 2011. 0, 1, 5
- [3] 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. 0, 1, 5
- [4] M.-E. Nilsback and A. Zisserman, “Automated flower classification over a large number of classes,” in *Computer Vision, Graphics & Image Processing, 2008. ICVGIP’08. Sixth Indian Conference on*. IEEE, 2008, pp. 722–729. 0, 1
- [5] L. Bossard, M. Guillaumin, and L. Van Gool, “Food-101—mining discriminative components with random forests,” in *European Conference on Computer Vision*. Springer, 2014, pp. 446–461. 0, 1, 5, 6
- [6] S. Bell and K. Bala, “Learning visual similarity for product design with convolutional neural networks,” *ACM Transactions on Graphics (TOG)*, vol. 34, no. 4, p. 98, 2015. 0
- [7] K. Hadi, H. Xufeng, L. Svetlana, B. Alexander, and B. Tamara, “Where to buy it: Matching street clothing photos in online shops,” in *Computer Vision (ICCV), 2015 IEEE International Conference on*, vol. 8, 2015. 0
- [8] N. Kumar, P. N. Belhumeur, A. Biswas, D. W. Jacobs, W. J. Kress, I. C. Lopez, and J. V. Soares, “Leafsnap: A computer vision system for automatic plant species identification,” in *Computer Vision—ECCV 2012*. Springer, 2012, pp. 502–516. 0
- [9] T. Berg, J. Liu, S. Woo Lee, M. L. Alexander, D. W. Jacobs, and P. N. Belhumeur, “Birdsnap: Large-scale fine-grained visual categorization of birds,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 2011–2018. 0, 1
- [10] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105. 0, 1
- [11] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014. 0, 1, 2
- [12] 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. 0, 1, 2
- [13] 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. 0, 1, 2, 5
- [14] J. Liu, A. Kanazawa, D. Jacobs, and P. Belhumeur, “Dog breed classification using part localization,” in *European Conference on Computer Vision*. Springer, 2012, pp. 172–185. 0, 1, 2
- [15] S. Branson, G. Van Horn, S. Belongie, and P. Perona, “Bird species categorization using pose normalized deep convolutional nets,” *arXiv preprint arXiv:1406.2952*, 2014. 0, 1, 6
- [16] V. Mnih, N. Heess, A. Graves *et al.*, “Recurrent models of visual attention,” in *Advances in neural information processing systems*, 2014, pp. 2204–2212. 1, 2, 3, 4, 5
- [17] P. Sermanet, A. Frome, and E. Real, “Attention for fine-grained categorization,” *arXiv preprint arXiv:1412.7054*, 2014. 1, 2, 5, 6, 7
- [18] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine learning*, vol. 8, no. 3–4, pp. 229–256, 1992. 1
- [19] R. Girshick, “Fast r-cnn,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448. 1, 5
- [20] Y. Cui, F. Zhou, Y. Lin, and S. Belongie, “Fine-grained categorization and dataset bootstrapping using deep metric learning with humans in the loop,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1153–1162. 1
- [21] 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. 1
- [22] J. Krause, H. Jin, J. Yang, and L. Fei-Fei, “Fine-grained recognition without part annotations,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5546–5555. 1, 6
- [23] 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. 1
- [24] S. Kong and C. Fowlkes, “Low-rank bilinear pooling for fine-grained classification,” *arXiv preprint arXiv:1611.05109*, 2016. 1, 5, 6
- [25] 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. 1, 5, 6
- [26] Q. Qian, R. Jin, S. Zhu, and Y. Lin, “Fine-grained visual categorization via multi-stage metric learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3716–3724. 1
- [27] X. Wang, T. Yang, G. Chen, and Y. Lin, “Object-centric sampling for fine-grained image classification,” *arXiv preprint arXiv:1412.3161*, 2014. 1, 2
- [28] E. Rosch, C. B. Mervis, W. D. Gray, D. M. Johnson, and P. Boyes-Braem, “Basic objects in natural categories,” *Cognitive psychology*, vol. 8, no. 3, pp. 382–439, 1976. 1

- [29] R. Farrell, O. Oza, N. Zhang, V. I. Morariu, T. Darrell, and L. S. Davis, “Birdlets: Subordinate categorization using volumetric primitives and pose-normalized appearance,” in *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE, 2011, pp. 161–168. 2
- [30] N. Zhang, R. Farrell, F. Iandola, and T. Darrell, “Deformable part descriptors for fine-grained recognition and attribute prediction,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 729–736. 2
- [31] N. Zhang, M. Paluri, M. Ranzato, T. Darrell, and L. Bourdev, “Panda: Pose aligned networks for deep attribute modeling,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1637–1644. 2
- [32] 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. 2, 6
- [33] L. Bourdev and J. Malik, “Poselets: Body part detectors trained using 3d human pose annotations,” in *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2009, pp. 1365–1372. 2
- [34] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, “Object detection with discriminatively trained part-based models,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 32, no. 9, pp. 1627–1645, 2010. 2
- [35] T. Berg and P. Belhumeur, “Poof: Part-based one-vs.-one features for fine-grained categorization, face verification, and attribute estimation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 955–962. 2
- [36] S. Yang, L. Bo, J. Wang, and L. G. Shapiro, “Unsupervised template learning for fine-grained object recognition,” in *Advances in neural information processing systems*, 2012, pp. 3122–3130. 2
- [37] E. Gavves, B. Fernando, C. G. Snoek, A. W. Smeulders, and T. Tuytelaars, “Fine-grained categorization by alignments,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 1713–1720. 2, 6
- [38] Y. Chai, V. Lempitsky, and A. Zisserman, “Bicos: A bi-level co-segmentation method for image classification,” in *Computer Vision (ICCV), 2011 IEEE International Conference on*. IEEE, 2011, pp. 2579–2586. 2
- [39] M. Simon and E. Rodner, “Neural activation constellations: Unsupervised part model discovery with convolutional networks,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1143–1151. 2, 6
- [40] J. Krause, B. Sapp, A. Howard, H. Zhou, A. Toshev, T. Duerig, J. Philbin, and L. Fei-Fei, “The unreasonable effectiveness of noisy data for fine-grained recognition,” in *European Conference on Computer Vision*. Springer, 2016, pp. 301–320. 2, 6
- [41] J. Ba, V. Mnih, and K. Kavukcuoglu, “Multiple object recognition with visual attention,” *arXiv preprint arXiv:1412.7755*, 2014. 2, 3
- [42] M. Jaderberg, K. Simonyan, A. Zisserman *et al.*, “Spatial transformer networks,” in *Advances in Neural Information Processing Systems*, 2015, pp. 2017–2025. 2, 5, 6
- [43] R. S. Sutton, D. A. McAllester, S. P. Singh, Y. Mansour *et al.*, “Policy gradient methods for reinforcement learning with function approximation.” in *NIPS*, vol. 99, 1999, pp. 1057–1063. 4
- [44] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, “Caffe: Convolutional architecture for fast feature embedding,” in *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 2014, pp. 675–678. 5
- [45] Y. Zhang, X.-s. Wei, J. Wu, J. Cai, J. Lu, V. Nguyen, and M. Do, “Weakly supervised fine-grained image categorization. arxiv preprint,” *arXiv*, vol. 1504, p. 7, 2015. 6
- [46] Y. Chai, V. Lempitsky, and A. Zisserman, “Symbiotic segmentation and part localization for fine-grained categorization,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 321–328. 6
- [47] P.-H. Gosselin, N. Murray, H. Jégou, and F. Perronnin, “Revisiting the fisher vector for fine-grained classification,” *Pattern Recognition Letters*, vol. 49, pp. 92–98, 2014. 6
- [48] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587. 6
- [49] Y. Wang, J. Choi, V. Morariu, and L. S. Davis, “Mining discriminative triplets of patches for fine-grained classification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1163–1172. 6
- [50] A. Meyers, N. Johnston, V. Rathod, A. Korattikara, A. Gorbun, N. Silberman, S. Guadarrama, G. Papandreou, J. Huang, and K. P. Murphy, “Im2calories: towards an automated mobile vision food diary,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1233–1241. 6