# Stochastic Region Pooling: Make Attention More Expressive

Mingnan Luo, Guihua Wen,* Yang Hu,† Dan Dai, Yingxue Xu
School of Computer Science & Engineering, South China University of Technology
Panyu, Guangzhou, Guangdong, China
{csluomingnan@mail.,crghwen@,cssuperhy@mail.,csdaidan@mail.,201530381885@mail.}scut.edu.cn

## Abstract

*Global Average Pooling (GAP) is used by default on the channel-wise attention mechanism to extract channel descriptors. However, the simple global aggregation method of GAP is easy to make the channel descriptors have homogeneity, which weakens the detail distinction between feature maps, thus affecting the performance of the attention mechanism. In this work, we propose a novel method for channel-wise attention network, called Stochastic Region Pooling (SRP), which makes the channel descriptors more representative and diversity by encouraging the feature map to have more or wider important feature responses. Also, SRP is the general method for the attention mechanisms without any additional parameters or computation. It can be widely applied to attention networks without modifying the network structure. Experimental results on image recognition datasets including CIAFR-10/100, ImageNet and three Fine-grained datasets (CUB-200-2011, Stanford Cars and Stanford Dogs) show that SRP brings the significant improvements of the performance over efficient CNNs and achieves the state-of-the-art results.*

## 1. Introduction

Convolutional neural network (CNN) is an effective method to solve the computer vision tasks [21, 35, 8]. Furthermore, combining it with attention mechanisms can better solve them [1, 40, 46, 32], such as channel-wise attention networks [13, 27, 12, 42, 54]. They usually use Global Average Pooling (GAP) to squeeze the entire feature map into a descriptor [13, 14, 27, 42]. However, GAP tends to ignore the detail area in feature map with lower magnitude, which easily leads to the homogeneity of the channel descriptors. To alleviate this problem, some researchers combine channel-wise attention with spatial attention to make the attention module to pay attention to the spatial details of

the feature map [27, 29]. And some researcher even design a 3D-like attention module to extract the channel descriptors with spatial information [12]. However, these methods bring a lot of additional parameters and computational consumption. Therefore, it is expected that the attention mechanism is paid attention to the details of feature maps under the framework of the channel-wise attention mechanism but with only a little additional cost.

This paper attempts to improve the representation of descriptors extracted by GAP through providing the feature maps with higher quality. The proposed method is called as Stochastic Region Pooling (SRP), which does not brings extra parameters and computation in test phase. SRP emphasizes more local features in the convolutional layer, making the channel descriptors extracted by GAP more representative and thus making the channel-wise attention mechanism works better. In more details, it stochastically selects the region from the feature map and used GAP to obtain the region descriptor, where the descriptor is the accurate representation of the region. Subsequently, the region descriptors are used in the follow-up attention structure. In such case, the back propagation [24] will encourage these regions to have more important feature responses to represent its original entire feature map.

This paper proposes a simple method to implement SRP, named as Single Square SRP (SS-SRP), which stochastically selects a single square region from feature map to extract the descriptor. In order to consider local response of irregular shape in feature map, another method named as Multiple Squares SRP (MS-SRP) is proposed that stochastically selects multiple square regions from the feature map and then extracts the descriptor from their union regions. These two methods are illustrated in figure.1. On the other hand, in residual networks, most of attention mechanisms only act on the residual branch of residual block. In order to make the feature maps of identity branch also have more feature responses, we use SRP to extract the channel descriptors of both identity branch and residual branch, which are then applied to serve the follow-up attention structure. The main contributions are as follows:

---

*corresponding author
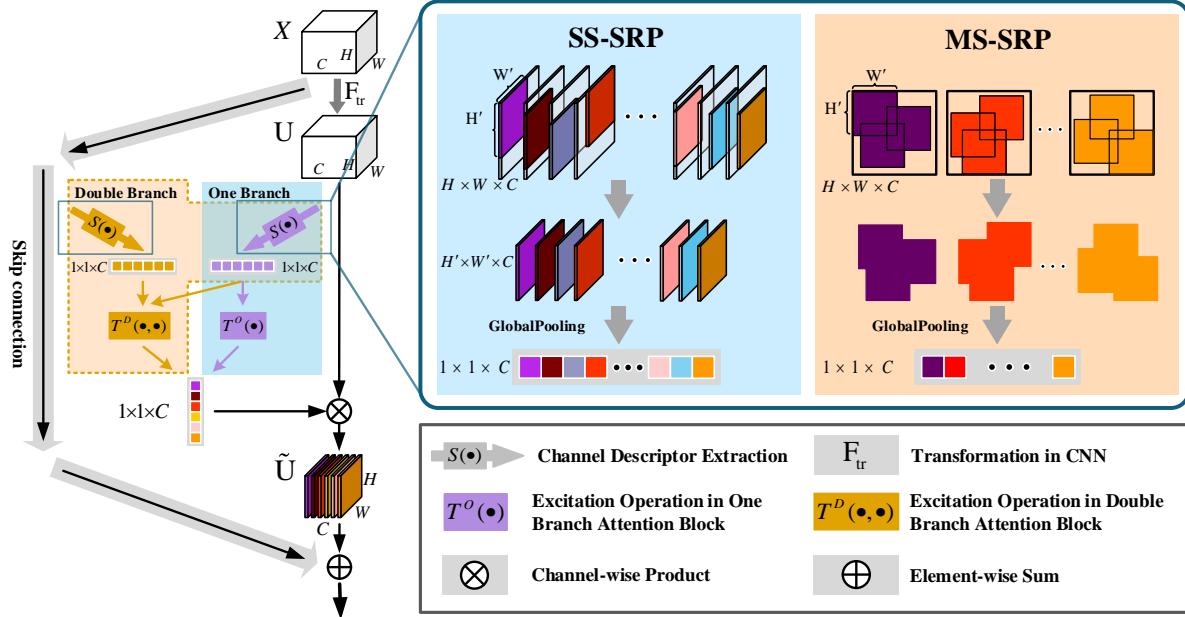†equal contribution with Guihua Wen

Figure 1: The framework of Stochastic Region Pooling (SRP) with applications to one branch or double branch attention block, where SS-SRP and MS-SRP are considered.

- A new method SRP is proposed that make more local regions in the feature map to have more important feature responses, so that the channel descriptors extracted later by GAP are more representative and robust.

- SRP obtains the significant improvement of the performance for one branch or double branch block of channel-wise attention structures without any additional parameters and computation in test phase. It can also work well with some augmentation methods to further improve the performance.

- A linear strategy is proposed to make SRP work well that gradually reduces the scale ratio of region as the depth of the layer increases.

- Experiments are conducted on serval datasets, including CIFAR-10/100, ImageNet and three Fine-grained datasets (CUB-200-2011, Stanford Cars and Stanford Dogs), which verified the effectiveness of SRP.

## 2. Related Work

**Improving representation of feature maps**. A common way to obtain high-quality feature maps is to find efficient network structures, such as [35, 38, 8, 15, 45] to extract more and better features. However, the feature maps learned by these networks are still not diverse enough. Another way

is regularization. Some regularization for channel can maintain high quality channels by removing or retraining the inefficient channels [10, 52, 5, 11]. For [10, 52], they will change the network structure. And for [5, 11], we have not found the evidence to prove that they are suitable for channel-wise attention neural networks. Other regularizations such as dropout [36], droppath [22], dropblock [6], cutout [3] can enhance the robustness of the feature by introducing randomness. And our method is closely related to Dropblock [6] which drops spatially correlated information to promote the network to reconstruct the important features from its surrounding. However, our method is aims to solve the problem that the descriptor in the channel-wise attention network contains few detailed information of the feature map, such as by promoting the feature map to have more or wider important feature responses.

**Extracting descriptors by spatial feature pooling**. The idea of spatial feature pooling was proposed by Hubel and Wiesel [16], and then Yann Lecun [25] successfully applied it to CNN. Furthermore, Spatial Pyramid Matching(SPM) [23, 44] manually designed the pooling weights to obtain spatial feature pyramid, Malinowski [31] parameterized pooling operator to learn the pooling regions, and Lee [26] combined the max pooling and the average pooling to obtains a generalized pooling function. Some researchers also use the second-order pooling even the third-order pooling instead of the first-order pooling (i.e., GAP) to collect richer statistics of the last convolution layer in

**Algorithm 1:** Stochastic Region Pooling

**Input:** Feature map:$U$, height and width of the feature map: $H$ and $W$, scale ratio:$\lambda$, the number of square regions:$M$, $mode$.

**Output:** The channel descriptors $z$ ($z_c$ is the $c^{th}$ channel descriptor of $z$).

**1  if** $mode\ != TrainingStage$ **then**

**2**  $\quad$ $\forall c: z_c = F_{sq}(u_c)$ $\qquad\qquad\qquad$ $\triangleright F_{sq}$ computes Eq. (1);

**3**  $\quad$ return $z$;

**4  end**

**5**  Calculate the height and width of square region: $H' = \lfloor \lambda H + \frac{1}{2} \rfloor, W' = \lfloor \lambda W + \frac{1}{2} \rfloor$;

**6**  Stochastically sample $M$ positions $P_{i,j}^m$ from the feature map $U$ where $1 \leq i \leq H - H' + 1, 1 \leq j \leq W - W' + 1$;

**7**  Crop $M$ square regions from the feature map $U$ with the left-top position as $P_{i,j}^m$, where the width and height of square region are $W'$ and $H'$ respectively; $\forall c: z_c = F_{sq}(u_c, P_{i,j}, H', W')$ $\qquad$ $\triangleright F_{sq}$ computes Eq. (3) or Eq. (5);

**8**  return $z$ ;

---

CNN [2, 28, 41]. Introducing randomness can also improve the performance of spatial pooling. For example, stochastic pooling [48] randomly selects the activation value based on a multinomial distribution formed by activations of each pooling region to regularize the network, and S3Pool [49] randomly picks feature map's rows and columns and then performs the max pooling operation to implicitly introduce data augmentation. However, using the above method to extract the channel descriptor will bring a lot of extra consumption, or it will still not make the descriptor's representation stronger. We select a simplest way, which takes GAP to extract the descriptor because it is widely used, does not bring any extra parameters, and has the potential to get global spatial information.

**Methods of spatial pooling in attention mechanism**. Channel-wise attention mechanisms have developed rapidly in recent years [32, 1, 27, 13, 12, 42] and the channel descriptors are crucial to them. In order to extract more representative descriptors, CBAM [42] combines the output of the global max pooling and the global average pooling as the pooling method, and GEnet [12] uses the depth-wise convolution with large kernel to replace GAP. However, the global max pooling in CBAM is prone to network overfitting [48], and CBAM also cannot enhance the channel descriptor with more spatial details. Besides, GENet will brings a lot of additional parameters or computation. Different from them, SRP is a training method that can encourage the descriptors to have more information about the feature map details. The reason why SRP uses GAP instead of the above methods to extract the descriptor is not only because GAP is simple and does not bring any additional parameters, but also that GAP is widely used for attention networks. This enables SRP to be conveniently used on these networks without modifying the network structure.

## 3. Stochastic Region Pooling

Many channel-wise attention block applied GAP operation to obtain the descriptor of feature map. Formally, given the feature maps $U = [u_1, u_2, \cdots, u_c] \in \mathbb{R}^{H \times W \times C}$ in the convolutional layer and the function $S(\cdot)$ squeezes the global spatial information into a channel descriptor $z = [z_1, z_2, \cdots, z_C] \in R^C$, the $c^{th}$ channel descriptor of $z$ can be calculated in GAP by

$$z_c = S(u_c) = \frac{1}{H \times W} \sum_{i=1}^{H} \sum_{j=1}^{W} u_c(i,j). \qquad (1)$$

From the Eq.1, it can be concluded that GAP regards each position of the space to make the same contribution, even if the elements of some local regions have the low magnitude [48]. This will weaken the details of the feature map and easily leads to high similarity between descriptors. Here we propose Stochastic Region Pooling (SRP) method that stochastically selects the region from the feature map instead of the whole map to extract descriptors during the training stage. SRP is presented as Algorithm 1.

**SS-SRP** is a simple method to be implemented, which stochastically selects a single square region from map. Supposing that the scale ratio $\lambda$ controls the size of square region, the width and height of the square region can be formulated as follows,

$$H' = \lfloor \lambda H + \frac{1}{2} \rfloor, W' = \lfloor \lambda W + \frac{1}{2} \rfloor. \qquad (2)$$

For each feature map, we stochastically select a position $P(a,b)$ as the upper left corner of the square region $R \in \mathbb{R}^{H' \times W'}$, where the spatial position $P(a,b)$ subject to $1 \leq i \leq H - H' + 1, 1 \leq j \leq W - W' + 1$. Now we use the average pooling as the squeeze operator $S(\cdot)$ to extract the descriptors $z$, and the $c^{th}$ channel descriptor of $z$ can be

calculated by

$$z_c = S(u_c) = \frac{1}{H' \times W'} \sum_{i=a}^{a+H'-1} \sum_{j=b}^{b+W'-1} u_c(i,j). \quad (3)$$

The module of SS-SRP can be seen in the Figure 1.

**MS-SRP** is applied to consider the non-regular shape of the local region in feature map, which stochastically selects multiple square regions from the feature map and then extracts the descriptor from their union regions. Suppose that we stochastically choose $M$ square regions from the feature map, defined as $R \in \mathbb{R}^{M \times H \times W}$, the target region we want is their union area $R^*$.

Let $\Omega^*$ be the set of all the points in $R^*$ and $\Omega_m$ be the set of all points in the $m^{th}$ square region $R_m$, we have

$$\Omega^* = \bigcup_{m=1}^{M} \Omega_m = \bigcup_{m=1}^{M} \left\{ (x,y) \mid (x,y) \in \Omega_m \right\}. \quad (4)$$

The global average pooling is used as we did before to squeeze the regional spatial information. Thus the $c^{th}$ channel descriptor of $z$ can be calculated by

$$z_c = S(u_c, \Omega^*) = \frac{1}{|\Omega^*|} \sum_{(i,j) \in \Omega^*} u_c(i,j), \quad (5)$$

where $|\Omega^*|$ is the number of elements in $\Omega^*$,i.e. the number of points in region $R^*$. The module of MS-SRP can be seen in the Figure 1.

**One branch or double branch attention block.** After calculating the channel descriptors, a one branch attention block applies an excitation operation $T^O(\cdot)$ to obtain the relationship $\alpha \in \mathbb{R}^C$ between the channels of the residual branch, which is $\alpha = T^O(z^r)$ where $z^r$ is the channel descriptor of the residual branch. And a double branch attention block utilizes $T^D(\cdot, \cdot)$ to compute the relationship among channels of the residual branch, which is $\alpha = T^D(z^{id}, z^r)$, where $z^{id}$ and $z^r$ are the channel descriptors of the identity branch and the residual branch respectively.

For the one branch attention block, we use two fully connected(FC) layers as the function $T^O(z^r)$ as described in SENet [13]. For the double branch attention block, we fold two branch's descriptors and use convolution $3 \times 3$ to model their relationship as the function $T^D(z^{id}, z^r)$ as described in CMPE-SENet [14].

Finally, these recalibrated feature maps $\widetilde{U} \in \mathbb{R}^{H \times W \times C}$ can be calculated by $\widetilde{U} = \alpha \cdot U$, where $\cdot$ is the element-wise multiplication.

**Scheduled SRP.** The neurons in the shallow layers of the network have smaller receptive field. In order to maintain a majority of responses in the region selected by SRP in the shallow layers, we gradually reduce $\lambda$ from 1 to the smaller value as the depth of the layer increases, instead of setting $\lambda$

| Datasets | #Class | #Train | #Test |
|----------|--------|--------|-------|
| CUB-200-2011 | 200 | 5,994 | 5,794 |
| Stanford Cars | 196 | 8,144 | 8,041 |
| Stanford Dogs | 120 | 12,000 | 8,580 |

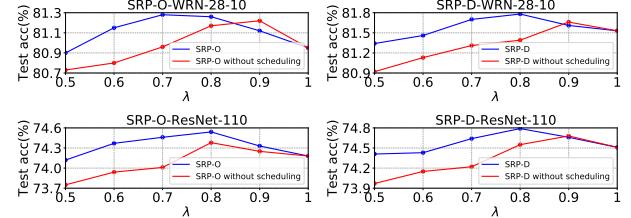Table 1: Statistics of three common Fine-grained datasets.



Figure 2: Testing acc (%) of SS-SRP (with or without scheduled) applied to One or Double branch block on CIFAR-100 data with the different $\lambda$.
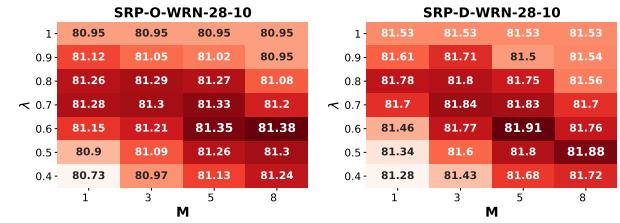


Figure 3: Test acc (%) of MS-SRP applied to One or Double Branch block on CIFAR-100 dataset under different hyperparameters $(\lambda, M)$.

to a fixed value. In our experiments, we use linear strategy to reduce $\lambda$, which is inspired by ScheduledDropPath [55], but ScheduledDropPath changes its parameters over training time.

## 4. Experiments

Some experiments are conducted to validate the proposed method on the CIFAR-10/100[20], ImageNet[33], and three fine-grained datasets (CUB-200-2011[39], Stanford Dogs[17] and Stanford Cars[19]). For experiments on CIFAR and fine-grained datasets, we report the average accuracy by running five times, while on ImageNet, we report the average accuracy by doing three times due to the limitation of computational resources. In the following subsections, SRP-O indicates that SRP is applied to the one branch attention block, and SRP-D indicates that SRP is applied to the double branch attention block.

## 4.1. Experiment Settings

**CIFAR.** Following [9, 47, 43], we use stochastic gradient descent (SGD) with 0.9 Nesterov momentum and batchsize of 128. The learning rate is set to be 0.1, which is then divided by 10 at epochs 100, 150 for ResNet, divided by 5 at epochs 60,120,160 for WRN, and divided by 10 at epochs 150,225 for ResNeXt. We train the model by 200 epochs for ResNet and WRN, and 300 epochs for ResNeXt. The weight decay is 0.0001 for ResNet, 0.0005 for WRN and ResNeXt. we use the standard data augmentation (translation/mirroring) for the training sets.

**ImageNet.** The ILSVRC 2012 contains 1.2 million training images and 50K validation images with 1K classes. We adopt the standard data augmentation for the training sets, which randomly samples a 224×224 crop from the original images or their horizontal flip, and applies a single-crop with the size 224×224 at testing stage. We train our models for 100 epochs and drop the learning rate by 0.1 at epoch 30, 60, and 90, and use SGD with the mini-batch size of 256 on 4 GPUs (64 each). The weight decay is 0.0001 and Nesterov momentum is 0.9. In experiments, we report the classification accuracy on the validation set.

**Fine-grained Datasets.** We conduct experiments on three fine-grained datasets, including CUB-200-2011,Stanford Dogs and Stanford Cars. The detailed statistics of each dataset are shown in Table 1. For all fine-grained datasets, we resize the input images to 512×512 and randomly crop the smaller images with 448×448 from it, and then generates the horizontal flip of the cropped images for training. At the testing stage, we only use a single cropped image with 448×448 from the input image which have been resized to 512×512. We fine-tune networks (pre-trained on ImageNet) using SGD with the batch size of 16, momentum of 0.9 and weight decay of 0.00001. For all fine-grained datasets, we train the networks for 90 epochs. The learning rate begins with 0.001 and then divided by 10 at epoch 30 and 60.

## 4.2. Impact of Hyper-parameters

In order to demonstrate the influence of two hyper-parameters (scale ratio $\lambda$ and the square regions number $M$) on the performance of our model, experiments are conducted on CIFAR-100, where SRP applied on attention networks with different hyper-parameter settings.

**SS-SRP**. It only randomly selects one square region, which becomes standard mothod when we set $\lambda = 1$. It can be seen from figure 2 that an appropriate $\lambda$ can improve the network performance and the lower $\lambda$ will result in the poor results. In the following experiments of SS-SRP, we use $\lambda = 0.8$ unless specified elsewhere, because SS-SRP obtains the better results at $\lambda = 0.8$. From figure 2, it can be also observed that SRP with the fixed scale ratio $\lambda$ can effectively improve the network performance, but the sched-

| Model | depth | params | C10 | C100 |
|---|---|---|---|---|
| FractalNet [22] | 21 | 38.6M | 95.40 | 76.27 |
| WRN-28-10 [47] | 28 | 36.5M | 96.00 | 80.75 |
| ResNeXt-29(8x64d) [43] | 29 | 34.4M | 96.35 | 82.23 |
| ResNeXt-29(16x64d) [43] | 29 | 68.1M | 96.42 | 82.69 |
| DenseNet(k=24) [15] | 100 | 27.2M | 96.26 | 80.75 |
| DenseNet-BC(k=40) [15] | 190 | 25.6M | 96.54 | 82.82 |
| PyramidNet(bottleneck,$\alpha = 270$) [7] | 272 | 27.0M | 96.52 | 82.99 |
| *mixup* [50], WRN-28-10 [47] | 28 | 36.5M | 97.30 | 82.50 |
| *mixup* [50], DenseNet-BC(k=40) [15] | 190 | 25.6M | 97.30 | 83.20 |
| *mixup* [50], SE-WRN-28-10 [13] | 28 | 36.8M | 97.32 | 83.23 |
| **SS-SRP-O**, WRN-28-10 | 28 | 36.8M | 96.28 | 81.38 |
| **SS-SRP-O**, ResNeXt-29(8x64d) | 29 | 34.9M | 96.52 | 82.59 |
| **SS-SRP-D**, WRN-16-8 | 16 | 11.1M | 96.02 | 80.89 |
| **SS-SRP-D**, WRN-28-10 | 28 | 36.9M | 96.50 | 81.78 |
| **MS-SRP-O**, WRN-28-10 | 28 | 36.8M | 96.34 | 81.38 |
| **MS-SRP-O**, ResNeXt-29(8x64d) | 29 | 34.9M | 96.53 | **82.62** |
| **MS-SRP-D**, WRN-16-8 | 16 | 11.1M | 96.10 | 80.98 |
| **MS-SRP-D**, WRN-28-10 | 28 | 36.9M | 96.61 | 81.91 |
| **MS-SRP-O**, *mixup*, WRN-28-10 | 28 | 36.8M | 97.48 | 84.08 |
| **MS-SRP-D**, *mixup*, WRN-16-8 | 16 | 11.1M | 96.84 | 82.71 |
| **MS-SRP-D**, *mixup*, WRN-28-10 | 28 | 36.9M | 97.56 | 84.12 |

Table 2: Comparison of test accuracy (%) with different methods on the CIFAR-10 and CIFAR-100. The best results are highlighted in red, and the best records of our models are in **bold**. Combined with the augmentation method of *mixup* [50], SRP can challenge state-of-the-art results.

| Model | params | top-1 | top-5 |
|---|---|---|---|
| CliqueNet-S3 [45] | 14.4M | 75.95 | 92.85 |
| ResNet-50 [8] | 25.6M | 75.30 | 92.20 |
| ResNet-101 [8] | 44.6M | 76.40 | 92.90 |
| SE-ResNet-50 [13] | 28.1M | 76.71 | 93.38 |
| ResNet-152 [8] | 28.1M | 77.00 | 93.30 |
| DenseNet-201 [15] | 20.2M | 77.42 | 93.66 |
| SE-ResNet-101 [13] | 49.4M | 77.62 | 93.93 |
| CBAM-ResNet-50 [42] | 25.9M | 77.34 | 93.69 |
| GE-$\theta^+$-ResNet-50 [12] | 33.7M | **78.12** | 94.20 |
| **SS-SRP-O**-ResNet-50 | 28.1M | 77.43 | 93.81 |
| **MS-SRP-O**-ResNet-50 | 28.1M | 77.58 | 93.88 |
| **SS-SRP-D**-ResNet-50 | 29.2M | 77.94 | 94.35 |
| **MS-SRP-D**-ResNet-50 | 29.2M | 78.09 | **94.40** |

Table 3: Comparison of test accuracy (%) between SRP and other different methods on the large ImageNet, where a single crop method is applied.

uled SRP makes the network work better. In all following experiments of SRP, scheduled SRP is used.

**MS-SRP**. In MS-SRP, there are two hyper-parameters that affects the network performance, i.e., the scale ratio $\lambda$ and the square number $M$. By observing the results in figure 3, MS-SRP outperforms standard method SRP (i.e., $\lambda = 1$) and SS-SRP within a wider range of hyper-parameters. Since $\lambda = 0.6$ and $M = 5$ can obtain the better results, we use these hyper-parameters for all experiments of MS-SRP unless stated elsewhere.

| Model | Anno. | 1-Stage | Acc. |
|---|---|---|---|
| DVAN [53] | × | × | 79.0 |
| Part-RCNN [51] | ✓ | × | 81.6 |
| PA-CNN [18] | ✓ | ✓ | 82.8 |
| RAN [40] | × | × | 82.8 |
| FCAN [30] | ✓ | ✓ | 84.7 |
| RACNN [4] | × | × | 85.3 |
| VGG-19 [35] | × | ✓ | 77.8 |
| ResNet-50 [8] | × | ✓ | 81.7 |
| DenseNet-161 [15] | × | ✓ | 84.2 |
| FCAN [30] | × | ✓ | 84.3 |
| ResNet-101 [8] | × | ✓ | 84.5 |
| **SS-SRP-D**-ResNet50 | × | ✓ | 84.9 |
| **MS-SRP-D**-ResNet50 | × | ✓ | **85.6** |

(a) CUB-200-2011.

| Model | Anno. | 1-Stage | Acc. |
|---|---|---|---|
| DVAN [53] | × | × | 87.1 |
| RAN [40] | × | × | 91.0 |
| FCAN [30] | ✓ | ✓ | 91.3 |
| RACNN [4] | × | × | 92.5 |
| PA-CNN [18] | ✓ | ✓ | **92.8** |
| VGG-19 [35] | × | ✓ | 84.9 |
| FCAN [30] | × | ✓ | 89.1 |
| ResNet-50 [8] | × | ✓ | 89.8 |
| DenseNet-161 [15] | × | ✓ | 91.8 |
| ResNet-110 [8] | × | ✓ | 91.9 |
| **SS-SRP-D**-ResNet50 | × | ✓ | 92.3 |
| **MS-SRP-D**-ResNet50 | × | ✓ | **92.8** |

(b) Stanford Cars.

| Model | Anno. | 1-Stage | Acc. |
|---|---|---|---|
| DVAN [53] | × | × | 81.5 |
| RAN [40] | × | × | 83.1 |
| VGG-16 [35] | × | ✓ | 76.7 |
| ResNet-50 [8] | × | ✓ | 81.1 |
| DenseNet-161 [15] | × | ✓ | 81.2 |
| FCAN [30] | × | ✓ | 84.2 |
| MAMC-ResNet50 [37] | × | ✓ | 84.8 |
| ResNet-101 [8] | × | ✓ | 84.9 |
| **SS-SRP-D**-ResNet50 | × | ✓ | 85.9 |
| **MS-SRP-D**-ResNet50 | × | ✓ | **86.3** |

(c) Stanford Dogs.

Table 4: Comparison results on three Fine-grained datasets including CUB-200-2011, Stanford Cars and Stanford Dogs. "Anno." represents using extra annotation in training. "1-Stage" represents whether the training can be done in one stage. "Acc." represents the test set accuracy (%)
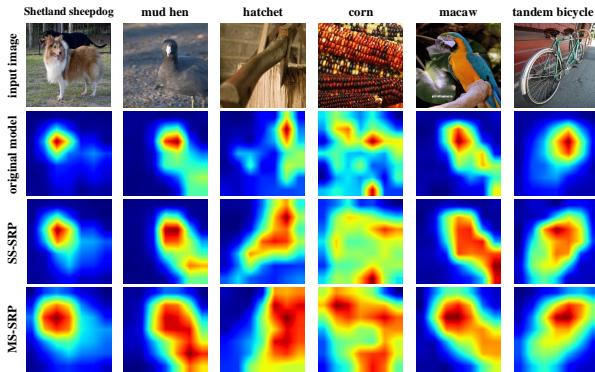


Figure 4: In **ImageNet**, the Grad-CAM [34] visualization for double branch block of attention ResNet50 model trained without SRP and trained with SRP. Best viewed in color.

## 4.3. CIFAR Classification

Table 2 presents the results of SRP and compared state-of-the-art CNN architectures on CIFAR. It can be observed that SRP method consistently achieve the better effective performance when applied it to other networks. Furthermore, the small network trained with SRP can achieve the comparable accuracy to some larger models, such as MS-SRP-O-ResNext-29 (34.9M) vs ResNext-29 (68.1M) [43], MS-SRP-D-WRN-16-8 (11.1M) vs DenseNet (27.2M) [15] or WRN-28-10 (36.5M) [47]. On the other hand, SRP may further improve network performance when augmentation methods are applied, such as *mixup* [50]. Combined with *mixup* in the training stage, MS-SRP-D surpasses all the comparison methods, and can challenge state-of-the-art results.

## 4.4. ImageNet Classification

We next investigate the effectiveness of SRP on large dataset, the ILSVRC 2012 dataset.

**Comparison with state-of-the-arts CNNs.** It can be observed from Table 3 that SRP can still improve the network performance effectively on the large data sets and achieved very competitive accuracy. For example, it exceeds all the methods in terms of top-5 accuracy, while it is much close to the state-of-the-art method by top-1. For SE-ResNet-50, the accuracy improvement of SRP is 0.87% on top-1 and 0.5% on top-5. And the best result of SRP surpasses the basic ResNet-50 and SE-ResNet-50 more than 2% and 1% respectively, by both top-1 and top-5. Compared with the GE-$\theta^+$ [12], SRP can achieve comparable or better performance but obviously uses fewer parameters. Moreover, compared with the counterparts, SRP can achieve higher accuracy by top-5 while their accuracy are similar by top-1. This is because SRP promotes the feature maps to contain more features response, making the network have better generalization ability.

**SRP learns more and wider regions.** The model trained with SRP will promote more local feature responses in the convolutional layer, making the network focus on the more and wider regions. Here, we use Grad-CAM [34] to visualize the importance of the spatial position in the convolutional layer. The visualization results of SRP networks (SS-SRP-D-ResNet50, MS-SRP-D-ResNet50) and the baseline (ResNet50 with double branch attention block) can been seen in Figure 4. It can be clearly observed that the Grad-CAM masks of SRP cover the target object better and wider than the baseline. That is, the model trained with SRP tends to focus on several spatially distributed regions, and aggregate information from multiple or wider regions.
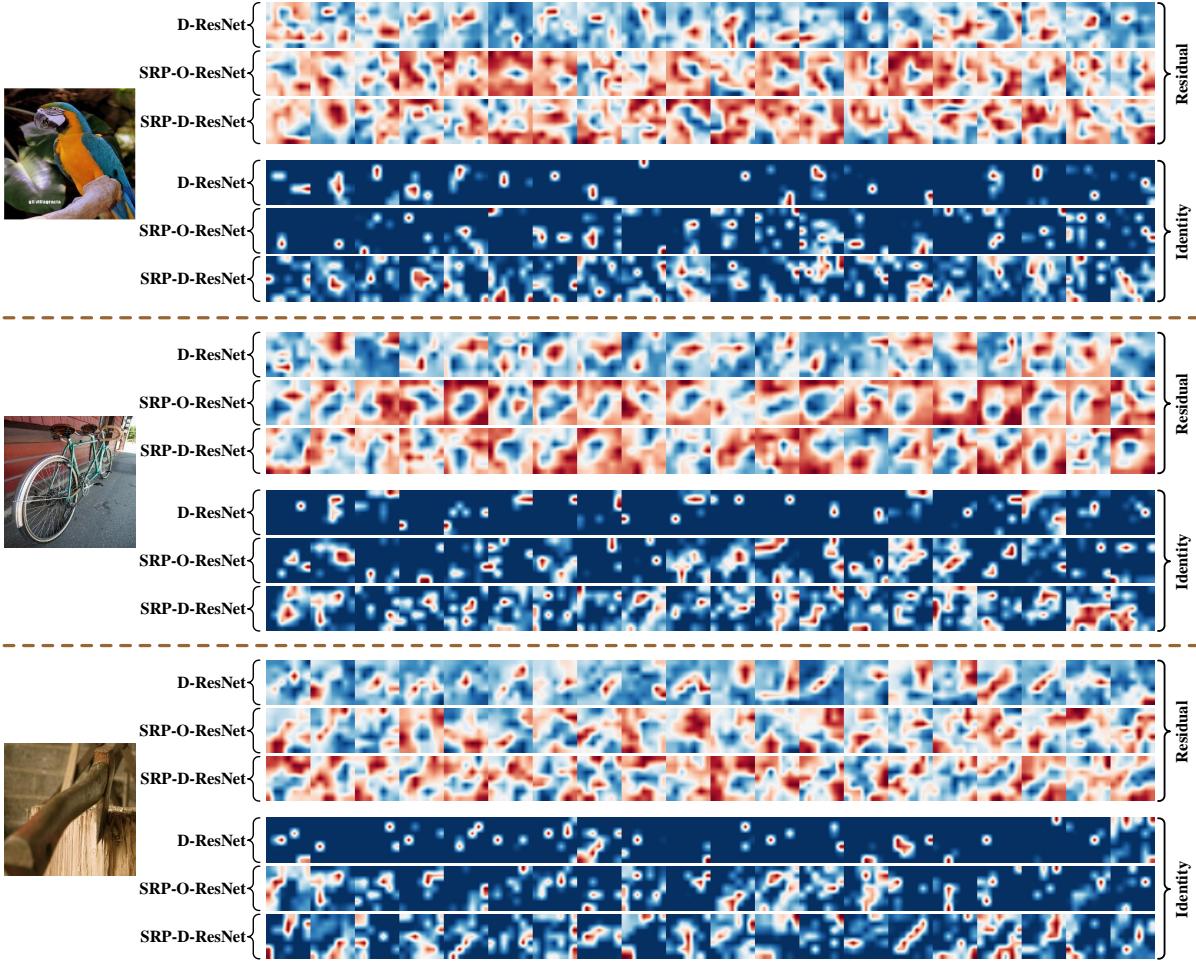
Figure 5: The feature maps of ResNet-50 with different mode. D-ResNet means ResNet-50 with double branch attention block. SRP-O-ResNet and SRP-D-ResNet means ResNet-50 with one or double branch attention block and train with MS-SRP. These feature maps are from the residual branch or identity branch of the 14th block of the network. We only display the first 20 feature maps. Best viewed in color.

## 4.5. Fine-Grained Classification

The difficulty of fine-grained classification tasks is that even with different categories of objects, they are still very similar. This requires the network to have the ability to learn multiple and more accurate region features. In this section, we investigate the performance of SRP on fine-grained classification tasks.

We first analyze the results on CUB-200-2011 dataset, as shown in Table 4a. Our method achieves strong performance with ResNet-50, which surpasses some deeper or larger network such as DenseNet-161 [15] or ResNet-110 [8] more than 1%. Also, compared with the method using extra annotation (FCAN [30]), the method using multiple training stage (Part-RCNN [51]), and the method using both extra annotation and multi-stage (RACNN [4]), our

method outperforms them by 0.9%, 4.0% and 0.3% respectively.

Our method also obtains the good performance on the Stanford Cars and Stanford Dogs, as shown in both Table 4b and Table 4c. On Stanford Cars, SRP outperforms all the comparison methods, except PA-CNN that uses extra annotation. On Stanfor Dogs, SRP surpasses the best result of other methods about 1.4% in Table 4c. Also, SRP outperforms its deeper or larger counterparts such as ResNet-110 and DenseNet-161 by about 1.0% on Stanford Cars and 3.3% on Stanford Dogs averagely. Furthermore, on these two datasets, SRP exceeds the efficient methods like DVAN [53] and RAN [40] by about 5.2% and 2.5% respectively, while DVAN and RAN both use extra annotation and multi-stage.

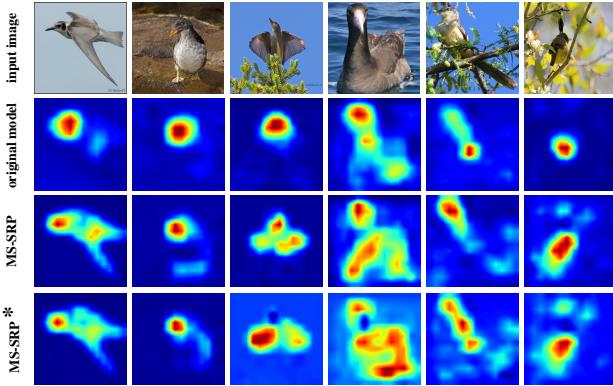These facts indicate that SRP can obtain the significant

Figure 6: In **Fine-Grained** dataset, the Grad-CAM [34] visualization datasets for double branch block of attention ResNet50 model trained without SRP, trained with MS-SRP($M = 5, \lambda = 0.6$), and trained with MS-SRP*($M = 5, \lambda = 0.2$). It indicates that SRP can promote the network to learn more detailed features of object, while SRP with too smal $\lambda$ will make the network pay attention to some unimportant fragment regions. Best viewed in color.

improvement on fine-grained datasets, which even challenge the state-of-the-art results. It is worth mentioning that our method is a general method for the attention mechanisms. It may further improve the classification accuracy by combining with the better base network or some methods which are specifically for the fine-grained image classification.

## 5. Discussion and Analysis

**Scheduled SRP.** The results in figure 2 shows that the scheduled SRP is superior to SRP with fixed $\lambda$ within a wide range of hyper-parameters. This indicates that the scheduled SRP is effective and practical. The possible reason is that the receptive field of shallow neurons in CNN are small. When SRP takes the fixed $\lambda$, the attention mechanism will receive a over-fragmented information about the feature map in shallow layer, which will in turn disturb the learning of the attention mechanism. Scheduled SRP avoided this problem by reducing $\lambda$ from 1 to the target value as the depth of layer increases.

**Feature maps analysis.** The purpose of SRP is to improve the representation of channel descriptors by increasing or widening the important responses in the feature map. Hence we output some feature maps of SRP-O-ResNet, SRP-D-ResNet and D-ResNet to analyze the effects of SRP visually, as shown in Figure 5. For SRP-O-ResNet, SRP only acts on the residual branch. For SRP-D-ResNet, SRP acts on both the residual and identity branchs. For D-ResNet, SRP is not used. It can be seen that in the residual branch, the feature map of SRP-O-ResNet and SRP-D-ResNet con-

tains more and wider responses than D-ResNet (double branch attention of ResNet-50). In the identity branch, the feature map of SRP-D-ResNet contains more and wider responses than both SRP-O-ResNet and D-ResNet. These facts indicate that due to the effect of SRP, the feature maps from corresponding branch will have more and wider object responses.

**SRP on fine-grained Recognition.** In order to investigate the reason why SRP is more effective in fine-grained classification tasks, we use Grad-CAM to compare the visualization results of network train with or without SRP in CUB-200-2011. As we can see in Figure 6, SRP can promote the network to focus on more details of the bird, such as the tip of the wing, the tail or the claw, while the network without SRP mostly tends to focus on one region. This indicates that SRP can promote the network to learn more detailed features of object, which becomes the key to the success of SRP in fine-grained classification tasks. It can be also observed that the network will pay attention to some unimportant fragment regions when $\lambda$ of SRP is too small. We conjecture that a too small $\lambda$ will cause SRP to obtain the over-fragmented information about the feature map, which may affect the performance of the attention mechanism.

## 6. Conclusion

In this paper, we propose a new method called Stochastic Region Pooling(SRP) for channel-wise attention networks. SRP stochastically selects the region from the feature map to extract descriptor in the training stage, promoting convolutional layer have more important feature responses, and making the network to focus on more and wider spatially distributed regions. Besides, SRP is the general method that can be applied to attention network without modifying the network structure and increasing any additional parameters. Our experiments show that the channel-wise attention network trained with SRP can achieve significant performance improvements on various image classification tasks and challenge the state-of-the-art methods. It is also proved that gradually decreasing the scale ratio of region as the depth of layer leads to the better accuracy.

## References

[1] L. Chen, H. Zhang, J. Xiao, L. Nie, J. Shao, W. Liu, and T.-S. Chua. Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning. In *CVPR*, 2017.

[2] Y. Cui, F. Zhou, J. Wang, X. Liu, Y. Lin, and S. J. Belongie. Kernel pooling for convolutional neural networks. In *CVPR*, 2017.

[3] T. DeVries and G. W. Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.

[4] J. Fu, H. Zheng, and T. Mei. Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *CVPR*, 2017.

[5] X. Gao, Y. Zhao, ukasz Dudziak, R. Mullins, and C. zhong Xu. Dynamic channel pruning: Feature boosting and suppression. In *ICLR*, 2019.

[6] G. Ghiasi, T.-Y. Lin, and Q. V. Le. Dropblock: A regularization method for convolutional networks. In *NIPS*, 2018.

[7] D. Han, J. Kim, and J. Kim. Deep pyramidal residual networks. In *CVPR*, pages 5927–5935, 2017.

[8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

[9] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, 2016.

[10] Y. He, X. Zhang, and J. Sun. Channel pruning for accelerating very deep neural networks. In *ICCV*, pages 1389–1397, 2017.

[11] S. Hou and Z. Wang. Weighted channel dropout for regularization of deep convolutional neural network. In *AAAI*, 2019.

[12] J. Hu, L. Shen, S. Albanie, G. Sun, and A. Vedaldi. Gather-excite: Exploiting feature context in convolutional neural networks. In *NIPS*, 2018.

[13] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *CVPR*, 2018.

[14] Y. Hu, G. Wen, M. Luo, and D. Dai. Competitive inner-imaging squeeze and excitation for residual network. *arXiv preprint arXiv:1807.08920*, 2018.

[15] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *CVPR*, volume 1, page 3, 2017.

[16] D. H. Hubel and T. N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of physiology*, 1962.

[17] A. Khosla, N. Jayadevaprakash, B. Yao, and F.-F. Li. Novel dataset for fine-grained image categorization: Stanford dogs. In *CVPR Workshop*, 2011.

[18] J. Krause, H. Jin, J. Yang, and L. Fei-Fei. Fine-grained recognition without part annotations. In *CVPR*, 2015.

[19] J. Krause, M. Stark, J. Deng, and L. Fei-Fei. 3d object representations for fine-grained categorization. In *ICCV Workshop*, 2013.

[20] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, 2009.

[21] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.

[22] G. Larsson, M. Maire, and G. Shakhnarovich. Fractalnet: Ultra-deep neural networks without residuals. In *ICLR*, 2017.

[23] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*, 2006.

[24] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.

[25] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

[26] C.-Y. Lee, P. W. Gallagher, and Z. Tu. Generalizing pooling functions in convolutional neural networks: Mixed, gated, and tree. In *Artificial Intelligence and Statistics*, 2016.

[27] W. Li, X. Zhu, and S. Gong. Harmonious attention network for person re-identification. In *CVPR*, 2018.

[28] T.-Y. Lin, A. RoyChowdhury, and S. Maji. Bilinear cnn models for fine-grained visual recognition. In *ICCV*, 2015.

[29] D. Linsley, D. Scheibler, S. Eberhardt, and T. Serre. Global-and-local attention networks for visual recognition. *arXiv preprint arXiv:1805.08819*, 2018.

[30] X. Liu, T. Xia, J. Wang, Y. Yang, F. Zhou, and Y. Lin. Fully convolutional attention networks for fine-grained recognition. *arXiv preprint arXiv:1603.06765*, 2016.

[31] M. Malinowski and M. Fritz. Learning smooth pooling regions for visual recognition. In *BMVC*, 2013.

[32] T. V. Nguyen, Q. Zhao, and S. Yan. Attentive systems: A survey. *IJCV*, 2018.

[33] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 2015.

[34] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017.

[35] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2014.

[36] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.

[37] M. Sun, Y. Yuan, F. Zhou, and E. Ding. Multi-attention multi-class constraint for fine-grained image recognition. In *ECCV*, 2018.

[38] 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*, pages 1–9, 2015.

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

[40] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang. Residual attention network for image classification. In *CVPR*, 2017.

[41] Q. Wang, Z. Gao, J. Xie, W. Zuo, and P. Li. Global gated mixture of second-order pooling for improving deep convolutional neural networks. In *NIPS*, 2018.

[42] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon. Cbam: Convolutional block attention module. In *ECCV*, 2018.

[43] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017.

[44] J. Yang, K. Yu, Y. Gong, and T. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *CVPR*, 2009.

[45] Y. Yang, Z. Zhong, T. Shen, and Z. Lin. Convolutional neural networks with alternately updated clique. In *CVPR*, 2018.

[46] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola. Stacked attention networks for image question answering. In *CVPR*, pages 21–29, 2016.

[47] S. Zagoruyko and N. Komodakis. Wide residual networks. In *BMVC*, 2016.

[48] M. D. Zeiler and R. Fergus. Stochastic pooling for regularization of deep convolutional neural networks. *arXiv preprint arXiv:1301.3557*, 2013.

[49] S. Zhai, H. Wu, A. Kumar, Y. Cheng, Y. Lu, Z. Zhang, and R. S. Feris. S3pool: Pooling with stochastic spatial sampling. In *CVPR*, 2017.

[50] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.

[51] N. Zhang, J. Donahue, R. Girshick, and T. Darrell. Part-based r-cnns for fine-grained category detection. In *ECCV*, 2014.

[52] X. Zhang, H. Xiong, W. Zhou, W. Lin, and Q. Tian. Picking deep filter responses for fine-grained image recognition. In *CVPR*, pages 1134–1142, 2016.

[53] B. Zhao, X. Wu, J. Feng, Q. Peng, and S. Yan. Diversified visual attention networks for fine-grained object classification. *IEEE Transactions on Multimedia*, 19(6):1245–1256, 2017.

[54] L. Zhu, S. Zhan, and H. Zhang. Stacked u-shape networks with channel-wise attention for image super-resolution. *Neurocomputing*, 2019.

[55] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le. Learning transferable architectures for scalable image recognition. In *CVPR*, 2018.

# Appendix

## 1. Area ratio of region in SRP

The area ratio of region selected by SRP is different under different hyper-parameters$(M, \lambda)$; the area ratio is equal to the area of region divided by the area of the feature map. Figure 7 shows the curve of the area ratio of region selected by SRP on different depths of network (ResNet-110, WRN-28-10 and ResNet-50). Under the sampe depth, the area ratio of region takes value in a large range in the MS-SRP, but is a fixed value in the SS-SRP. We speculate that this smoothing is one of the reasons that make MS-SRP better than SS-SRP and allows SRP steadily promote network learning.
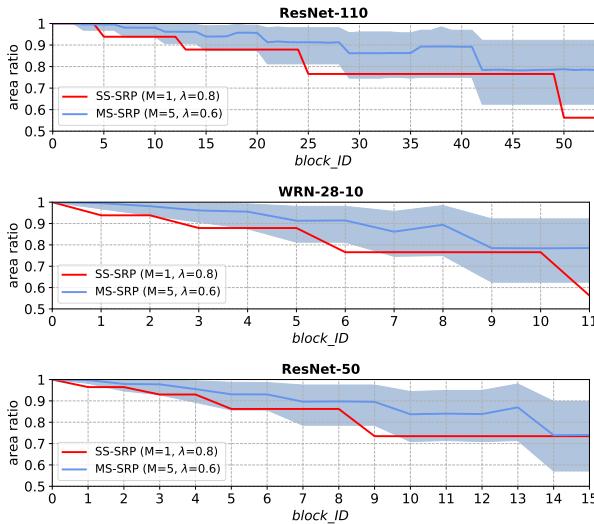


Figure 7: The curve of region area ratio on different depth of netwrok blocks, where the region is seleted by SRP. The red line is the value of the region area ratio in SS-SRP. The blue line is the mean of the region area ratio in the MS-SRP, and the value of area ratio in MS-SRP has a probability of 95% on the blue shadow. Best viewed in color.

## 2. Feature maps of network block at different depths

Figure 8 shows the visualization of feature maps from the network which trained with or without SRP at different depth. We can observe that the deeper the network block, the more obvious the SRP characteristic that makes the feature map contains more and wider feature responses. At the same time, we found that too small $\lambda$ values tend to generate more but messy responses in deep layer.
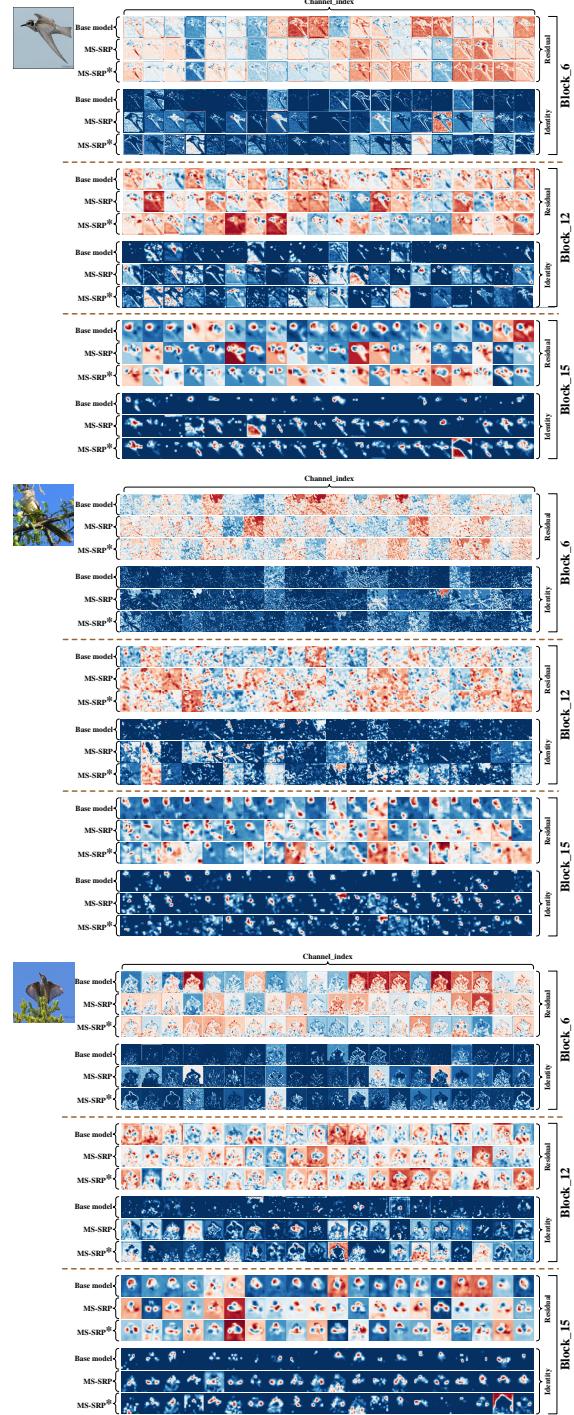


Figure 8: The visualization of feature maps from network block of different depths in **Fine-Grained** dataset. The Base model means D-ResNet-50, the MS-SRP means D-ResNet-50 trained with SRP($M = 5, \lambda = 0.6$) and the MS-SRP* means D-ResNet-50 trained with SRP($M = 5, \lambda = 0.2$). We only display the first 20 feature maps. Best viewed in color.