

Weakly Supervised Attention Pyramid Convolutional Neural Network for Fine-Grained Visual Classification

Yifeng Ding¹ Shaoguo Wen¹ Jiyang Xie¹ Dongliang Chang¹
 Zhanyu Ma¹ Zhongwei Si¹ Haibin Ling²

¹Beijing University of Posts and Telecommunications, China ²Stony Brook University, New York

Abstract

Classifying the sub-categories of an object from the same super-category (e.g., bird species, car and aircraft models) in fine-grained visual classification (FGVC) highly relies on discriminative feature representation and accurate region localization. Existing approaches mainly focus on distilling information from high-level features. In this paper, however, we show that by integrating low-level information (e.g., color, edge junctions, texture patterns), performance can be improved with enhanced feature representation and accurately located discriminative regions. Our solution, named Attention Pyramid Convolutional Neural Network (AP-CNN), consists of a) a pyramidal hierarchy structure with a top-down feature pathway and a bottom-up attention pathway, and hence learns both high-level semantic and low-level detailed feature representation, and b) an ROI guided refinement strategy with ROI guided dropblock and ROI guided zoom-in, which refines features with discriminative local regions enhanced and background noises eliminated. The proposed AP-CNN can be trained end-to-end, without the need of additional bounding box/part annotations. Extensive experiments on three commonly used FGVC datasets (CUB-200-2011, Stanford Cars, and FGVC-Aircraft) demonstrate that our approach can achieve state-of-the-art performance. Code available at <http://dwz1.cc/ci8so8a>

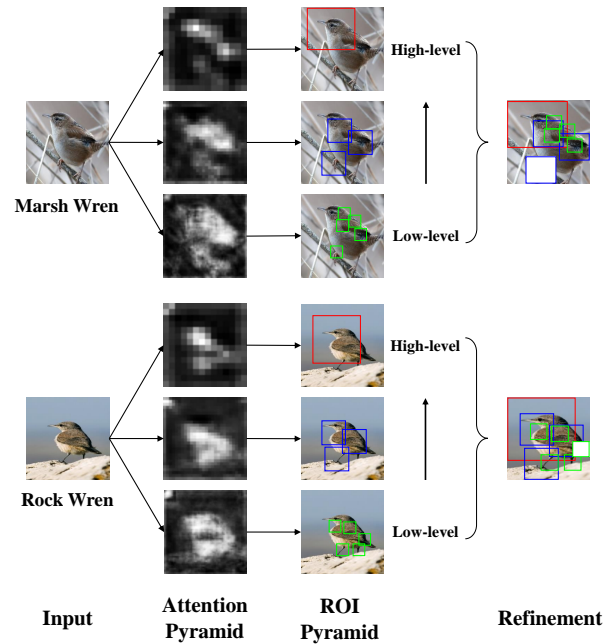


Figure 1. The discriminative regions of interest from different pyramidal hierarchy learned by AP-CNN for two bird species of “wren”. It can be observed that the low-level information can capture more subtle parts to distinguish the birds, e.g., the texture of wings and the shape of claws. Refinement is conducted on features with background noises eliminated and discriminative parts enhanced.

1. Introduction

The fine-grained visual classification (FGVC) task focuses on differentiating sub-categories of the objects from the same super-category (e.g., bird species, cars and aircrafts models). It has attracted extensive attention recently due to a wide range of applications such as expert-level image recognition [30], rich image captioning [13], intelligent retail [1], and intelligent transportation [31]. Different from the traditional image classification task, images from different sub-classes in the FGVC problems share close sim-

ilarities. At the same time, it differs from the face recognition task, because the faces are aligned into similar directions in face recognition while different poses are often occurred in FGVC. As a result, the challenging and distinctive keystones of the FGVC problem are a) high intra-class variance: objects that belong to the same category usually present significantly different poses and viewpoints; and b) low inter-class variance: the visual differences among the subordinate classes are often subtle as they belong to the same super-category.

In order to address the aforementioned challenges, early solutions [4, 35, 12] introduce additional bounding box/part annotations to help locate the target object and align the components in an image. Although effective, these methods are not optimal for FGVC as human annotations can be very time consuming and need professional knowledge. Recent methods solve this problem in a weakly supervised manner. These approaches can be classified into two categories: a) feature encoding methods [19, 14, 33] that extract the fine-grained features by encoding a highly parameterized representation of the features, and b) region locating methods [36, 6, 32, 8] that figure out the discriminative regions by learning part detectors, and then conduct refinement such as cropping and amplifying the attended parts on multiple stages.

Although promising results have been reported in the above studies, further improvement suffers from lack of using low-level information. Our study shows that low-level information (*e.g.*, color, edge junctions, texture patterns [34]), is indeed essential in the FGVC task. Specifically, as the structure of CNN getting deeper, the neurons in high layers are strongly respond to entire images and rich in semantics, but inevitably lose detailed information from small discriminative regions. Figure 1 shows examples of differences in activations extracted from diverse layers. such detailed information, *i.e.* the low-level information, is helpful in the FGVC task as it reflects the subtle difference within various sub-classes, and is invariable no matter how the pose or viewpoint changes. Existing FGVC methods pay much attention to high-level features, and in this work, we use additional enhanced low-level information as supplement.

The motivation of this paper is to effectively integrate both the high-level semantic and the low-level detailed information for fine-grained classification. To this end, we propose a novel attention pyramid convolutional neural network (AP-CNN), which jointly learns multi-level information and refined feature representations without using bounding box/part annotations. The proposed method can accurately locate the discriminative local regions as well as reduce the background noise. The main contribution can be summarized as follows:

1) We propose a novel attention pyramid convolutional neural network (AP-CNN) by building an enhanced pyramidal hierarchy, which combines a top-down pathway of features and a bottom-up pathway of attentions, and thus learns both high-level semantic and low-level detailed feature representations.

2) We propose ROI guided refinement consisting of ROI guided dropblock and ROI guided zoom-in to further refine the features. The dropblock operation helps to locate more discriminative local regions, and the zoom-in operation aligns features with background noises eliminated.

3) We conduct extensive experiments on three commonly used FGVC datasets (CUB-200-2011 [28], Stanford-Cars [16], and FGVC-Aircraft [23]). Visualization and ablation studies are further conducted to draw insights into our method. The results demonstrate that our model can significantly improve the accuracy of fine-grained classification.

2. Related Work

Methods Using Multi-level Features. Features from different layers are commonly used in detection and segmentation tasks. NoCs [25] extracts feature maps on input images with different scales and then conduct feature fusion. FCN [22], U-Net [26], FPN [18] fuse information from lower layers to high-level features through skip-connections. Meanwhile, HyperNet [15], ParseNet [21], and ION [2] concatenate features of multiple layers before computing predictions. SSD [20] and MS-CNN [3] predict individual target locations at multiple layers without combining features.

In this paper, we use the multi-level features in the FGVC task for better classification and weakly supervised detection. Besides, we further enhance the multi-level features by establishing strong correlations between them. This is done through a top-down feature pathway delivering the semantic information from high levels to low levels, combined with a bottom-up attention pathway carrying low-level information back to the top.

Weakly Supervised Fine-grained Classification. We define weakly supervised learning in FGVC task as methods without bounding box/part annotations. This setting is generally applied in recent methods due to its feasibility in real-world scenarios.

Feature encoding-based approaches [19, 7, 14, 33] encode higher order information on features. The classical benchmark, Bilinear-CNN [19], utilizes a bilinear pooling operation to aggregate the pairwise interactions between features in two independent CNNs, and computes the outer product over the output feature channels of the two streams to capture the second-order information. Further works [7, 14] reduce the huge computation cost of the outer product by a single-stream output and adopting low-rank approximation to the covariance matrix, respectively. Moreover, [33] proposes a cross-layer bilinear pooling approach to capture the inter-layer part feature relations.

Region locating methods [6, 36, 32, 8], another common way in FGVC, generally apply localization networks to detect discriminative regions in images. RA-CNN [6] is proposed to zoom in discriminative local regions learned by a novel attention proposal network. Meanwhile, MA-CNN [36] generates multiple object parts by clustering channels of feature maps into different groups. NTS [32] enables a navigator agent as the region proposal network to detect multiple most informative regions under the guidance from

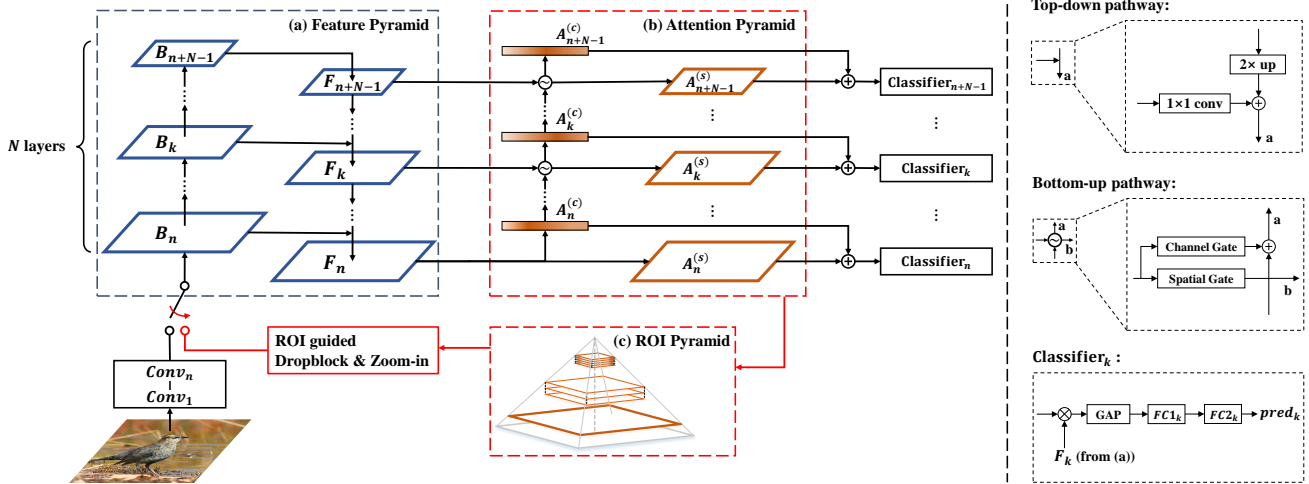


Figure 2. Illustration of AP-CNN. (a) FPN backbone. (b) Attention pyramid module. (c) ROI pyramid. The refinement operation is shown in the red flow and conducted on the low-level features B_n . In this figure, the feature maps are indicated by blue outlines and the channel/spatial attentions are indicated in orange. The structure details of the top-down and the bottom-up pathways, and the classifiers are illustrated on the right. \oplus represents broadcasting addition and \otimes represents element-wise multiplication.

a teacher agent. WS-CPM [8] develops a novel pipeline combined with object detection and segmentation using Mask R-CNN, followed by a bi-directional long short-term memory (LSTM) network to integrate and encode the partial information.

The most relevant work to ours comes from NTS [32], which also applies pyramidal features to the FGVC task. However, the NTS network only learns the region localizers by simply applying the FPN [18] structure on the CNN, which poses challenges to accurate region localization in the weakly supervised way. Besides, it ignores the fact that these pyramidal features can also contribute to the classification.

Compared with NTS, the advantages of our work are two-folds. First, we introduce the pyramid hierarchy to the FGVC task and further enhance its representation. We use these multi-level information not only for precise region localization but also for better classification. Second, we conduct refinement that takes full advantage of multi-level features, by using small ROIs learned from low-level features for dropblock operation, and using bounding rectangle merged by ROIs from all levels for zoom-in operation.

3. Attention Pyramid Convolutional Neural Network

In this section, we introduce the proposed Attention Pyramid Convolutional Neural Network (AP-CNN) for fine-grained classification. AP-CNN is a two-stage network that respectively takes coarse full images and refined features as input. These two stages, which we define as the raw-stage and the refined-stage, share the same network ar-

chitecture with the same parameters to extract information from both the coarse and the refined inputs.

An overview of the proposed AP-CNN is illustrated in Figure 2. Feature and attention pyramid structure, and ROI guided refinement are conducted for improving performance. First, the feature and attention pyramid structure takes coarse images as input, which generates the pyramidal features and the pyramidal attentions by establishing hierarchy on the basic CNN following a top-down feature pathway and a bottom-up attention pathway. Second, once the spatial attention pyramid has been obtained from the raw input, the region proposal network (RPN) proceeds to generate the pyramidal regions of interest (ROIs) in a weakly supervised way. Then the ROI guided refinement is conducted on low-level features with a) the ROI guided dropblock which erases the most discriminative regions selected from small-scaled ROIs, and b) the ROI guided zoom-in which locates the major regions merged from all ROIs. Third, the refined features are sent into the refined-stage to distill more discriminative information. Both stages set individual classifiers for each pyramid level, and the final classification result is averaged over the raw-stage predictions and the refined-stage predictions. Note that the AP-CNN can be trained end-to-end, and the framework is flexible in the CNN backbone structure (e.g., AlexNet [17], VGG [27] and ResNet [10]). In this paper, we present results using VGG16 and ResNet50.

3.1. Attention Pyramid Model

Motivation. Our goal is to leverage both the semantic and the detailed information for improving FGVC per-

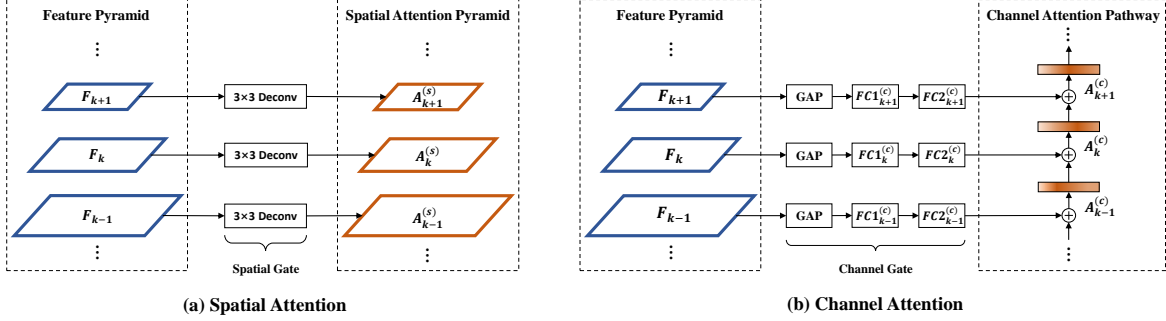


Figure 3. Process of getting (a) spatial attention pyramid and (b) channel attention pathway.

formance. Specifically, CNN backbones apply a series of convolution blocks, we denote the output feature maps of blocks with the different spatial size as $\{B_1, B_2, \dots, B_l\}$, where l indicates the number of blocks. Traditional methods use the last feature map B_l for classification. It contains strong semantics but lacks detailed information, which is adverse to the FGVC task. FPN chooses part of these features and generates N corresponding feature hierarchy $\{F_n, F_{n+1}, \dots, F_{n+N-1}\}$ ($1 \leq n \leq n+N-1 \leq l$) by applying a) a top-down pathway after B_{n+N-1} , which upsamples spatially coarser but semantically stronger feature maps from higher pyramid levels to lower pyramid levels, and b) lateral connections between corresponding features $B_k \rightarrow F_k$ ($k = n, n+1, \dots, n+N-1$) to maintain the backbone information. The pyramidal features can locate samples on different scales, which is also beneficial to the FGVC task, to focus on the subtle differences of objects from different scales.

We further enhance the FPN structure by introducing an additional attention hierarchy $\{A_n, A_{n+1}, \dots, A_{n+N-1}\}$ upon the pyramidal features, which consists of (a) pyramidal spatial attentions $\{A_n^{(s)}, A_{n+1}^{(s)}, \dots, A_{n+N-1}^{(s)}\}$ to locate discriminative regions from different scales, and (b) pyramidal channel attentions $\{A_n^{(c)}, A_{n+1}^{(c)}, \dots, A_{n+N-1}^{(c)}\}$ to embed channel correlations and deliver local information from lower pyramid levels to higher pyramid levels in an additional bottom-up pathway.

Spatial Gate and Spatial Attention Pyramid. As shown in Figure 3(a), each building block takes the corresponding feature map F_k as input and generates a spatial attention mask $A_k^{(s)}$. Specifically, each feature map F_k first goes through a 3×3 deconvolution layer with one output channel to squeeze spatial information. Then each element of the spatial attention mask $A_k^{(s)}$ is normalized to the interval (0,1) using the sigmoid function to reflect the spatial importance:

$$A_k^{(s)} = \sigma(v_c * F_k). \quad (1)$$

Here σ refers to the sigmoid function, while $*$ de-

notes deconvolution and v_c represents convolution kernel. As a result, we get spatial attention pyramid $\{A_n^{(s)}, A_{n+1}^{(s)}, \dots, A_{n+N-1}^{(s)}\}$ based on multi-scale feature maps. We use these spatial activations to generate the ROI pyramid and conduct further refinement on features, which is described as below.

Channel Gate and Channel Attention Pathway. Inspired by SE-Net [11], channel attentions $\{A_n^{(c)}, A_{n+1}^{(c)}, \dots, A_{n+N-1}^{(c)}\}$ are gained from corresponding feature maps in the feature pyramid by operating the global average pooling (GAP) combined with two fully-connected (FC) layers. The channel attention mask can be represented as:

$$A_k^{(c)} = \sigma(W_2 \cdot \text{ReLU}(W_1 \cdot \text{GAP}(F_k))), \quad (2)$$

where $\text{GAP}(\cdot)$ is the global average pooling function:

$$\text{GAP}(F_k) = \frac{1}{W \times H} \sum_{i=1}^H \sum_{j=1}^W F_k(i, j). \quad (3)$$

Here σ and ReLU refer to the sigmoid and the ReLU function, respectively. The dot product denotes element-wise multiplication. W and H represent the spatial dimensions of F_k . W_1 and W_2 are the weight matrices of two FC layers. In our framework, channel attentions play a different role from the spatial attention pyramid as they are settled for delivering low-level detailed information in a bottom-up pathway from lower pyramid levels to higher pyramid levels. Figure 3(b) shows the flow diagram.

Classifier. We use the learned attentions to weight features F_k , and get F'_k for classification:

$$F'_k = F_k \cdot (A_k^{(s)} \oplus A_k^{(c)}), \quad (4)$$

where \oplus represents the addition operation using broadcasting semantics. Individual classifiers with a GAP layer and two FC layers are settled to make final predictions.

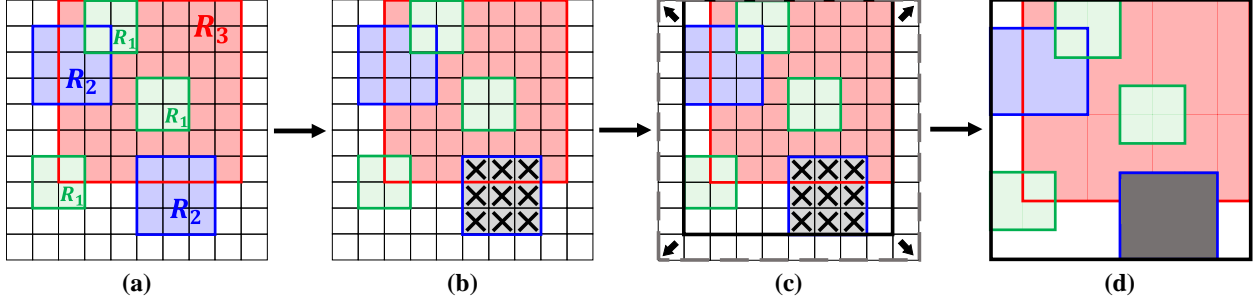


Figure 4. ROI guided refinement based on Algorithm 1. (a) Low-level feature map B_n with ROI pyramid R_{all} (R_1, R_2, R_3 for example). (b) ROI guided Dropblock. (c) ROI guided zoom-in. (d) Refined feature map Z_n with local regions masked and background noises eliminated. Rectangles with different colors are from distinctive levels in the pyramidal hierarchy, while the gray one represents the dropped region.

Algorithm 1 ROI guided refinement algorithm.

Input: Low-level feature map B_n ; ROI pyramid $R_{all} = \{R_n, R_{n+1}, \dots, R_{n+N-1}\}$; Dropblock candidates selection probabilities $P = \{p_n, p_{n+1}, \dots, p_{n+N-1}\}$;
Output: Refined feature map Z_n ;

- 1: **if** training **then**
- 2: Randomly select R_s from R_{all} ($n \leq s \leq n+N-1$) with probabilities P ;
- 3: Randomly select r_s from $R_s = \{r_{s,1}, r_{s,2}, \dots, r_{s,\xi_s}\}$ with equal probability;
- 4: Calculate dropblock feature D_n according to Eq. 6;
- 5: **else**
- 6: $D_n = B_n$;
- 7: **end if**
- 8: Merge R_{all} and obtain the bounding box;
- 9: Crop the region and enlarge as Z_n by Eq. 7;
- 10: **return** Z_n ;

3.2. ROI Guided Refinement

ROI Pyramid. Region proposal network (RPN) [24] is a widely used structure in visual detection to locate possible informative regions. Recent RPN designs [18] are conducted on a single-scale convolutional feature map or multi-scale convolutional feature maps. Anchors of multiple pre-defined scales and aspect ratios are designed to cover objects of different shapes.

In this work, we adapt RPN on the spatial attention pyramid, rather than the multi-scale feature maps. We assign the anchors with a single scale and ratio to each pyramidal level according to its convolutional receptive field, and adopt non-maximum suppression (NMS) on the region proposals to reduce region redundancy. Specifically, for each pyramid level k , we select the top- ξ_k informative regions $R_k = \{r_{k,1}, r_{k,2}, \dots, r_{k,\xi_k}\}$ based on the responding activation values of anchors in the spatial attention mask, and then form the region pyramid $R_{all} =$

$\{R_n, R_{n+1}, \dots, R_{n+N-1}\}$. As a result, we get an ROI pyramid in the raw-stage, and then conduct ROI guided refinement (ROI guided dropblock and ROI guided zoom-in) on the pyramid bottom features B_n to further improve the performance in the refined-stage. Algorithm 1 shows the main procedure of the refinement operations and Figure 4 illustrates a visualized example.

ROI Guided Dropblock. Overfitting is a common problem in deep-learning, especially in the FGVC task as each species only contains small number of images. [9] proposes the dropblock strategy by dropping continuous regions randomly on feature maps to remove certain semantic information and consequently enforcing network to learn information on the remaining units. In this paper, we randomly select an ROI union R_s from R_{all} ($n \leq s \leq n+N-1$) with probabilities $P = \{p_n, p_{n+1}, \dots, p_{n+N-1}\}$. Then we randomly choose an informative region $r_s \in R_s$ with equal probability from the selected R_s . We scale r_s into the same sampling rate as B_n , and obtain drop mask M by setting activations in the ROI region to zero:

$$M(i, j) = \begin{cases} 0, & (i, j) \in r_s \\ 1, & \text{otherwise,} \end{cases} \quad (5)$$

and obtain dropped feature maps D_n by applying the mask on the low-level features B_n with normalization:

$$D_n = B_n \times M \times \text{Count}(M) / \text{Count_ones}(M), \quad (6)$$

where $\text{Count}(\cdot)$ and $\text{Count_ones}(\cdot)$ denote the number of all elements and the number of elements with value one, respectively.

Different from the original random dropblock, our ROI guided dropblock can directly erase the informative part and encourage the network to find more discriminative regions, which achieves higher accuracy. Table 4 shows the comparison of the random dropblock and our ROI guided dropblock. Note that we only conduct dropblock refinement in

| Method | Base | Pre-trained | Image resolution | Accuracy (%) | | |
|---------------------|-----------|---------------|---------------------|--------------|---------------|---------------|
| | | | | CUB-200-2011 | Stanford Cars | FGVC-Aircraft |
| FT VGGNet [29] | VGG19 | ImageNet | 448 × 448 | 77.8 | 84.9 | 84.8 |
| FT ResNet [27] | ResNet50 | ImageNet | 448 × 448 | 84.1 | 91.7 | 88.5 |
| B-CNN [19] | VGG16 | ImageNet | 448 × 448 | 84.1 | 91.3 | 84.1 |
| MA-CNN [36] | VGG19 | ImageNet | 448 × 448 | 86.5 | 92.8 | 89.9 |
| DFL [29] | ResNet50 | ImageNet | 448 × 448 | 87.4 | 93.1 | 91.7 |
| NTS [32] | ResNet50 | ImageNet | 448 × 448 | 87.5 | 93.9 | 91.4 |
| DCL [5] | ResNet50 | ImageNet | 448 × 448 | 87.8 | 94.5 | <u>93.0</u> |
| TASN [37] | ResNet50 | ImageNet | 448 × 448 | 87.9 | 93.8 | - |
| WS-CPM [8] | GoogLeNet | ImageNet+COCO | Shorter side is 800 | 90.3 | - | - |
| AP-CNN (one stage) | VGG19 | ImageNet | 448 × 448 | 85.4 | 93.2 | 91.5 |
| AP-CNN (two stages) | VGG19 | ImageNet | 448 × 448 | 86.7 | <u>94.6</u> | 92.9 |
| AP-CNN (one stage) | ResNet50 | ImageNet | 448 × 448 | 87.2 | <u>93.6</u> | 92.2 |
| AP-CNN (two stages) | ResNet50 | ImageNet | 448 × 448 | <u>88.4</u> | 95.4 | 94.1 |

Table 1. Comparison results on CUB-200-2011, Stanford Cars, and FGVC-Aircraft datasets. The best and second-best results are respectively marked in bold and underlined fonts.

the training process, but in the test process we skip this operation.

ROI Guided Zoom-in. We merge ROIs from all pyramid levels to learn the minimum bounding rectangle of the input image in a weakly-supervised way, and get $[t_{x1}, t_{x2}, t_{y1}, t_{y2}]$ denoting the minimum and maximum coordinates in terms of x and y axis of the merged bounding rectangle, respectively. Then we extract this region out from the dropped feature maps D_n , and enlarge it to the same size as D_n to get the zoom-in features Z_n :

$$Z_n = \varphi(D_n[t_{y1} : t_{y2}, t_{x1} : t_{x2}]), \quad (7)$$

where φ represents the bilinear upsample operation. The refined features Z_n is sent to the refined-stage to conduct further prediction. The final prediction is made by averaging the raw-stage prediction and refined-stage prediction.

4. Experimental Results and Discussions

We conduct experiments on three FGVC benchmark datasets, including CUB-200-2011, Stanford-Cars, and FGVC-Aircraft. All the datasets contain a set of sub-categories of the same super-category. The following is a brief description of these datasets:

CUB-200-2011 [28] has 11,778 images from 200 classes officially split into 5,994 training and 5,794 test images.

Stanford-Cars [16] has 16,185 images from 196 classes officially split into 8,144 training and 8,041 test images.

FGVC-Aircraft [23] has 10,000 images from 100 classes officially split into 6,667 training and 3,333 test images.

4.1. Implement Details

We implement AP-CNN on 50-layer ResNet [10] pre-trained on ImageNet. Specifically, we choose the last output feature of the residual block conv3, conv4 and conv5

in ResNet50 to establish pyramidal hierarchy, denoting as B_3, B_4, B_5 respectively. We do not include conv1 and conv2 into the pyramid because of their large memory footprint. The refinement operation is conducted on B_3 (for detailed information of this choice, please refer to Table 5). The input images are resized into 448×448 , which is standard in the literature. We do not use extra bounding box/part annotations and compare our method with other weakly supervised approaches. We respectively assign anchors with single scales of 64, 128, 256 and 1:1 ratio for each pyramidal level and choose the top 5, 3, 1 anchors with the highest activation value as potential refinement candidates. The IOU threshold in NMS operation is set as 0.05 and the drop-block rate in Algorithm 1 is set as $\{30\%, 30\%, 0\%\}$. Note that most of the hyperparameters of AP-CNN are involved in the process of getting anchors and the dropblock operation, we set them with empirical experiences.

We use open-sourced Pytorch as our code-base, and train all the models on a single GTX 1080Ti GPU. Optimization is performed using Stochastic Gradient Descent with momentum 0.9 and a minibatch size of 16. The initial learning rate is set to 0.001 and drops to 0 using cosine anneal schedule. All models are trained for 100 epochs.

4.2. Comparison with State-of-the-Art Methods

Table 1 lists the performance evaluations on three aforementioned benchmark datasets. Each column includes 7 to 9 representative weakly supervised methods that have reported evaluation results on the corresponding datasets, including fine-tuned baselines, feature encoding methods and region locating methods. We display the results of our model based on the VGG16 and ResNet50 backbone. Compared with the above FGVC works, our AP-CNN achieves significant performance improvement on all the three datasets. The evaluation results can be summarized

as follows:

- On the CUB-200-2011 dataset, our AP-CNN (one stage) achieves a significant improvement from the corresponding backbones, with clear margins of 7.6% and 3.1% on VGG19 and ResNet50, respectively. By applying ROI guided refinement, AP-CNN (two stages) reaches 88.4% accuracy on ResNet50, which outperforms the existing methods using the same backbone, pre-trained datasets, and image resolution. Note that the WS-CPM model currently gets the highest classification accuracy with 90.3%, which mainly benefits from the extra pre-trained data and the high input resolution.
- On the Stanford Cars dataset, currently the state-of-the-art accuracy is achieved by the DCL model with 94.5%. Our method outperforms DCL for a clear margin (0.9% relative gain) with accuracy 95.4%.
- On the FGVC-Aircraft dataset, our method again reaches the best accuracy of 94.1%. Compared with the leading result achieved by DCL, the relative accuracy gain is 1.1%, which confirms the significance of our method.

Overall, the proposed AP-CNN benefits from two aspects: 1) By establishing the pyramidal hierarchy on CNN backbones, we extract multi-scale features guided by individual attention activations, which can distill both the high-level semantic and the low-level detailed information for better classification and precise localization. 2) Conducting the ROI guided refinement (aligning features with background noise excluded by ROI guided zoom-in, and enhancing discriminative parts on features by ROI guided dropblock) on the refined-stage can also contribute to performance improvement.

4.3. Ablation Studies

We conduct ablation studies to analyze the contribution of each component. The following experiments are all conducted on the CUB-200-2011 dataset and we use ResNet50 as the backbone if not particularly mentioned.

Effect of the Pyramidal Hierarchy. We investigate the effect of constructing the pyramidal hierarchy on CNN backbones by comparing the performances obtained by the backbone, by the feature pyramid structure (FP), and by the attention pyramid structure (AP) on the VGG and ResNet network. As shown in Table 2, FP leads significant performance improvement compared to the baseline, and AP further raises the accuracy by enhancing the correlations between features. The results confirm that the pyramidal architecture with multi-level information is essential in the FGVC task.

| Method | Base Model | Accuracy (%) |
|----------|------------------|--------------------|
| Baseline | VGG19 / ResNet50 | 77.8 / 84.1 |
| FP | VGG19 / ResNet50 | 83.3 / 86.6 |
| FP + AP | VGG19 / ResNet50 | 85.4 / 87.2 |

Table 2. Comparison results on backbones with/without pyramidal hierarchy structure. FP: Feature pyramid. AP: Attention pyramid.

| Method | Accuracy (%) | mIoU (%) | Recall (%) |
|--------------|--------------|-------------|-------------|
| FPN | 86.6 | - | - |
| FPN + C | 86.8 | - | - |
| FPN + S | 86.7 | 54.9 | 73.6 |
| FPN + S + C | 86.7 | 54.9 | 74.5 |
| FPN + C + SP | 86.9 | 54.1 | 72.0 |
| FPN + S + CP | 87.2 | 56.4 | 74.0 |

Table 3. Comparison results on different ways of getting Attention. C: Channel attention. S: Spatial attention. SP: Bottom-up spatial attention pathway. CP: Bottom-up channel attention pathway.

| Erasing | Zoom-in | Guidance | Accuracy (%) |
|---------|---------|--------------|--------------------|
| - | - | ROI / Random | 87.2 / 87.2 |
| ✓ | - | ROI / Random | 87.4 / 87.1 |
| - | ✓ | ROI / Random | 87.6 / 87.3 |
| ✓ | ✓ | ROI / Random | 88.4 / 87.6 |

Table 4. Contribution of each refinement component.

Arrangement of the Attention Components. The attention union in AP-CNN consists of two parts: the spatial attention pyramid and the channel attention pathway. We conduct experiments to demonstrate their effectiveness by evaluating the classification accuracy, the mean Intersection-over-Union (mIoU) and the recall rate with ground-truth bounding boxes. Table 3 shows that using the spatial attention set or the channel attention set alone or their combination can only contribute limited improvement in classification accuracy. We add the additional bottom-up pathways to enhance the relationship between the neighboring pyramidal levels, which can be constructed by either channel attention or spatial attention. The spatial attention based bottom-up pathway will make the activation maps become similar, which is not consistent to our motivation and therefore yields poor mIoU and recall rate. The channel attention based bottom-up pathway is reasonable, as the FPN has made the features from different levels aligned in channels. The experimental results confirm that the channel attention pathway is an appropriate choice.

Contribution of the Refinement Components. As described above, our ROI guided refinement mainly consists of two operations, including the ROI guided dropblock and the ROI guided zoom-in. We conduct the ablation studies upon them, and compare the ROI guided methods with the corresponding random operations. As shown in Table 4,

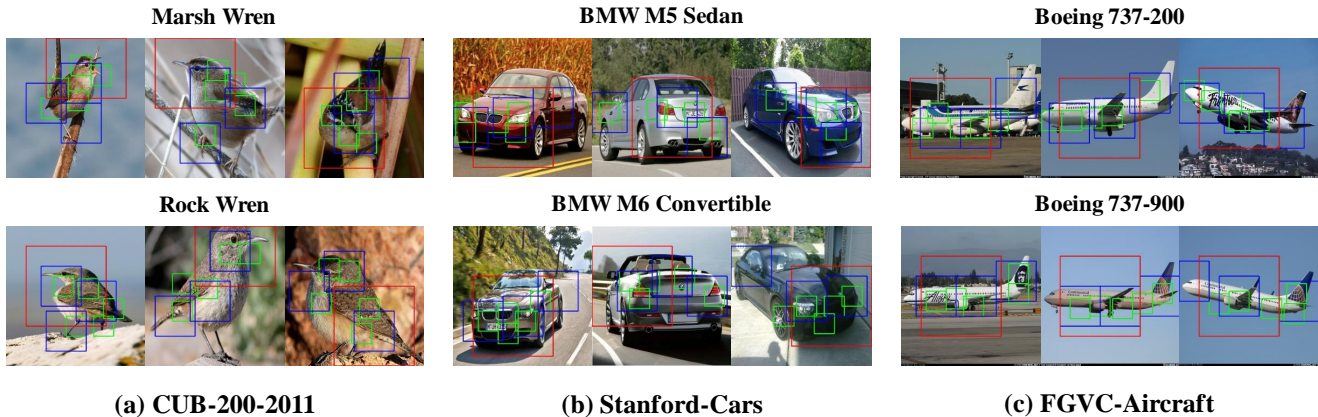


Figure 5. Visualization of ROI pyramid from two similar sub-categories on (a) CUB-200-2011, (b) Stanford-Cars, and (c) FGVC-Aircraft dataset. Rectangles with different colors are from distinctive levels in the pyramidal hierarchy.

| Refinement position | Accuracy (%) | Time cost |
|---------------------|--------------|---------------------|
| Input image | 88.1 | 527s / epoch |
| Conv1 feature | 87.9 | 525s / epoch |
| Conv2 feature | 88.1 | 458s / epoch |
| Conv3 feature | 88.4 | 389s / epoch |

Table 5. Comparison between refinement on input image and low-level features.

| Method | GFlops | Params | Accuracy (%) | | |
|----------|--------|--------|--------------|-------------|-------------|
| | | | Birds | Cars | Airs |
| Baseline | 16.48 | 25.56M | 84.1 | 91.7 | 88.5 |
| A | 19.64 | 27.96M | 86.7 | 93.1 | 91.5 |
| B | 19.64 | 27.96M | 87.2 | 93.6 | 92.2 |
| C | 31.93 | 27.96M | 88.4 | 95.4 | 94.1 |

Table 6. Detailed information of the most contributed parts. A: AP-CNN with two most contributed parts (the bottom-up pathway and the ROI guided refinement) removed, B: A + bottom-up pathway, C: B + ROI guided refinement.

ROI guided methods have advantages compared with the random ones in many aspects, and both the two refine components are effective in our refinement process.

Model Complexity. Our refinement operations can be theoretically conducted on any low-level position of the network (*e.g.*, the input and the feature maps from conv1, conv2) by sampling the ROIs into different scales. Table 5 compares the classification accuracies and the training time costs (on a single GTX 1080Ti GPU) among different refinement positions. We consequently refine the low-level features from Conv3 with the consideration of both accuracy and efficiency.

In summary, the efficiency of the proposed method can benefit from three aspects: a) the bottom-up pathway only sums two channel attention masks without additional parameters, b) the raw-stage and the refined-stage are based on the same network with shared parameters, and c) the re-

finement is conducted on low-level features rather than inputs. Table 6 shows consistent improvements yielded by the most contributed parts (the bottom-up pathway and the ROI guided refinement) on the aforementioned FGVC datasets, with limited increase of the amount of parameters and computational cost.

4.4. Visualization

Figure 5 visualizes the ROI pyramid learned by the AP-CNN. In each line, we randomly select three test images from one specific “wren” species from the CUB-200-2011 dataset, one “BMW” series from the Stanford-Cars dataset, and one “Boeing” series from the FGVC-Aircraft dataset. We use the red, blue and green boxes to denote the most activated regions, with the red ones representing the high-level ROIs with big anchor size, to the green ones representing the low-level ROIs with small anchor size. It can be intuitively observed that the localized regions are indeed informative for fine-grained classification, and the ROIs from different pyramidal levels can focus on more distinctive parts due to their particular receptive field.

5. Conclusion

In this paper, we propose an attention pyramid network for fine-grained image classification without extra annotations. This is conducted via building a pyramidal hierarchy upon CNN which consists of a top-down feature pathway and a bottom-up attention pathway to deliver both the high-level semantic and the low-level detailed information. ROI guided refinement, which enhances the discriminative local activations and aligns the features with the background noises eliminated, is conducted to further improve the performance. Experiments on CUB-Bird, Stanford-Cars, and FGVC-Aircraft demonstrate the superiority of our method.

References

- [1] Ipek Baz, Erdem Yoruk, and Mujdat Cetin. Context-aware hybrid classification system for fine-grained retail product recognition. In *IVMSP*, pages 1–5. IEEE, 2016. 1
- [2] Sean Bell, C Lawrence Zitnick, Kavita Bala, and Ross Girshick. Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks. In *CVPR*, pages 2874–2883, 2016. 2
- [3] Zhaowei Cai, Quanfu Fan, Rogerio S Feris, and Nuno Vasconcelos. A unified multi-scale deep convolutional neural network for fast object detection. In *ECCV*, pages 354–370. Springer, 2016. 2
- [4] Yuning Chai, Victor Lempitsky, and Andrew Zisserman. Symbiotic segmentation and part localization for fine-grained categorization. In *ICCV*, pages 321–328, 2013. 2
- [5] Yue Chen, Yalong Bai, Wei Zhang, and Tao Mei. Destruction and construction learning for fine-grained image recognition. In *CVPR*, pages 5157–5166, 2019. 6
- [6] Jianlong Fu, Heliang Zheng, and Tao Mei. Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *CVPR*, pages 4438–4446, 2017. 2
- [7] Yang Gao, Oscar Beijbom, Ning Zhang, and Trevor Darrell. Compact bilinear pooling. In *CVPR*, pages 317–326, 2016. 2
- [8] Weifeng Ge, Xiangru Lin, and Yizhou Yu. Weakly supervised complementary parts models for fine-grained image classification from the bottom up. In *CVPR*, pages 3034–3043, 2019. 2, 3, 6
- [9] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Dropblock: A regularization method for convolutional networks. In *NIPS*, pages 10727–10737, 2018. 5
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 3, 6
- [11] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, pages 7132–7141, 2018. 4
- [12] Shaoli Huang, Zhe Xu, Dacheng Tao, and Ya Zhang. Part-stacked cnn for fine-grained visual categorization. In *CVPR*, pages 1173–1182, 2016. 2
- [13] Justin Johnson, Andrej Karpathy, and Li Fei-Fei. Densecap: Fully convolutional localization networks for dense captioning. In *CVPR*, pages 4565–4574, 2016. 1
- [14] Shu Kong and Charless Fowlkes. Low-rank bilinear pooling for fine-grained classification. In *CVPR*, pages 365–374, 2017. 2
- [15] Tao Kong, Anbang Yao, Yurong Chen, and Fuchun Sun. Hypernet: Towards accurate region proposal generation and joint object detection. In *CVPR*, pages 845–853, 2016. 2
- [16] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *ICCV Workshops*, pages 554–561, 2013. 2, 6
- [17] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012. 3
- [18] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, pages 2117–2125, 2017. 2, 3, 5
- [19] Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. Bilinear cnn models for fine-grained visual recognition. In *ICCV*, pages 1449–1457, 2015. 2, 6
- [20] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, pages 21–37. Springer, 2016. 2
- [21] Wei Liu, Andrew Rabinovich, and Alexander C Berg. Parsenet: Looking wider to see better. *arXiv preprint arXiv:1506.04579*, 2015. 2
- [22] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015. 2
- [23] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013. 2, 6
- [24] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, pages 91–99, 2015. 4
- [25] Shaoqing Ren, Kaiming He, Ross Girshick, Xiangyu Zhang, and Jian Sun. Object detection networks on convolutional feature maps. *PAMI*, 39(7):1476–1481, 2016. 2
- [26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241. Springer, 2015. 2
- [27] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 3, 6
- [28] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 2, 6
- [29] Yaming Wang, Vlad I Morariu, and Larry S Davis. Learning a discriminative filter bank within a cnn for fine-grained recognition. In *CVPR*, pages 4148–4157, 2018. 6
- [30] Tianjun Xiao, Yichong Xu, Kuiyuan Yang, Jiaying Zhang, Yuxin Peng, and Zheng Zhang. The application of two-level attention models in deep convolutional neural network for fine-grained image classification. In *CVPR*, pages 842–850, 2015. 1
- [31] Linjie Yang, Ping Luo, Chen Change Loy, and Xiaoou Tang. A large-scale car dataset for fine-grained categorization and verification. In *CVPR*, pages 3973–3981, 2015. 1
- [32] Ze Yang, Tiange Luo, Dong Wang, Zhiqiang Hu, Jun Gao, and Liwei Wang. Learning to navigate for fine-grained classification. In *ECCV*, pages 420–435, 2018. 2, 3, 6
- [33] Chaojian Yu, Xinyi Zhao, Qi Zheng, Peng Zhang, and Xinge You. Hierarchical bilinear pooling for fine-grained visual recognition. In *ECCV*, pages 574–589, 2018. 2
- [34] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *ECCV*, pages 818–833. Springer, 2014. 2

- [35] Ning Zhang, Jeff Donahue, Ross Girshick, and Trevor Darrell. Part-based r-cnns for fine-grained category detection. In *ECCV*, pages 834–849. Springer, 2014. [2](#)
- [36] Heliang Zheng, Jianlong Fu, Tao Mei, and Jiebo Luo. Learning multi-attention convolutional neural network for fine-grained image recognition. In *ICCV*, pages 5209–5217, 2017. [2](#), [6](#)
- [37] Heliang Zheng, Jianlong Fu, Zheng-Jun Zha, and Jiebo Luo. Looking for the devil in the details: Learning trilinear attention sampling network for fine-grained image recognition. In *CVPR*, pages 5012–5021, 2019. [6](#)